

# Demystifying Power and Performance Variations in GPU Systems through Microarchitectural Analysis <sup>\*</sup>

Burak Topcu<sup>\*\*1</sup> Deniz Karabacak<sup>2</sup> and Işıl Öz<sup>3</sup>

<sup>1</sup> The Pennsylvania State University

Department of Computer Science and Engineering, State College, PA, USA

topcuuburak@gmail.com

<sup>2</sup> Izmir Institute of Technology

Electrical and Electronics Engineering Department, Izmir, Turkey

denizkarabacak@std.iyte.edu.tr

<sup>3</sup> Izmir Institute of Technology

Computer Engineering Department, Izmir, Turkey

isiloz@iyte.edu.tr

**Abstract.** Graphics Processing Units (GPUs) serve efficient parallel execution for general-purpose computations at high-performance computing and embedded systems. While performance concerns guide the main optimization efforts, power issues become significant for energy-efficient and sustainable GPU executions. Profilers and simulators report statistics about the target execution; however, they either present only performance metrics in a coarse kernel function level or lack visualization support that can enable microarchitectural performance analysis or performance-power consumption comparison. Evaluating runtime performance and power consumption dynamically across GPU components enables a comprehensive tradeoff analysis for GPU architects and software developers. In this work, we present a novel memory performance and power monitoring tool for GPU programs, GPPRMon, which performs a systematic metric collection and provides useful visualization views to guide power and performance analysis for target executions. Our simulation-based framework dynamically gathers SM and memory-related microarchitectural metrics by monitoring individual instructions, and reports achieved performance and power consumption information at runtime. Our visualization interface presents spatial and temporal views of the execution. While the first demonstrates the performance and power metrics across GPU memory components, including global memory, cache, and SMs, the latter shows the corresponding information at the instruction granularity in a timeline. Based on our framework, we demonstrate performance and power analysis for memory-bound graph applications and resource-critical embedded programs from GPU benchmark suites. Our case studies reveal potential usages of our tool in memory-bound kernel identification, performance bottleneck analysis of a memory-intensive workload, performance-power evaluation of an embedded application, and the impact of input size on the memory structures of an embedded system.

**Keywords:** GPU Computing, Performance monitoring, Power consumption

<sup>\*</sup> This article is an extended version of our previously published workshop paper, "GPPRMon: GPU Runtime Memory Performance and Power Monitoring Tool, Burak Topçu, Işıl Öz, Workshop on Resource Awareness of Systems and Society (RAW), co-located with International European Conference on Parallel and Distributed Computing (Euro-Par), 2023."

<sup>\*\*</sup> Work was done when the author was at Izmir Institute of Technology.

## 1. Introduction

As high-performance and energy-efficient computation requirements increase in data processing tasks, parallel GPU architectures, with heterogeneous components, play an essential role in accelerating parallel workloads such as streaming and graph applications [40, 47]. Since GPU devices have evolved into more complex systems with the recent technological developments, efficiency and throughput utilization require more detailed research effort. As a result of high computational capacity, energy consumption and power issues have become crucial in GPU-based systems [22].

While GPU devices have large computational power with multiple processing units, the performance and energy efficiency may decline for memory-intensive workloads due to the high pressure on memory units by concurrently executing multiple threads. While several works discuss the impacts of the memory wall problem for the GPUs [10, 12], there are also studies explaining the root causes of memory bottleneck points for GPU applications [20, 32] and proposing various improvements for the problem [13, 16, 46, 50]. Additionally, the researchers propose energy-efficient methods by considering power consumption of the memory operations on GPU programs [6, 36]. While both performance and energy improvements contribute to the efficient execution of GPU programs, they may compete with each other, and the design decisions become critical and get complicated, requiring the evaluation of the tradeoffs between performance and energy efficiency [4, 9, 17, 42].

Performance and energy efficiency analysis for GPU execution requires low-level measurements at runtime and detailed evaluations of the performance bottlenecks and power consumption. Besides data-parallel programs utilizing SIMD parallelism in GPU resources, general-purpose task-parallel applications introduce more sophisticated algorithm implementations. Therefore, evaluating the execution for performance and power consumption at the kernel function level hides most of the clear evidence for conducting a baseline analysis from many perspectives. More fine-grained level analysis, where the individual instructions belonging to a warp, the smallest execution element in the GPU execution context, should be performed throughout the execution on SM resources. However, profiling and simulation tools collect and report GPU performance results at the kernel level. For instance, the NVIDIA Nsight Compute Tool [27], which presents occupancy, IPC, and memory utilization metrics, operates on a kernel basis. Similarly, the state-of-the-art GPU simulation tools [15, 45] report the performance and hardware metrics at the kernel level. None of the tools directly reports GPU programs' dynamic performance, memory access behavior, and power consumption at runtime. While profiling tools and simulation frameworks report runtime statistics, the software developers and researchers spend additional effort to collect related metrics from the experiments. In other words, several in-house target-specific works still exist to monitor the runtime performance and power consumption for a GPU execution, and repetitive studies cause redundant effort.

While a set of profiling tools helps software developers to understand performance and power consumption information about GPU execution, microarchitecture-level simulators are quite significant in modeling hardware and monitoring the runtime application behaviors. The simulators target micro and macro scales by collecting performance and energy metrics. In GPU-related research, GPGPU-Sim [15] and Multi2Sim [45] simulators have been prominent in offering accurate hardware models for NVIDIA and AMD GPUs, re-

1 spectively, among the other simulators [3, 7, 34]. Moreover, the developer communities  
2 of these simulators have provided continuity by incorporating new GPU architectures and  
3 optimizations introduced in GPU hardware and software. Among the top five hundred  
4 computer systems [41], 179 use NVIDIA-based GPUs as the accelerator/co-processor  
5 technology; specifically, 84 include NVIDIA Volta, and 64 use NVIDIA Ampere devices.

6 In this work, we build and implement, **GPPRMon**, a performance and power moni-  
7 toring tool, for GPU programs executing on top of a simulation environment. [We aim to](#)  
8 [close the gap between the profilers' high-level static results and the cycle-accurate sim-](#)  
9 [ulators' large-volume raw data about the execution. Not only do we set up GPPRMon](#)  
10 [built on configurable simulation execution, but we also generate abstractions and provide](#)  
11 [visualization views that are easy to capture and comprehend large amounts of data. Our](#)  
12 [tool presents both dynamic and configurable simulation execution, and rich architectural](#)  
13 [profiling views by combining the best of both worlds.](#)

14 As potential users of the GPPRMon, we target program developers, system architects,  
15 and researchers, who are working to optimize GPU software or hardware, considering  
16 both performance improvement and energy efficiency. Our simulation-based framework  
17 dynamically collects microarchitectural metrics by monitoring individual instructions'  
18 issues/completions and reports achieved performance and power consumption informa-  
19 tion at runtime. Hence, it enables the users to analyze the dynamic behavior of memory  
20 accesses and thread blocks during the target program execution. Based on the detailed  
21 characteristics collected at runtime, our visualization interface presents both spatial and  
22 temporal views of the execution, where the first demonstrates the performance and power  
23 metrics for each hardware component, including memory units and SMs; and the latter  
24 shows the corresponding information at the instruction granularity in a cycle-based time-  
25 line. Our tool enables the users to perform a fine-granularity evaluation of the target execu-  
26 tion by observing instruction-level microarchitectural features related to performance and  
27 power consumption at runtime. In this article, we extend our previously published work-  
28 shop paper [44] by including additional case studies that demonstrate the usage scenarios  
29 of our tool. To the best of our knowledge, this is the first work monitoring a GPU kernel's  
30 performance by visualizing the execution of instructions for multiple user-configurable  
31 scenarios, relating memory hierarchy utilization with performance, and tracking power  
32 dissipation at runtime. Our main contributions are as follows:

- 33 – We design a systematic microarchitectural metric collection methodology that keeps  
34 track of instruction-per-cycle (IPC) per streaming multiprocessor (SM) as a perfor-  
35 mance metric, instruction execution records for each warp to observe issue and com-  
36 pletion cycles, memory statistics per each component in the memory hierarchy to  
37 understand the possible impacts on performance and power-related statistics per each  
38 GPU component at runtime. We build our configurable collection framework on top  
39 of the GPGPU-Sim simulation environment [15], which provides cycle-accurate in-  
40 formation about the execution.
- 41 – Based on the information gathered by our metric collection module, we design and  
42 build a visualization framework that executes in parallel to our collection module  
43 and generates displays and charts for each kernel execution with the following three  
44 perspectives:

- 1 1. *General View* displays the average IPC among SMs, access statistics of L1 and  
2 L2 caches, row buffer utilization of DRAM partitions, and dissipated power by  
3 the main components for any execution cycle interval.
- 4 2. *Temporal View* shows the details of the instructions with issue and completion cy-  
5 cles for each thread block at warp level. In addition to power consumption statis-  
6 tics for the sub-components in an SM, we include L1 Data (L1D) cache access  
7 statistics, which are local for each SM, to relate the thread block’s performance  
8 in the same execution interval.
- 9 3. *Spatial View* demonstrates the access information for each on-chip L1D cache,  
10 L2 cache in each sub-partition, and row buffers of DRAM banks in each memory  
11 partition. Additionally, it shows the power consumption distribution among the  
12 memory components in the execution interval.
- 13 – We present case studies to demonstrate the potential usages of our framework and its  
14 visualizations by performing experiments for memory-bound graph workloads and  
15 resource-critical embedded applications. Our tool enables us to perform detailed per-  
16 formance and power analysis for the target GPU executions.

17 The remainder of this paper is organized as follows: Section 2 presents some back-  
18 ground on GPU architecture, CUDA programming model, and the simulator. We explain  
19 our design and implementation details for tool construction in Section 3. Then, we present  
20 case studies in Section 4, demonstrating usage scenarios of the GPPRMon. Section 5  
21 presents the existent performance and power evaluation studies for GPGPU applications.  
22 Finally, Section 6 summarizes the work with some conclusive remarks.

## 23 **2. Background**

### 24 **2.1. GPU Hardware**

25 Modern GPU architectures employ a single instruction multiple thread (SIMT) execution  
26 in the Streaming Multiprocessor (SM) units. Each SM includes multiple warp schedulers,  
27 instruction dispatchers, a register file, multiple single and double precision ALUs, tensor  
28 cores, special function units (SFU), and load/store units with on-chip (on-SM) local mem-  
29 ory. An interconnection network connects SMs to off-chip memory partitions on which  
30 DRAM and Last-Level caches (LLC) are placed. While the cores inside the same SM can  
31 access the private L1 cache, all the cores can communicate via the L2 cache structure.  
32 Load/store instructions may require off-chip accesses whenever requested data cannot be  
33 found in the L1D cache. Furthermore, data access gets slower for memory instructions  
34 as moving down the hierarchy. GPU architectures have been evolving, with each new  
35 generation introducing enhancements in performance, power efficiency, and specialized  
36 features. This overview provides a general understanding of GPU architecture, but spe-  
37 cific details and terminology may vary based on the GPU device model and manufacturer.  
38 We specify the architectural and resource specifications for the GPUs in our experimental  
39 setup parts.

### 40 **2.2. CUDA Programming Model**

41 Compute Unified Device Architecture (CUDA) [26] is an API to execute a function,  
42 namely kernel function, on GPUs. A GPU kernel consists of a 3D grid space, where each

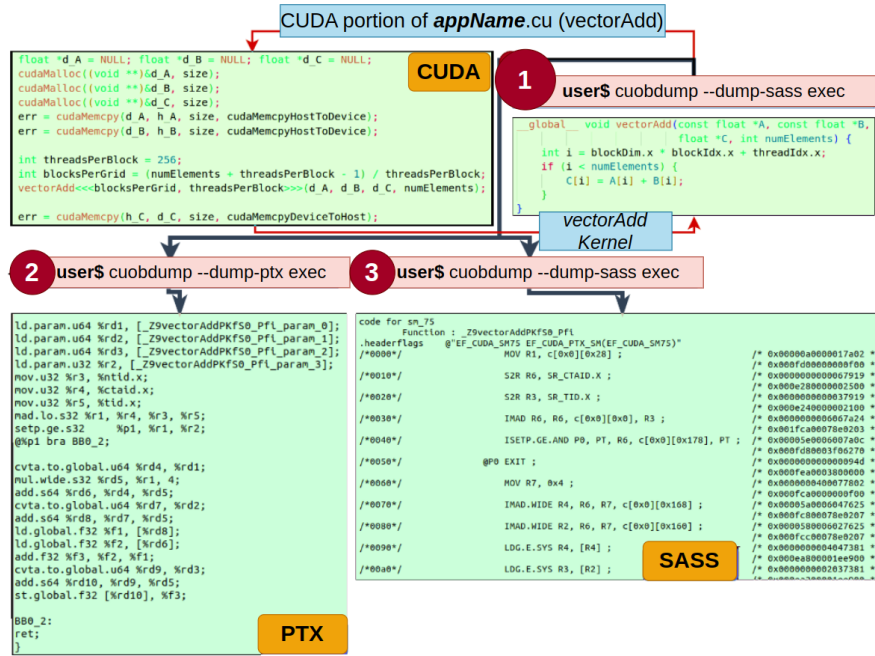


Fig. 1: CUDA, PTX, and SASS code snippets for `vectorAdd` code.

1 grid has a 3D thread block with multiple threads. Each sequential 32-thread set forms  
 2 a *warp* within a thread block in CUDA. When a kernel function launches, the Gigath-  
 3 read engine (thread block scheduler) schedules thread blocks to SMs in the Round-Robin  
 4 fashion. Register resources on SMs determine the number of active thread blocks, some  
 5 of which may be issued to the same SMs. Figure 1 demonstrates a GPU kernel code,  
 6 `vectorAdd`, for the vector addition operation. Part 1 presents a compilation command by  
 7 `nvcc` [29], which is the CUDA compiler to generate the executable. Part 2 demonstrates  
 8 the Parallel Thread Execution (PTX) [30] instructions, which represent a virtual machine  
 9 ISA generated by `nvcc`, and Part 3 includes the SASS machine instructions, which repre-  
 10 sent low-level machine assembly that compiles to binary code executing on NVIDIA  
 11 GPU hardware.

### 12 2.3. GPGPU-Sim Simulation Framework

13 A cycle-level microarchitectural simulator GPGPU-Sim [15] (hereafter referred to as the  
 14 simulator) has been heavily utilized by researchers working on GPU software and hard-  
 15 ware optimizations. Figure 2 displays the workflow of the simulator, which configures  
 16 a traditional NVIDIA GPU and simulates CUDA-based applications. The simulator pro-  
 17 vides functional and performance-driven modes such that functional mode enables devel-  
 18 opers to check the kernel’s functional correctness, while the performance mode simu-  
 19 lates the kernel for the configured GPU in a cycle-accurate manner. The simulator offi-  
 20 cially supports Volta-based Titan V, V100, RTX2060 GPU series, Pascal-based Titan X,  
 21 Kepler-based TITAN, and Fermi-based GTX480. Additionally, AcceleWattch developers  
 22 introduced Volta-based RTX2060 S and Ampere-based RTX3070 GPUs to the official

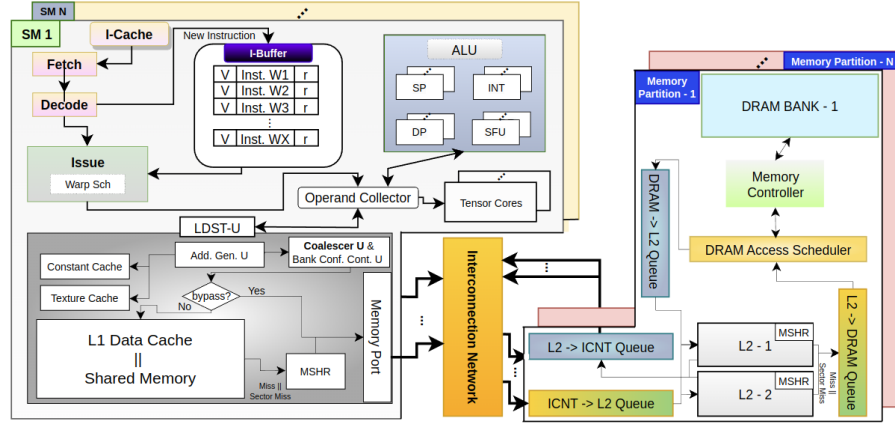


Fig. 2: A general workflow simulation model for a modern GPU.

1 simulator configurations. Beyond these, one can reconfigure any GPU with different re-  
 2 sources on top of these architecture models. Additionally, the simulator reports achieving  
 3 15% performance accuracy error rate with its virtual ISA implementation [15].

4 The AccelWatch [14], a power model extension of the simulator, is an analytical  
 5 model formulating power dissipation and utilizing validated power coefficients of mi-  
 6 croarchitectural components such as functional units, CUDA core lanes, and memory  
 7 components, which are gathered through a comprehensive set of experiments. Accel-  
 8 Watch, which supports Dynamic Voltage-Frequency Scaling (DVFS), estimates the en-  
 9 ergy consumption for V100 with 90% accuracy. In addition to V100, AccelWatch is val-  
 10 idated for TITAN X and Turing RTX GPU series through a large set of applications from  
 11 Rodinia, Parboil, CUTLASS, and DeepBench benchmark suites and NVIDIA CUDA  
 12 Samples, enables tracking detailed power dissipation of any GPU kernel execution.

13 Since we build our GPPRMon tool completely on top of GPGPU-Sim and its Accel-  
 14 Watch power model, its accuracy, architectural support and scalability limitations can be  
 15 considered parallel to the support of the simulation framework.

### 16 3. Methodology

17 GPPRMon tool enables monitoring and visualizing runtime GPU execution performance  
 18 and power consumption. Figure 3 displays GPPRMon workflow consisting of two main  
 19 parts: [1] Metric Collection, [2] Visualization. For any execution interval, GPPRMon sys-  
 20 tematically calculates IPC rates of SMs, records warp instruction’s issues and comple-  
 21 tions, collects memory access statistics across the memory hierarchy, and tracks dissipa-  
 22 ted power on sub-hardware units. Parts [1-a] and [1-b] demonstrate examples of the  
 23 power and performance metrics, respectively, and Section 3.1 details what these metrics  
 24 are and configuration options for users. Furthermore, Part [2] reveals GPPRMon’s visu-  
 25 alizer that contains three views to show general performance, memory access statistics,  
 26 instruction monitoring, and their power dissipation by processing the collected metrics.  
 27 We build our framework on top of the simulator, which is compatible with many GPU

- 1 models mentioned in Section 2.3. The GPPRMon framework is available as open-source  
 2 software <sup>4</sup>.

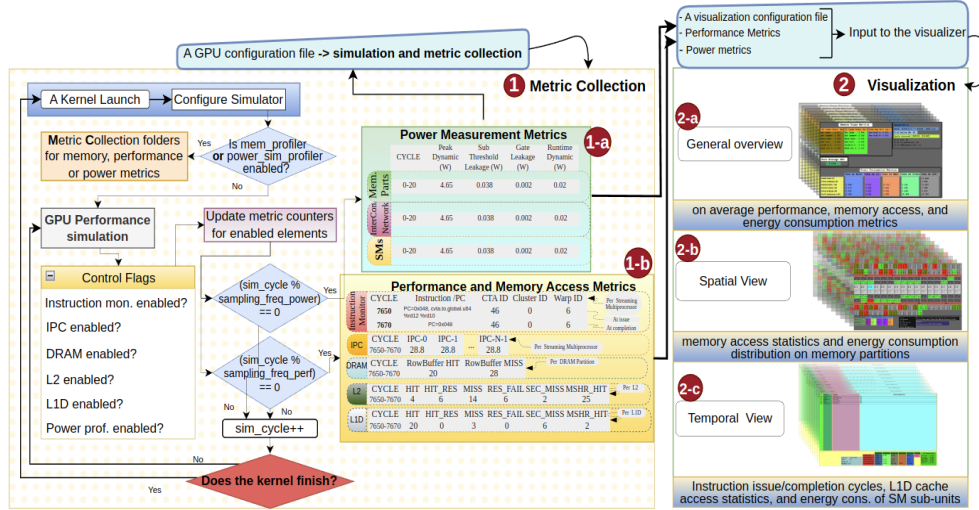


Fig. 3: A general workflow overview of GPPRMon framework.

### 3.1. Metric Collection

The *Metric Collection* phase of GPPRMon, shown in Part 1 in Figure 3, conducts the systematical recording of performance and power metrics during the execution. Performance metrics mainly consist of hardware utilization and execution statistics of the memory hierarchy and SMs, such as cache usage efficiency and SMs IPC values. The metric collector extension to the simulator cumulatively tracks the hardware performance, execution statistics, and power counters during a sampling interval, which is determined by *sampling\_freq\_perf* in Figure 3. At the end of each observation interval, the metric collector exports the tracked metrics to their respective files and clears counters. To enable the performance metric collection feature, the user needs to specify the *mem\_profiler* flag in the simulation configuration file. Similarly, power metrics, such as dynamic runtime and peak power, can be collected by setting the *power\_sim\_profiler* flag.

GPPRMon's metric collection is multi-functional, such that users can separately track warp instruction issue/completion for each thread block, runtime IPC values of each SM, and runtime memory hierarchy utilization statistics (i.e., L1, L2, and DRAM row buffers). GPPRMon allows users to either accumulate or independently collect metrics for each observation interval. Furthermore, users can discard store operations from the memory hierarchy utilization metrics for their runtime target observation. Users can configure these features through metric collection specifications, as depicted *Control Flags* in Figure 3. GPPRMon creates separate folders for each component's statistics and distinct files for

<sup>4</sup> <https://github.com/parsiye/GPPRMon>

1 each sub-component with IDs. To reduce storage and access overheads during the ker-  
 2 nel’s execution, we utilize CSV file format for recording metrics. Our visualizer, shown  
 3 in Part [2], can execute parallel to Part [1] and processes the collected metrics to gener-  
 4 ate runtime visualizations. The following subsections will briefly describe the collected  
 5 microarchitectural performance and power metrics and how to configure each metric col-  
 6 lection separately. More detailed information about how to build GPPRMon and deploy  
 7 multi-functional metric collection is available on the tool’s GitHub page.

## 8 Performance Metrics

9 ***L1 Data and L2 Caches.*** A memory request’s access status on caches can be one of  
 10 the possible states: i) Hit: Data resides on the cache line; ii) Hit Reserved: The data is re-  
 11 quested, and the corresponding cache line is allocated, but the data is still invalid; iii) Miss:  
 12 The corresponding cache line causes several sector misses, resulting in a line eviction; iv)  
 13 Reservation Failure: The situations in which a line allocation, MSHR entry allocation, or  
 14 merging an entry with an existing in MSHR fail, or the miss queue does not have any  
 15 slot to hold new requests result in reservation failure; v) Sector Miss: A memory request  
 16 cannot find the data in the sector of a cache line (i.e., a sector is 32B, whereas a cache  
 17 line is 128B); vi) MSHR hit: When the upcoming request’s sector miss has already been  
 18 recorded, and request can merge with the existing entry, MSHR hit occurs.

19 [GPPRMon can observe runtime access statistics of each L1D and L2 cache sepa-](#)  
 20 [rately. Users can activate the metric collection feature of GPPRMon for L1D and L2](#)  
 21 [caches with regarding control flags as depicted in Figure 3.](#) Tracking runtime cache uti-  
 22 lization helps researchers evaluate the application’s memory access behavior and relate  
 23 it to the overall application performance. While cache hits indicate small access laten-  
 24 cies, misses generally refer to longer latencies and increasing active memory traffic in the  
 25 hierarchy. Furthermore, reservation failures result in a memory pipeline stall, directly de-  
 26 laying the subsequent load/store operations. Handling intensive misses requires detailed  
 27 analytic observations to deduct behavioral interpretations across applications and cache  
 28 utilization. [In this manner, GPPRMon can meet the runtime analytic observation neces-](#)  
 29 [sity for micro-architectures on GPU to identify the memory performance and power issues](#)  
 30 [understandably.](#)

31 ***Row Buffers of DRAM Banks.*** A row buffer hit occurs when an L2 miss request finds  
 32 the requested data in the row buffer of the target DRAM bank. A row buffer miss results  
 33 in a longer access service time than those hits because handling a row buffer miss requires  
 34 scheduling a memory request to the correct address on DRAM and activation latencies.  
 35 Row buffer utilization with cache access statistics completes the runtime memory hierar-  
 36 chy behavior exploitation, crucial for describing overall memory performance, especially  
 37 for sparse data applications. [GPPRMon provides separate metric collections for the row](#)  
 38 [buffers in each DRAM partition, and users can activate this feature by enabling DRAM](#)  
 39 [control flag as depicted in Figure 3.](#)

40 ***SM IPC.*** An IPC rate mainly describes an SM’s performance, which consists of various  
 41 functional units with varying lane depths. For example, V100 SMs include four single-  
 42 precision (SP) ALUs with four pipe depths, as presented in Figure 2. When a thread



1 block occupies only SP-ALU lanes, assuming an operation takes one cycle, the ideal IPC  
2 for each SM should be sixteen without other functional units' contribution. However,  
3 IPC oscillates during the execution depending on SM and memory hierarchy utilization.  
4 GPPRMon tracks the runtime IPC rate for each SM separately for each sampling interval,  
5 and one can active per SM IPC tracking by enabling *IPC control flag* as in Figure 3. With  
6 GPPRMon, developers can analyze runtime IPC variations and investigate the root causes  
7 of IPC suffering.

8 **Instruction Monitor.** The instruction monitor feature of GPPRMon records the issue/completion  
9 cycles of warp instructions within each thread block together with their opcode, operand,  
10 and PC separately. Since GPUs execute instructions for multiple threads concurrently by  
11 a common PC within a warp, we design tracking instruction issues/completions at the  
12 warp level. While the first row of the *Instruction Monitor* in Part [1-b] shows the issue  
13 statistics, the second displays the completion. Even if the dispatcher unit may issue the  
14 same instruction multiple times for a warp, it is guaranteed that any two instructions,  
15 of which *CTA\_ID*, *SM\_ID*, *local Warp\_ID*, and *PC* are the same, cannot change the is-  
16 sue/completion sequence. Hence, we can obtain the correct issue/completion matching  
17 for each instruction. Users can activate the instruction monitoring utility for each thread  
18 block with *instruction mon. control flag* as in Figure 3.

## 19 Power Metrics

20 The comprehensive results of the dissipated power on GPU yield the analytical observa-  
21 tion that gains significance, especially on embedded systems. Therefore, we develop GP-  
22 PPRMon to systematically collect the power distribution on the sub-units of SMs, memory  
23 partitions, and the interconnection network, in addition to performance metrics. More-  
24 over, SMs are GPUs' most impactful hardware units in terms of runtime power dissipa-  
25 tion with their dense compute units. Thus, GPPRMon further classifies SM's power  
26 distribution through execution units, including the register file and the beginning of the  
27 instruction pipeline, functional units, and load/store units involving the L1D cache.

28 We implement the collection of power metrics utility on top of the AccelWattch [14],  
29 built upon McPAT [21], and the accuracy and theoretical limitations of AccelWattch yet  
30 reside. However, GPPRMon's power metric collection feature is still a reliable and prac-  
31 tical tool since we extend various configurability options of AccelWattch to GPPRMon,  
32 such as DVFS. GPPRMon assures the following runtime power measurements for each  
33 component apart from idle SM: *Peak Dynamic(W)*, the maximum momentary power  
34 within the interval, *Sub-threshold Leakage (W)* and *Gate Leakage(W)*, the leaked power  
35 (due to current leakage) from the junctions of MOSFETs, and *Runtime Dynamic (W)*,  
36 the total consumed power. Moreover, GPPRMon supports collecting power metrics either  
37 cumulatively or distinctly for each sampling interval, starting from a kernel's execution.  
38 For instance, power dissipation results in Figure 4 show the *cumulatively* collected results  
39 for the interval between [55000, 56000], which means aggregating power results for each  
40 sampling interval. The cumulative power metric collection option eases determining aver-  
41 age power dissipation for different execution intervals. While V100's TDP is 300W [35],  
42 aggregated *Runtime Dynamic* power for 1000 cycles is 1095.5W, which corresponds to an  
43 average of 109.5W per sampling interval with the cycling frequency of 100 on the sim-

1 ulator. Users can configure the power metric collection by first enabling power profiling  
 2 and determining other configurations as detailed in the tool’s source code.

### 3 3.2. Visualization

4 By processing the collected metrics (Part 1 in Figure 3), GPPRMon depicts performance  
 5 and power dissipation with three perspectives at runtime, as represented in Part 2 in  
 6 Figure 3, and enables pointing out detailed interaction of the program with the underlying  
 7 hardware.

- 8 i **General View**, Part 2-a, presents the overall IPC of GPU, the average memory ac-  
 9 cess statistics, and dissipated power among the major components with application-  
 10 and architecture-specific information;
- 11 ii **Spatial View**, Part 2-b, presents the detailed access statistics of the GPU memory  
 12 hierarchy and dissipated overall power among the memory partitions by enabling the  
 13 monitoring of the entire GPU memory space;
- 14 iii **Temporal View**, Part 2-c, demonstrates instruction execution statistics with activa-  
 15 tion intervals at warp-level for user-specified thread blocks, L1D cache access char-  
 16 acteristics, and power distribution among the sub-components of SMs by activating  
 17 the execution monitoring feature.

On Average Memory Access Statistics						
<b>L1D Cache Stats (Av)</b>		<b>L2 Cache Stats (Av)</b>		<b>DRAM Row Util. (Av)</b>		
Hit Rate	0.003	Hit Rate	0.312	Row Buffer H	0.383	
Hit Reserved R	0.001	Hit Reserved R	0.000	Row Buffer M	0.617	
Miss Rate	0.038	Miss Rate	0.464	Kernel ID: 0		
Reserv. Failure R	0.944	Reserv. Failure R	0.000	Cycle Interval: [55000, 56000]		
Sector Miss R	0.013	Sector Miss R	0.223	Grid:(1784,1,1) Block:(256,1,1)		
MSHR Hit R	0.008	MSHR Hit R	0.005	# of active SMs: 80		
Average IPC on SMs : 1.08						
<b>Dissipated Power</b>		InterCon. Net	L2	Mem Part.	SMs	GPU
Peak Power (W)						185.63
Total Leakage (W)						17.346
Peak Dynamic (W)	0.338	4.687	137.55	25.704	168.264	
Sub-Threshold Leak (W)	0.067	0.138	1.316	13.474	14.995	
Gate Leakage (W)	0.011	0.013	0.016	2.168	2.352	
Runtime Dynamic (W)	64.618	2.537	823.184	205.174	1095.513	

Fig. 4: *General View* with average performance and power consumption metrics collected at runtime.

18 GPPRMon includes three configuration options for different visualization perspec-  
 19 tives and an interval sampling cycle to divide runtime execution into portions. Since GP-  
 20 PPRMon’s *Temporal View* may require scanning of many statistics for all thread blocks,  
 21 especially in large architectures, GPPRMon provides a *Temporal View* option among the  
 22 thread blocks determined with IDs. Depending on the configuration, GPPRMon starts  
 23 tracking collected metrics for each kernel and systematically saves images in PNG format  
 24 per execution interval.

Figure 4, an example of our *General View*, presents the overall measurement of the first kernel for the *SPMV* from *Gardenia* benchmark [48] executed on V100 GPU [23]. It displays average memory access statistics among the active L1D caches, L2 caches, and DRAM banks; average IPC value among the active SMs; dissipated power on major sub-GPU components within the [55000, 60000] cycle interval. The view includes grid (i.e., 1764 thread blocks) and block dimensions (i.e., 256 threads per block) with the number of actively used SMs so that the users can relate active SMs and workloads at runtime. To illustrate, Figure 4 shows that the kernel executes with the IPC rate of 1.08 and utilizes the memory hierarchy inefficiently due to high miss and reservation failure rates on V100 in this interval, where the Volta SM architecture supports concurrent execution of 2048 threads per SM. Considering that V100 SMs contain 16 SP/INT/DP ALUs, SFUs, and Tensor Cores, we can notice the low performance since the ideal IPC should be much more than 1.08 with 640 thread blocks (8 thread blocks per SM) in the given interval. Long-latency memory operations may slow instruction completion and result in low IPC. Moreover, memory partitions consume 75% of total power dissipation, which validates that SMs mostly stay idle for the execution interval given in Figure 4.

Our *General View* supports two additional visuals that show the time spanning of memory access statistics and the relationship between IPC and power metrics for longer runtime intervals as in Figure 11 and Figure 14 (the examples given as part of our case studies), respectively.

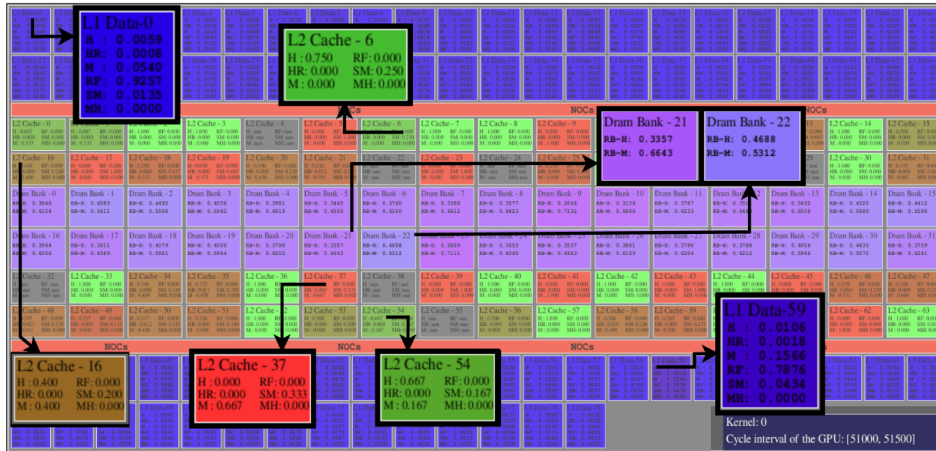


Fig. 5: *Spatial View* displaying memory access statistic at runtime.

Figure 5, an example of our *Spatial View*, shows the memory access statistics across the GPU memory hierarchy on L1D caches, L2 caches, and DRAM. On caches, the green emphasizes hit and hit reserved accesses, the red indicates miss and sector miss, and the blue states the reservation failures through miss queues or MSHR. Similarly, DRAM bank pixels are colored with a mixture of red and blue to specify the row buffer misses and hits, respectively. GPPRMon supports memory hierarchy visualization for all official GPU configurations of the simulator, even if some disable the L1D cache for load/store operations. *Spatial View* includes detailed data about access statistics, resource quanti-

ties, architecture types, and dissipated power on memory partitions. To detail the view, we zoom in on some memory units in Figure 5, which presents statistics in the cycles of [51000, 51500] for the kernel execution. Different L1D caches behave similarly in that interval, such that almost all L1D caches turn blue due to reservation failure concentration. To illustrate, the reservation failure rates are 0.92 and 0.78 on *L1D-0* and *L1D-1*, respectively. On the contrary, statistics on L2 caches imply heterogeneous data accesses. While *L2 Cache-6* and *L2 Cache-54* bring hit rates of 0.75 and 0.67, respectively, *L2 Cache-37* causes 0.67 miss and 0.33 sector miss rates, and L2 cache is intermediate in utilizing data locality among the others. Gray color among the units regards no access occurring on the *Spatial View*.

PC	OPCODE	OPERAND	ISSUE / COMPLETION	Performance and Power Metrics					
352	fma.rm.f32	%f21 %f20 %f19 %f18	1-8044 1-8053	<b>Dissipated Power</b>	Execution U.	Func. U.	LD/ST U.	IDLE	TOTAL
360	st.global.f32	[%rd1] %f21	1-8053 1-8107	Peak Dynamic (W)	18.158	1.000	6.546		25.704
368	ld.global.f32	%f22 [%rd16 + 24]	1-8071 1-8179	Sub-Threshold Leak (W)	10.808	0.587	1.465		13.474
376	ld.global.f32	%f23 [%rd14 + 24]	1-8072 1-8178	Gate Leakage (W)	0.038	0.074	1.983		2.168
448	add.s32.f32	%r15 %r15 8	2-8072 1-8326 2-8082 1-8337	Runtime Dynamic (W)	75.928	1.645	437.529	0.000	515.102
456	setp.ne.s32%p2	%r15 0	2-8082 1-8337 2-8088 1-8343	<b>IPC Rate on SM : 3.776</b>					
464	@ %p2	bra BBO_2	2-8088 1-8343 2-8093 1-8348	<b>SM ID : 2</b>					
168	add.s64	%rd14 %rd15 %rd2	2-8090 1-8345 2-8099 1-8354	<b>Thread Block ID : 2</b>					
				<b>Kernel ID : 0</b>					
				<b>Cycle Interval : [8000, 8500]</b>					
				<b>L1D Cache Stats (Av)</b>					
				Hit Rate 0.429					
				Hit Reserved R 0.000					
				Miss Rate 0.437					
				Reserv. Failure R 0.000					
				Sector Miss R 0.134					
				MSHR Hit R 0.000					

Fig. 6: *Temporal View* monitoring instruction executions together with the performance and power consumption of the corresponding SM.

Figure 6, an example of our *Temporal View*, displays a thread block’s execution statistics at warp-level, L1D cache statistics, and dissipated power of core components configurable execution intervals. It presents each warp’s PTX instruction sequence, with opcodes, operands (source/destination registers and immediate values if they exist), and the program counter (PC). The *Issue/Completion* column indicates the execution start and writeback times of warp instruction segments within any thread block. For instance, Figure 6 reveals the execution monitoring of the *Thread Block 2* on *SM 2* in the cycle range of [8000, 8500]. The instruction dispatcher unit issues two SP global loads with PC=368 and PC=376 at *Cycle 8071* and *Cycle 8072*, and they are completed at *Cycle 8179* and *Cycle 8178*, respectively. *Temporal View* allows tracking execution duration per instruction in this manner. Since L1D cache hits result in low latency, these load instructions lasting above 100 cycles are among the misses or sector misses of L1D cache statistics. In addition, the view enables us to relate IPC, instruction statistics, and power metrics. The fact that the rate of memory instructions is 0.37 and inefficient use of the L1D cache within the given interval significantly degrades the IPC on SM2. Furthermore, the load/store unit dissipates nearly 92% of SMs total power in the corresponding interval because of the pressure on the L1D cache. As a result, one can analyze the data locality in a multi-perspective by utilizing access statistics of caches and row buffers on *Spatial View*, tracing the issue/completion times of memory operations on *Temporal View* at runtime.

**GPPRMon Overheads.** GPPRMon execution performance mainly depends on *i. the metric collection sampling frequency*, *ii. visualizer intervals*, and *iii. view types*. Firstly,

1 each sampling operation during simulation conducts multiple I/O operations to the met-  
2 ric collection files, such as performance, memory, and power trace files, which are sig-  
3 nificantly slow compared to the simulation execution. Hence, widening the simulation  
4 runtime sampling interval directly reduces the number of I/O operations for exporting  
5 results to output files, and the simulation duration decreases accordingly. For instance,  
6 simulating the *SPMV* benchmark [48] takes 98 minutes for the *Higgs Twitter Mention*  
7 data by recording both power and performance metrics per 5000 simulation cycles, while  
8 the baseline simulation (i.e., not collecting runtime performance and power metrics) com-  
9 pletes the execution at 88 minutes in our local infrastructure. Furthermore, the visualiza-  
10 tion interval mainly determines the spanning range of the collected metrics. That is, lower  
11 visualization intervals search for fewer sampled metrics, reducing generating views to  
12 demonstrate runtime performance and power observations on the GPU. Since each view  
13 type requires different metrics, composing each figure takes various amounts of time.  
14 For example, generating *General View* in Figure 4 necessitates spanning all statistics on  
15 L1D and L2 caches, row buffers, SMs, and power dissipation for a given interval while  
16 *Spatial View* only displays memory access statistics, and is much easier to compose. Ad-  
17 ditionally, tracking the instruction monitoring of limited thread blocks through *Temporal*  
18 *View* is comparably easy, whereas increasing the number of thread blocks for instruction  
19 monitoring raises the spanning overhead and takes longer to generate those views.

## 20 4. Usage Scenarios

21 We evaluate a set of CUDA programs from two benchmark suites and demonstrate the  
22 case studies that can use our framework. Specifically, we utilize graph workloads from the  
23 Gardenia benchmark suite [48] and embedded applications from the GPU4S embedded  
24 benchmark suite [37]. Since graph workloads target large-scale systems, we configure the  
25 GPPRMon to simulate Volta architecture-based V100, commonly used in HPC systems  
26 and GPU architecture research. Table 1 presents salient characteristics of the V100 device.  
27 On the other hand, for embedded applications, we configure our tool to simulate the GPU  
28 device on Jetson AGX Xavier, which provides a System-on-Module with a Volta-based  
29 GPU and contains an 8GB/16GB unified memory with a high-bandwidth interface to the  
30 GPU, and a 512KB shared L2 cache exists in the memory hierarchy [25]. Its GPU includes  
31 512 cores corresponding to 8 SMs involving a 128KB on-chip cache per SM.

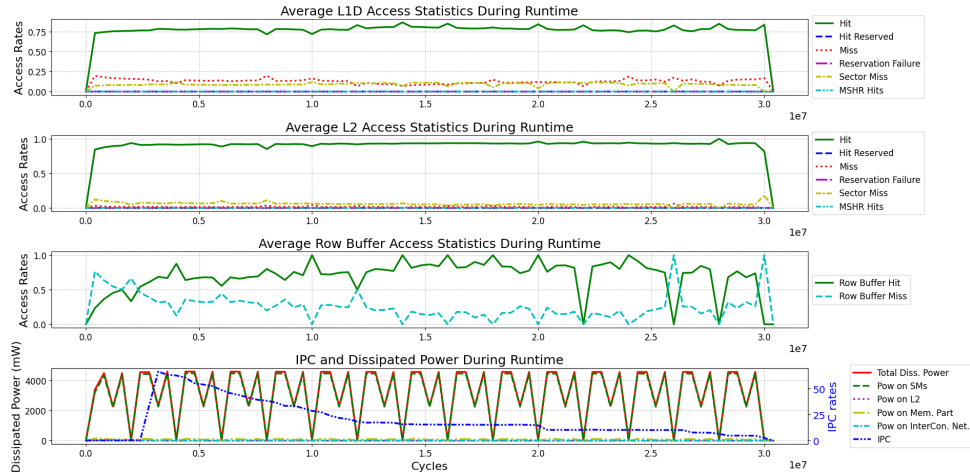
### 32 4.1. Determining Memory-Bound Kernels and Spatial Characteristics

33 We evaluate the *Betweenness Centrality (BC)* program (from Gardenia suite [48]) with  
34 multiple kernel executions. By analyzing the kernels with the largest cycles, we classify  
35 them as kernels with no memory pressure and kernels with memory bottleneck. Addition-  
36 ally, we review SM distribution and utilization of the target executions.

37 Figure 7 presents memory-related data and IPC values during the execution time in-  
38 tervals for the kernel function, *Kernel 1*, and demonstrates L1, L2, and row buffer statis-  
39 tics and the corresponding IPC/power behavior. While our spatial and temporal views  
40 present fine-grained values for each hardware component across the GPU device, those  
41 statistics (Figure 7 and Figure 9 afterward) demonstrate the dynamic average values for  
42 overall GPU SMs. While there are a few oscillations in the rates, the general behavior

Table 1: V100 GPU configuration specifications based on Volta architecture.

SM (80) Specifications	Registers, Register Banks	65536, 16
	SP-U, SF-U, DP-U, INT-U, TC-U, LD/ST-U (WB-PipeDepth)	4,4,4,4,4,1(8)
	Warp Scheduler	4
Memory Partiton (32) Specifications	on-chip L1I Cache, NoF banks, access latency, cache line	128KB (64 sets, 16-way), 1, 20 cycles, 128B
	L2 Cache, NoF banks, access latency, cache line	96KB (32 sets and 24-way), 2, 160 cycles, 128B
	DRAM, NoF banks, access latency (after L2)	1GB, 16 banks, 100 cycles, 128B
	DRAM scheduler	First-ready, first-come first-service

Fig. 7: Average memory and power consumption statistics for *BC-Kernel 1*.

1 demonstrates steadily high hit rates. Moreover, we can see the peak IPC value during the  
 2 intensive computations at the beginning of the kernel execution. Then, IPC values tend to  
 3 decrease as the computations end. We can observe that the kernel has no memory pressure  
 4 during the execution. After warming the caches at the very beginning of the execution, the  
 5 kernel could perform its operations with high hit rates and compatible IPC rates.

6 After getting the memory behavior of our target kernel based on our average tim-  
 7 ing statistics, we can utilize our *Spatial View* to understand its SM utilization. Figure 8  
 8 presents the partial view that includes the L1 cache structures of each SM at time inter-  
 9 vals [5000-10000], [15000-20000], and [245000-250000], respectively. The application  
 10 launches the kernel function with 64 thread blocks and 256 threads per block, such that  
 11 the thread blocks are scheduled to 64 SMs (out of 80 in the V100 device). Therefore, only  
 12 the corresponding L1 caches exhibit non-zero values. Since all thread blocks are active  
 13 and the caches do not hold data at the beginning (5K-10K time interval), all 64 L1 cache  
 14 structures employ low hit rates (red color in our representation, while inactive L1s are  
 15 gray). At the next time interval, the execution starts warming the caches, and hit rates

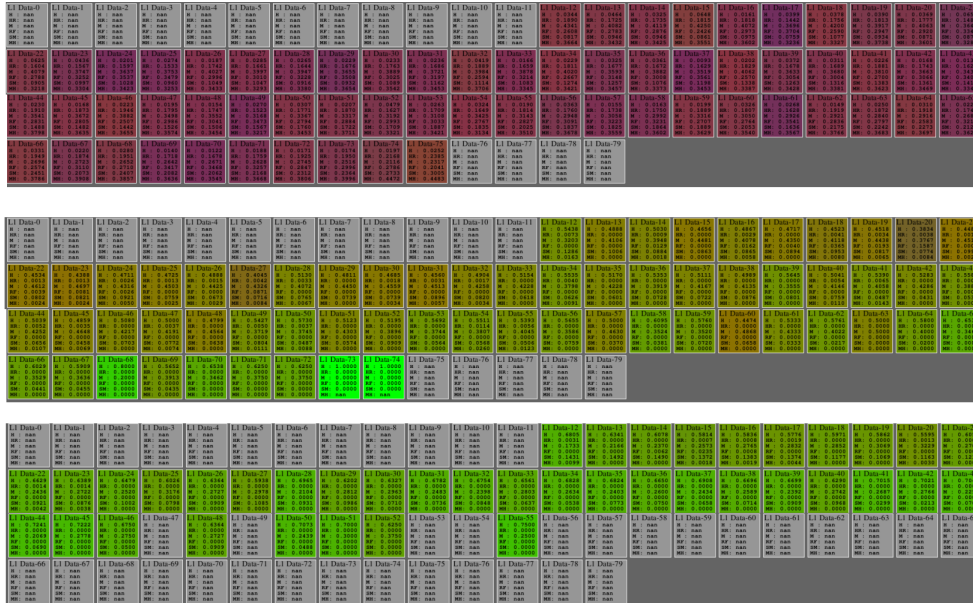


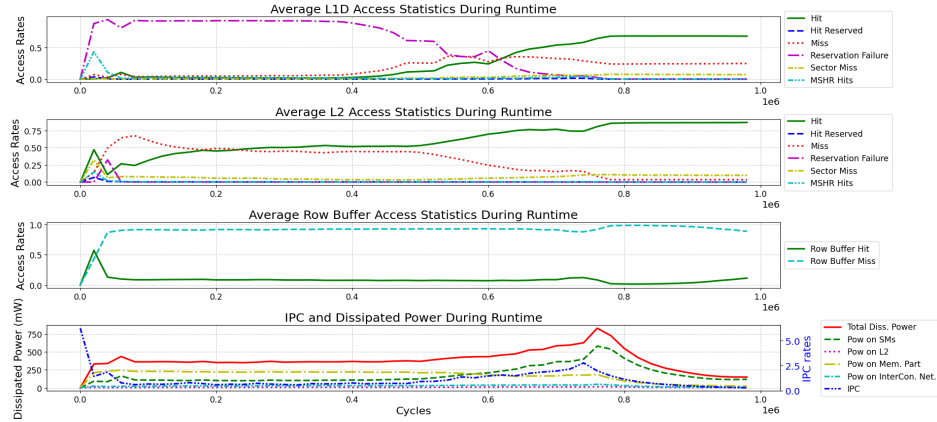
Fig. 8: L1 cache status for *BC-Kernel 1* for successive time intervals.

1 increase (green color in our representation). Eventually, as the thread blocks complete  
 2 execution and they retire, L1 caches on the corresponding SMs do not collect statistics  
 3 (gray as the other L1 caches at the beginning). We can track all thread blocks' scheduling  
 4 and their L1 cache utilization during the execution for the given time intervals. For the  
 5 specific kernel execution, observing the high hit rates (green colors) in all L1 caches of  
 6 the active SMs and corresponding IPC rates demonstrates efficient L1 utilization and low  
 7 memory pressure.

8 Figure 9 demonstrates memory-related data and IPC values for another kernel, *Kernel*  
 9 *2*. Since *Kernel 2* has been launched by much more threads than available SMs, all SMs  
 10 are active during the execution, and the threads compete for both L1 and L2 caches.  
 11 The reservation failure and miss rates at the first part of the execution are significant for  
 12 private L1 and shared L2 cache, in turn, the kernel exhibits low IPC values. Whenever L2  
 13 and especially L1 hit rates become more than 0.5 (around cycle 600K), IPC values also  
 14 increase, and the kernel function performs most of the target computations (until cycle  
 15 800K). We can observe the corresponding cache status in our *Spatial View*. For this case,  
 16 we include both L1 and L2 cache structures in our aggregated visualization in Figure 10.  
 17 Different from *Kernel 1*, we can see that all L1 caches exhibit non-zero values from the  
 18 beginning of the execution since the available threads utilize all SMs in the device. The  
 19 figure demonstrates the warming of the caches through the execution cycles, where the  
 20 kernel function has pressure on both L1 and L2 by highly utilizing those structures.

## 21 4.2. Performance Bottleneck Analysis and its Power Impacts for a 22 Memory-Intensive Workload

23 We evaluate CUDA implementation of the *Page Ranking (PR)* algorithm given in Gar-  
 24 denia suite [48] to analyze a memory-bound GPU program and irregular memory access

Fig. 9: Average memory and power consumption statistics for *BC-Kernel 2*.

1 statistics through GPPRMon. *PR* assigns weights to graph nodes describing relative im-  
 2 portance among nodes.

3 The *PR* execution iterates with the *Contribution Step (K0)*, *Pull Step (K1)*, and *Linear*  
 4 *Normalization (K2)* kernels, and the total number of iterations varies depending on the  
 5 data size. Table 2 presents the performance overview of the kernel statistics. At the begin-  
 6 ning of the execution, *K0* causes a high miss rate on caches since all memory operations  
 7 are completed before warming up. No row buffer locality information is provided across  
 8 DRAM accesses, as the required data mostly fits into the L2 cache. In other words, L2  
 9 misses do not access the same row of DRAM banks within any memory partition. Since  
 10 total elapsed cycles indicate that *K1* dominates the execution at 99.7% and directly affects  
 11 the application performance and overall power consumption, we focus on that kernel ex-  
 12 ecution.

Table 2: *Page Ranking (PR)* kernel performance statistics.

Kernel	GPU IPC	GPU Oc- cupancy	L1D		L2		DRAM		Total Cycle
			Miss Rate	Res. Fail Rate	Miss Rate	Res. Fail Rate	Row B. Loc. (L+S)	Row B. Loc. (Load)	
Page Ranking - Contrib K0	715.59	82.76%	1.000	0.819	0.333	0.0	NA	NA	8670
Page Ranking - PullStep K1	3.007	5.55%	0.584	0.400	0.156	0.011	0.658	0.667	8677889
Page Ranking - Lin-Norm K2	1297.68	77.108%	0.501	0.285	0.457	0.001	0.724	0.732	11718

13 V100 includes 870GB/s bandwidth migrating 217.5 Giga SP float to the SMs and  
 14 14.8-SP/7.4-DP TFLOPS peak computational power [23]. With these specifications, we  
 15 can state that the time consumed for one load complies with the execution of 68 SP float  
 16 operations on SMs, ideally (i.e., without L1D and L2 caches). On the other hand, *K1*  
 17 has a memory instruction intensity of around 0.2 in its PTX code, which validates that



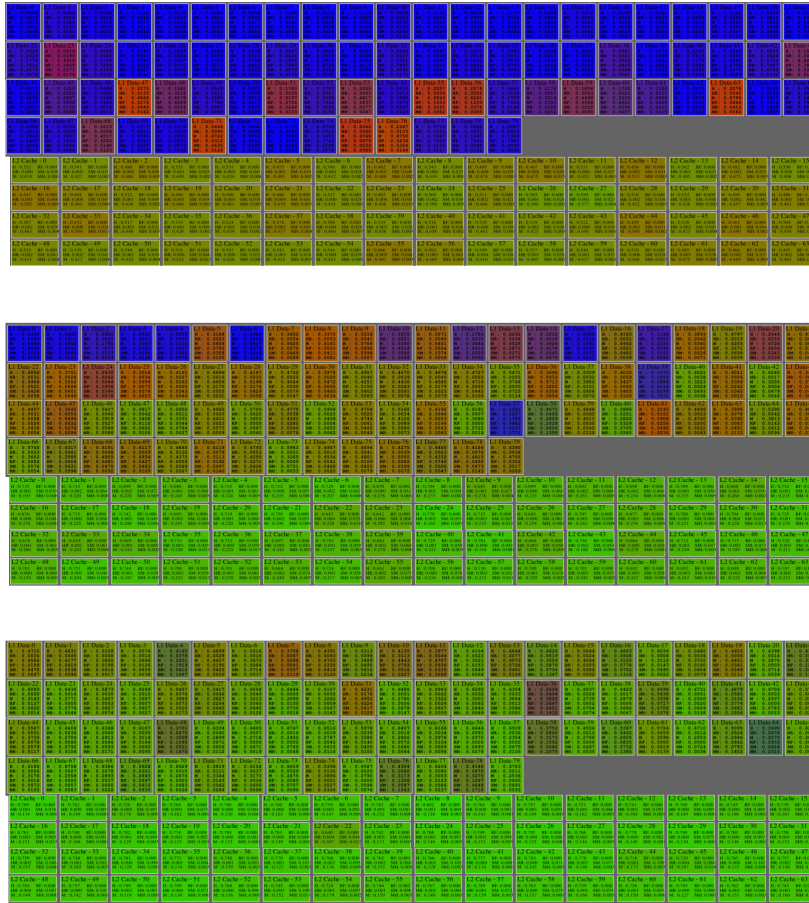


Fig. 10: L1 and L2 cache status for *BC-Kernel 2* from cycles 400K, 600K, 700K.

1 *K1* bounds the performance due to intrinsic memory workload. Not only is *K1* memory-  
 2 bound, but also inefficient memory usage severely impacts performance. To see performance-  
 3 critical points in the execution, we monitor runtime performance-degrading factors with  
 4 low-frequency execution snapshots provided as part of our tool.

5 While the overall IPC is around 0.3 for *K1*, *General View* results obtained for every  
 6 500 cycles during the execution indicate that IPC oscillates in the range of [0.1, 9.54],  
 7 with the highest peak of 53.44. Figure 11, which is part of our *General Overview*, shows  
 8 average access statistics on memory units in [5000, 100000] cycles. After caches warm-  
 9 ing up (around cycles 10000), while the average miss rate on L1D caches oscillates in  
 10 [0.14, 0.51], sector misses, which the simulator does not provide separately, vary in [0.05,  
 11 0.31] with the metrics collected in every 20 cycles. We can understand the data pollution  
 12 on the L1D caches, which prevents the execution from exploiting cache locality. For in-  
 13 stance, *K1* does not utilize spatial locality on the L1D cache since the MSHR hits oscillate

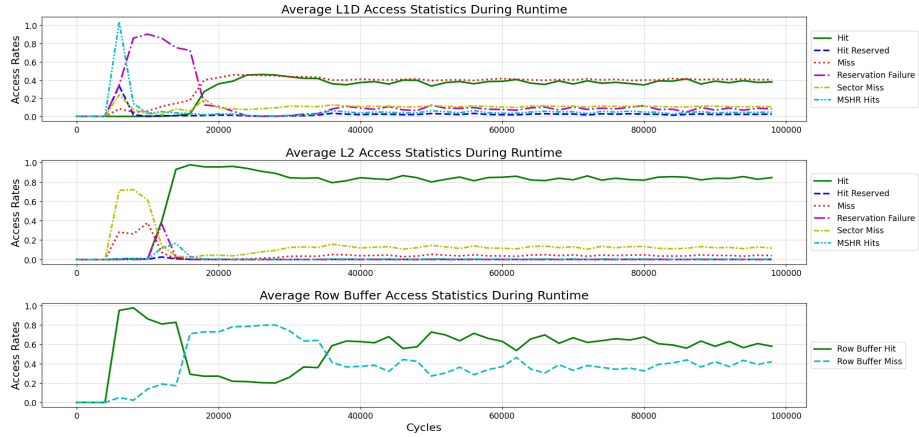


Fig. 11: Average memory access statistics in the cycle range of [5000, 100000] for *PR-K1*.

1 slightly in [0.03, 0.08] during the execution. Hence, the overall hit rate on L2 caches is  
 2 quite high when we use the web-Stanford dataset [19], which occupies memory space  
 3 five times larger than the L2 cache size. While the performance metrics in Table 2 hide  
 4 the L2 statistics as it counts misses before warming up at kernel launch, the actual L2 hit  
 5 rate oscillates in [0.82, 0.95] at runtime with sampling per 500 cycles. Additionally, the  
 6 row buffer hits and misses vary in [0.2, 0.85] in an unstable manner, which verifies data  
 7 sparsity throughout the execution.

PC -> 296	opcode -> ld.global.u64	operand-> %rd1 [%rd11]	Interval : 5000-30000
PC -> 312	opcode -> ld.global.u32	operand-> %rd3 [%rd11+8]	SM ID : 0
			Kernel ID : 1
CTA_ID=24	2-5455 0-5455 1-5455 3-5456 4-5469 5-5469 6-5469 7-5470	32-5565 33-5565 34-5565 35-5565 37-5566 38-5566 39-5641	CTA_ID=344
	2-5826 0-5864 1-6045 3-6075 4-6065 5-6026 6-6087 7-6074	32-6595 33-6719 34-6692 35-6567 37-7711 38-6877 36-7168 39-7214	
	2-5867 0-5868 5-6027 1-6046 4-6066 7-6076 3-6077 6-6088	35-6568 32-6569 34-6695 33-6720 38-6878 36-7169 39-7215 37-7714	
	2-5946 0-5948 5-6053 1-6072 4-6092 7-7696 3-6105 6-6114	35-6897 32-6622 34-6721 33-6746 38-6935 36-7379 39-8263 37-7740	
CTA_ID=104	9-5470 10-5470 8-5470 11-5471 13-5884 14-5891 15-5885 12-5887	41-5570 40-5640 45-5640 44-5642 42-5643 43-5644 47-5704 46-5704	CTA_ID=424
	9-6140 10-6141 8-6143 11-6163 13-8076 14-8063 15-8074 12-8122	41-7081 40-6976 45-7244 44-7272 42-7068 43-7039 47-7827 46-7499	
	9-6141 10-6143 8-6144 11-6164 14-8064 15-8075 13-8077 12-8123	40-6977 43-7040 42-7069 41-7082 45-7245 44-7273 46-7500 47-7830	
	9-6167 10-6170 8-6173 11-8055 14-8090 15-8867 13-8104 12-8149	40-7003 43-7066 42-7096 41-7108 45-7271 44-7299 46-7659 47-8802	
CTA_ID=184	16-5472 17-5472 18-5472 19-5473 20-5474 21-5475 22-5474 23-5475	49-5642 48-5704 53-5704 52-5705 51-5705 50-5708 55-5710 54-5714	CTA_ID=504
	16-6286 17-6251 18-6969 19-6639 20-6551 21-6492 22-6469 23-6445	49-7281 48-7878 53-7510 52-7468 51-7602 50-7848 55-7538 54-7543	
	16-6287 17-6252 23-6446 22-6470 21-6493 20-6552 19-6640 18-6971	49-7282 52-7469 53-7511 55-7539 54-7544 51-7603 50-7849 48-7879	
	16-6313 17-6618 23-7799 22-6496 21-6519 20-6578 19-6666 18-6997	49-7673 52-7495 53-7537 55-8402 54-7570 51-7629 50-7875 48-7905	
CTA_ID=264	24-5476 26-5476 25-5482 27-5482 28-5484 29-5484 30-5484 31-5486	56-5708 57-5708 61-5714 60-5714 59-5718 58-5720 62-5786 63-5786	CTA_ID=584
	24-6619 26-6607 25-6826 27-6579 28-6300 29-6556 30-6635 31-6627	56-7590 57-7630 61-7718 60-7951 59-7779 58-7832 62-7863 63-7919	
	28-6301 29-6558 27-6580 26-6608 24-6621 31-6629 30-6636 25-6827	56-7591 57-7631 61-7719 59-7780 58-7833 62-7864 63-7917 60-7952	
	28-6339 29-6584 27-6606 26-6634 24-6647 31-7968 30-6662 25-6853	56-7617 57-7657 61-7745 59-7806 58-7859 62-7890 63-8715 60-7978	

Fig. 12: Instruction monitoring for the load instructions on *SM 0* in the cycle range of [5000, 30000] for *PR-K1*.

8 Figure 12 presents the instruction issue/completion cycles of 8 *K1* thread blocks run-  
 9 ning on *SM 0*. We merge multiple snapshots of our *Temporal View* that belong to thread  
 10 blocks in *SM 0* to evaluate the performance of all load instructions together. The first and  
 11 second lines point to the load instructions of which the program counter (PC) equals 296

- 1 (loads DP) and 312 (loads SP), respectively. Figure 13, a snapshot of our *Spatial View*,  
 2 demonstrates the memory access statistics of representative components within the same interval. After the kernel launch, each thread collects thread-specific information from

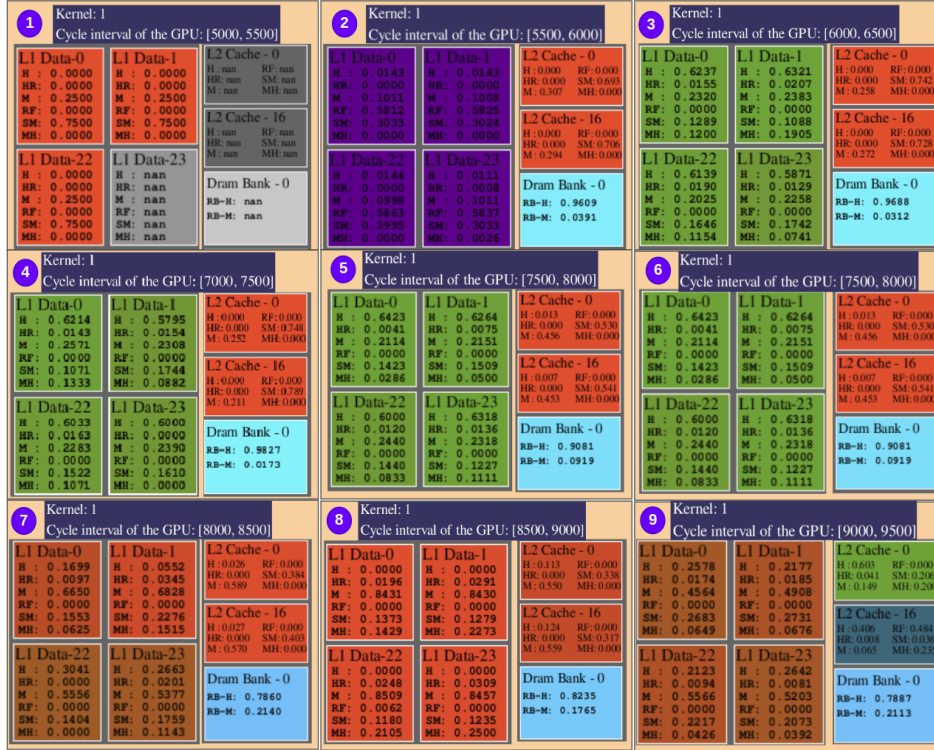


Fig. 13: Memory performance overview in the cycle range of [5000, 9500] for *PR-K1*.

- 3 parameter memory, which takes 250-450 cycles to process target data addressed with  
 4 the thread's private registers. The warp schedulers dispatch the load instructions (i.e., at  
 5 PC=296 *ld.global.u64*), and all eight warps at *Thread Block 24* start executing the instruc-  
 6 tion after *Cycle 5455*. Furthermore, Figure 12 reveals that SM dispatches load instructions  
 7 from the remaining thread blocks in the interval of [5470, 5786] after issuing the load in-  
 8 structions of *Thread Block 24*. Figure 13 reveals that no access occurs on some of the L1D  
 9 caches, while none of the L2 caches and DRAM banks are accessed during the preparation  
 10 time in [5000, 5500] in Part [1]. No data brought to the L1D cache of *SM 0* by the warps of  
 11 *Thread Block 24* after *Cycle 6087* (*Warp 6*) enables the early completion of the instruc-  
 12 tion pointed with PC=296, belonging to the thread blocks 104, 184, 264, 344, 424, 504, and  
 13 584. We highlight this observation with the bold fonts representing the earliest completion  
 14 times within each thread block in the second row of the first instructions. Additionally,  
 15 a high reservation failure and no MSHR hit rates on the L1D cache of *SM 0* in Part [2]  
 16 confirms that the locality utilization between thread blocks is quite low for the first load. If  
 17 the L1D cache locality was utilized, we should have observed larger MSHR hit rates and  
 18 completion of the same instructions just after *Cycle 6087*. Parts [2,3,4,5] reveal that the

1 reservation failed requests pointed by PC=296 cause a miss on L2 caches without MSHR  
 2 merging. Thus, memory requests of the same instruction from different SMs cannot ben-  
 3 efit from L2 locality and cause more traffic in the memory hierarchy. Additionally, Parts  
 4 3,4,5,6 reveal that the L1 access status mostly turns to the hit after *Cycle 6000*. Unlike  
 5 the load instructions at PC=296, ones at PC=312 (*ld.global.u32*) usually hit on L1. The  
 6 second line for each thread block shown in Figure 12 indicates that the completion takes  
 7 much fewer cycles for the loads at PC=312. To illustrate, while *Thread Block 504* com-  
 8 pletes the first load instructions within 2133 cycles, except *Warp 63*, it takes 26 cycles for  
 9 the second instructions, whose requests result in a miss on both the L1D and L2 caches.  
 10 While the loads at PC=296 complete the execution in the range of [350, 2250] cycles,  
 11 the loads at PC=312 take less than 50 cycles for most of the warps due to the increasing  
 12 hits on the L1D cache. However, the loads at PC=296 delay the issue of the second load  
 13 instructions due to long latency. The remaining thread blocks (from *Thread Block 641*)  
 14 are assigned to SMs after *Cycle 55000*. With the observation that a thread block occupies  
 15 50000 cycles on an SM (for the web-Stanford graph, which can easily fit into DRAM),  
 16 the schedule of any waiting thread block delays around 2000 cycles. Such delays affect  
 17 the performance of a thread block by 4% in addition to affecting the memory traffic neg-  
 18 atively and degrading the overall performance of *K1*. In this manner, an approach such as  
 19 adaptive thread block scheduling by throttling the load/store unit issue amount depend-  
 20 ing on the access statistics of caches can reduce the side effects and increase the overall  
 21 performance.

Table 3: Dissipated average power in milliWatt in the cycle range of [5000,10000] for *PR-K1*.

Cycles	Streaming Multiprocessors					Memory Partitions						NoC	TOTAL
	Exec. Units	Func. Units	LD/ST Unit	Idle	Total	MC FEE	PHY	MC TE	Dram	L2	Total		
5k, 5.5k	2637.5	54.3	35.6	23.7	2751.3	3.7	8.2	4.6	0	0	16.5	0.7	2768.5
5.5k, 6k	597	6.3	860	0	1463.7	177.2	17.8	9.4	557.2	3.4	764.9	26.9	2501.6
6k, 6.5k	614.4	12.4	399.9	0	1026.7	56.1	31.3	16.1	1346.4	3.0	1452.9	92.6	2577.3
6.5k, 7k	708.6	14.4	464.3	0	1187.3	65.5	31.4	16.2	1354.3	3.1	1470.5	94.2	2755.1
7k, 7.5k	686.8	13.9	463.9	0	1164.6	65.6	31.8	16.2	1354.9	3.1	1471.2	94.4	2733.2
7.5k, 8k	795.8	16.3	487.4	0	1299.4	69.9	29.7	15.3	1264.3	3.7	1383.1	96.7	2782.9
8k, 8.5k	543	10.2	335	0	888.2	60.4	30.5	15.7	1341.4	4.4	1452.0	124.7	2469.3
8.5k, 9k	354.5	5.6	249.3	0	609.3	52	31.1	16.1	1362.6	19	1480.7	148	2257.2
9k, 9.5k	474.9	5.3	455.4	0	935.7	80.2	27.8	14.4	1096.9	66.1	1284.4	216.8	2503
9.5k, 10k	446.8	4.8	475.5	0	927	78.2	23.8	12.4	843.2	41	968.7	153.8	2090.4

22 Table 3 presents power measurements for *K1* execution throughout the hardware struc-  
 23 tures. After kernel launch, (i.e., *Cycle 5000*), the threads load thread-identifier data from  
 24 the parameter memory, causing higher power consumption on SMs. The power values of  
 25 the memory partitions in the following cycles get higher than the SMs. Additionally, the  
 26 power consumption by the LD/ST unit is high due to the intense memory operations and  
 27 pressure on the L1D cache, and other units apart from the register file portion of the exe-  
 28 cution unit get lower after *Cycle 5500*. We can see that DRAM consumes the most power  
 29 in the memory partitions with intense usage of high-bandwidth [24].

### 4.3. Performance-Power Analysis of an Embedded Application

Embedded applications targeting artificial intelligence, computer vision, and advanced graphics require high computational power. Jetson AGX Xavier provides a System-on-Module that meets these demands with a Volta-based GPU. We execute CUDA implementation of the *Fast Fourier Transform (FFT)* from the GPU4S suite [37].

Table 4 presents the overall performance metrics for executing the first three kernels. The thread blocks consist of 256 threads, and a maximum of 8 thread blocks can run in parallel on the SMs since the register file limits the number of simultaneous thread block execution. Since the kernels other than the *K0* utilize the hardware similarly, they result in similar performance as shown for *K1* and *K2*. Thus, we explain the relationship between performance and power consumption for *K0*, which employs diverse characteristics.

Table 4: Fast Fourier kernel performance statistics.

Kernel	GPU IPC	GPU Occ.	L1D			L2			DRAM		Total Cycle
			Miss Rate	Res. Rate	Fail	Miss Rate	Res. Rate	Fail	Row Buffer Loc. (L+S)	Row Buffer Loc. (Load)	
FFT - K0	31.76	62.52 %	0.67	0.92		0.65	0		0.21	0.89	990276
FFT - K1	100.22	80.25 %	0.1	0.75		0.20	0		0.33	0.92	235390
FFT - K2	49.20	85.28 %	0.1	0.837		0.21	0		0.41	0.92	239755

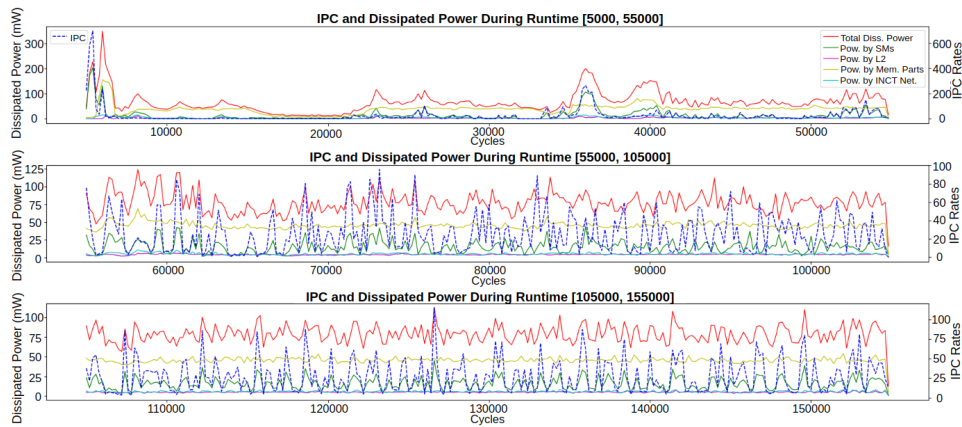


Fig. 14: IPC with dissipated power metrics in the cycle range of [5000, 155000] for *FFT-K0*.

Figure 14 displays the IPC and power metrics measurements for *K0* at different execution cycle intervals. To avoid losing observation details related to IPC and power metrics of the kernel, we report the relationship in three execution intervals separately. At cycles [5000, 6500], the power per cycle and IPC values increase significantly since each SM

1 registers thread-identifier information in their private registers, causing a high activation  
 2 on register files and functional units. Figure 15, (as part of our *General View*), shows four  
 3 loads, three data movements, and one multiply-add instruction executed by 16384 threads  
 4 concurrently (with 64 thread blocks where each contains 256 threads) in that interval.  
 5 Throughout the *K0* execution, we have high L1D and L2 miss rates. Figure 14 reveals the  
 6 effect of those accesses on power dissipation (yellow line). The memory partitions dissi-  
 7 pate half of the total power, around 50W, due to intensive activation at runtime. IPC and  
 8 power dissipation increase instantly with SM activation points. By tracking the points,  
 9 where IPC increases between [55000, 105000] and [105000, 155000] cycles, we can see  
 10 parallel increments in power metrics dissipated by SMs and GPU. While the overall IPC  
 11 value for *K0* equals 31.76 on SMs, runtime IPC varies in the range of [5, 75]. The concur-  
 12 rent load and store instructions cause small latencies and slight computational loads on  
 13 the SMs.

PC	OPCODE	OPERAND
0	ld.param.u64	%rd1 [_Z21binary_reverse_kernelPKfPfli_param_0]
8	ld.param.u64	%rd2 [_Z21binary_reverse_kernelPKfPfli_param_1]
16	ld.param.u64	%rd3 [_Z21binary_reverse_kernelPKfPfli_param_2]
24	ld.param.u32	%r2 [_Z21binary_reverse_kernelPKfPfli_param_3]
32	mov.u32	%r3 %ntid.x
40	mov.u32	%r4 %ctaid.x
48	mov.u32	%r5 %tid.x
56	mad.lo.s32	%r1 %r3 %r4 %r5

Fig. 15: Instructions preparing threads for execution with thread-specific data in the cycle range of [5000, 6500] for *FFT-K0*.

#### 14 4.4. Analyzing the Impact of the Input Size on Resource-Critical Embedded 15 System

16 Since embedded devices employ relatively smaller computational and memory resources,  
 17 the target programs' resource utilization becomes more important and influential on the  
 18 execution performance. Besides program characteristics, the input size significantly af-  
 19 fects the pressure on the memory structures and the obtained performance. To understand  
 20 the impacts of the input size on memory behavior and execution performance, we utilize  
 21 the memory statistics view of our tool. We execute *softmax* program from GPU4S suite  
 22 [37] for two matrix sizes (1024 and 4096).

23 Figure 16 presents the memory statistics of the first 500K cycles. The execution em-  
 24 ploys high L1 reservation failure rates and L2 miss rates for the first 300K cycles for both  
 25 executions. However, the small input (1024) can fit within the L2 and L1 caches around  
 26 *Cycle 300K* and *Cycle 350K*, respectively. On the other hand, the large input (4096) ex-  
 27 ecution still employs low hit L1 and L2 rates.

28 Figure 17 demonstrates the corresponding IPC behavior, which is compatible with  
 29 cache rates. While the IPC for the large input case oscillates at small values (0.2-0.8), the  
 30 IPC values for the small input case are steadily low but increase after L2 and L1 cache

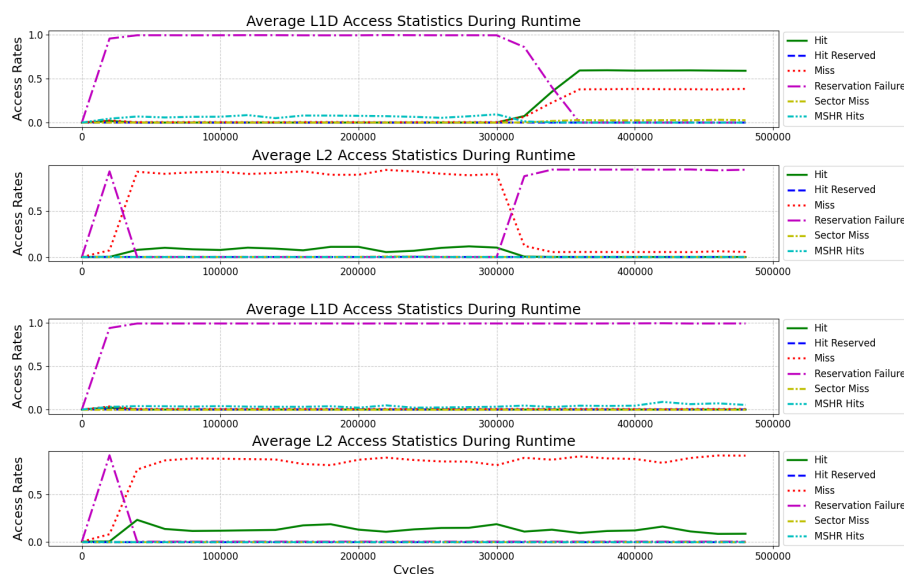


Fig. 16: Average memory access statistics at first 500K cycles for *softmax* kernel (matrix sizes with 1024 and 4096).

1 breakpoints. With this analysis, we can decide on the cache utilization of the target input  
 2 dimension for the program.

### 3 5. Related Work

4 AerialVision [1] visualizes runtime warp divergence, dynamic IPC, global memory access  
 5 statistics, active thread count, and a mapping window between source code and exposed  
 6 pipeline latency metrics of kernel execution. While GPPRMon enables developers to dig  
 7 into details of runtime GPU execution within specific microarchitectural units, AerialVi-  
 8 sion profiles the run time metrics on a much longer execution scale and visualizes overall  
 9 GPU performance without per-component performance analysis. GPPRMon is similar to  
 10 AerialVision in terms of providing runtime performance metrics, but GPPRMon can un-  
 11 cover more detailed behavioral insight inside kernel execution with more detailed runtime  
 12 observation opportunities. Furthermore, AerialVision does not support displaying dissipated  
 13 power.

14 Nsight Compute Tool [27] runs an application on an NVIDIA GPU device and col-  
 15 lects average hardware usage statistics for the main components on a kernel basis. Nsight  
 16 System Tool [28] is the other tool that is mostly preferred to analyze end-to-end ker-  
 17 nel execution performance, including CPU-GPU communications. Profiling through these  
 18 tools eases interpreting the overall kernel performance and hardware utilization with well-  
 19 designed GUIs. In addition to performance analysis tools, GPU users can track instant  
 20 power dissipation of deployed GPU through *System Management Interface (SMI)* li-  
 21 brary (i.e., *nvidia-smi*) [31]. Different from all NVIDIA tools and library frameworks,  
 22 GPPRMon can track and visualize execution and hardware utilization statistics at run-

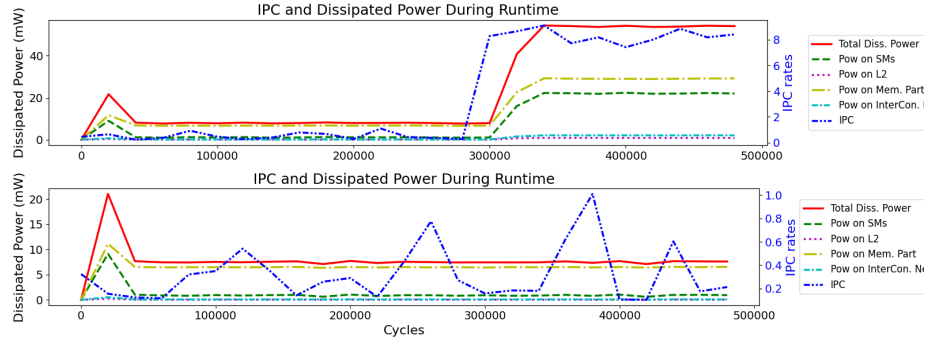


Fig. 17: IPC values for first 500K cycles for *softmax* kernel (matrix sizes with 1024 and 4096).

1 time for each component separately. GPPRMon further monitors warp-instruction is-  
 2 sues/completions for each thread block, which provides an on-point analysis opportunity  
 3 in a highly parallel execution. Hence, GPPRMon provides multi-functional performance  
 4 and power tracking options together with various configurations, and it supports these  
 5 features for all official GPU models described in Section 2.3.

6 TAU Performance System [39] is a performance profiling tool for hybrid parallel pro-  
 7 grams such as CUDA and OpenCL by intercepting the execution and inserting metric  
 8 collection calls. After gathering the performance results, TAU integrates them with data  
 9 through instrumentation to display a performance picture of the execution. Like Nsight  
 10 Systems Tool [28], which provides performance traces for high-level CUDA functions  
 11 such as *cudaMemcpy*, TAU does not address detailed hardware usage and performance  
 12 during runtime as GPPRMon does. While TAU is insufficient to reveal the hardware's  
 13 energy consumption, GPPRMon provides runtime power breakdown of GPU subcompo-  
 14 nents during the execution.

15 Daisen [43] displays the overall occupancy on SM pipeline stages and memory compo-  
 16 nents during the simulation. The authors aim to propose a performance-improving ar-  
 17 chitecture by iterating an algorithm that benefits from the previously collected perfor-  
 18 mance and hardware usage metrics. In other words, Daisen does not highlight the per-  
 19 formance degrading points analytically. Instead, their approach addresses general bottle-  
 20 necks, shares with the user, and tries to optimize performance in each iteration. Similar to  
 21 [33], their approach offers more systematic performance optimization points on GPU ex-  
 22 ecutions. On the other hand, GPPRMon focuses on monitoring the execution at the PTX  
 23 instruction level and relating the execution with the memory occupancy during execution.  
 24 That is, GPPRMon offers a runtime analytical observation environment to evaluate both  
 25 architecture and application inefficiencies. Moreover, while GPPRMon supports runtime  
 26 power tracking and overall energy consumption metrics, Daisen does not provide them.

27 CHAMPVis [33] offers a web-supported architectural performance monitoring tool  
 28 that provides a hierarchical analysis of trends and bottlenecks for varying applications.  
 29 The tool aims to analyze performance by evaluating metrics from a system view and auto-  
 30 matically generates predictive optimization speculations for the applications. Unlike GP-  
 31 PPRMon, CHAMPVis does not employ any dissipated power tracking. GPPRMon yields  
 32 detailed execution statistics, whereas CHAMPVis highlights the overall execution trends.



1 Francisco et al. [5] model a portion of the memory hierarchy of the AMD GPUs ac-  
2 curately by investigating the behavior of MSHRs, coalescing for vector (warp for our  
3 case) memory requests by extending Multi2Sim [45]. The authors find that the size and  
4 switching frequencies of MSHRs affect performance directly, especially for irregular  
5 workloads. Additionally, coalesced memory accesses, implemented in the simulator, re-  
6 duce the repeated overheads of memory requests, significantly affecting performance for  
7 global memory accesses. GPPRMon offers more comprehensive evaluation opportunity  
8 besides the impact of MSHRs and memory request coalescing by being capable of run-  
9 time monitoring of both all memory hierarchy and SMs. Furthermore, GPPRMon extends  
10 relating the performance analysis options together with the dissipated power for perfor-  
11 mance/power trade-off evaluations at any execution scale.

12 Tanzima et al. [11] present an analysis and visualization framework (DASHING) tar-  
13 getting exascale computing consisting of multi-core architectures. The authors provide  
14 user interaction to compare varying configuration models by providing environment con-  
15 figuration options for analysis and visualization. Similarly, GPPRMon supports official  
16 GPU configurations of the GPGPU-Sim for multiple performance analyses and includes  
17 multi-functional metric collection and visualization perspectives depending on the user's  
18 demand. However, we obtain a more low-level performance analysis tool as described  
19 in *Temporal* and *Spatial View* Sections. In addition to runtime monitoring support, GP-  
20 PPRMon allows tracking the power dissipation across microarchitectural GPU components  
21 at runtime, which enables evaluating performance/power together, especially for energy-  
22 critical applications.

23 MemAxes [8] proposes an analysis framework for memory performance with various  
24 inspections on multi-core architectures. Its interface displays the analysis by obtaining  
25 performance metric samples from different simulations and mapping them into a single  
26 visual. In addition, the authors offer memory utilization-based clustering research among  
27 the benchmark applications. On the contrary, GPPRMon enables the analysis of GPU  
28 execution at runtime with performance and power dissipation features through *Spatial*,  
29 *Temporal*, and *General Views*.

30 To the best of our knowledge, GPPRMon differs from existing works by targeting  
31 hardware from both embedded and large-scale GPUs and providing runtime performance  
32 and power analysis opportunities. GPPRMon offers a detailed micro-architectural and  
33 power dissipation analysis for any GPU application. With observations through GPPRMon,  
34 developers can identify performance and power issues for both applications and archi-  
35 tectures. For instance, data sparsity, irregular memory access behaviors, compute-bound  
36 behaviors, and the execution interval of these bottlenecks can be easily recognized via  
37 GPPRMon. Moreover, since GPPRMon supports runtime power dissipation together with  
38 performance monitoring, developers can analyze the energy impact of those performance  
39 issues and make their decisions. The studies conducting performance-energy tradeoff  
40 analysis [2, 38] and power capping strategies accordingly [49, 18] can potentially ben-  
41 efit from our fine-grained dynamic performance and power evaluations. We believe GP-  
42 PPRMon will be useful in fulfilling the micro-architectural performance and power analysis  
43 of GPUs for diverse application domains.

## 1 **6. Conclusion**

2 To conclude, we design and build a systematic runtime metric collection of instruction  
 3 monitoring, performance, memory access, and power consumption metrics. Our tool pro-  
 4 vides a multi-perspective visualization framework that displays performance, execution  
 5 statistics of the workload, occupancy of the memory hierarchy, and dissipated power re-  
 6 sults to conduct baseline analysis on GPUs at runtime. Since GPPRMon reliably reveals  
 7 all the interactions between hardware and application at runtime, it potentially helps the  
 8 researchers and software developers understand the dynamic behavior of the target GPU  
 9 execution. While the researchers can propose architectural optimizations based on the  
 10 detailed information, the software developers can deal with performance bottlenecks by  
 11 analyzing our visualizations and fixing the target GPU code accordingly. Additionally,  
 12 the low-power embedded system developers can utilize dynamic performance and power  
 13 considerations to balance between two important design points.

14 We believe that GPPRMon will help to conduct baseline analysis for the literature con-  
 15 cerning GPU performance and power dissipation and eliminate the need for additional in-  
 16 house efforts that involve real-time monitoring and profiling support. Since GPPRMon’s  
 17 metric collection and visualization components are independent, other microarchitectural  
 18 metrics obtained during runtime from either GPUs or other simulators can be displayed by  
 19 exploiting the GPPRMon visualizer. The integration of other simulators with GPPRMon  
 20 results in a system capable of runtime execution, memory statistics, and power dissipation  
 21 tracker among all the possible GPUs. Alternatively, extending the runtime visualization  
 22 by deepening the scope of the runtime metric collection and visualization phases of GP-  
 23 PPRMon is another future direction.

## 24 **7. Acknowledgement**

25 This work was supported by the Scientific and Technological Research Council of Turkey  
 26 (TÜBİTAK), Grant No: 122E395. This work is partially supported by CERCIRAS COST  
 27 Action CA19135 funded by COST Association.

## 28 **References**

- 29 1. Ariel, A., Fung, W.W.L., Turner, A.E., Aamodt, T.M.: Visualizing complex dynamics in many-  
 30 core accelerator architectures. In: 2010 IEEE International Symposium on Performance Anal-  
 31 ysis of Systems & Software (ISPASS). pp. 164–174 (2010)
- 32 2. Aslan, B., Yilmazer-Metin, A.: A study on power and energy measurement of nvidia jetson em-  
 33 bedded gpus using built-in sensor. In: 2022 7th International Conference on Computer Science  
 34 and Engineering (UBMK). pp. 1–6 (2022)
- 35 3. del Barrio, V., Gonzalez, C., Roca, J., Fernandez, A., E, E.: Attila: a cycle-level execution-  
 36 driven simulator for modern gpu architectures. In: 2006 IEEE International Symposium on  
 37 Performance Analysis of Systems and Software. pp. 231–241 (2006)
- 38 4. Becker, P.H.E., Arnau, J.M., González, A.: Demystifying power and performance bottlenecks  
 39 in autonomous driving systems. In: 2020 IEEE International Symposium on Workload Charac-  
 40 terization (IISWC). pp. 205–215 (2020)
- 41 5. Candell, F., Petit, S., Sahuquillo, J., Duato, J.: Accurately modeling the gpu memory subsystem.  
 42 In: 2015 International Conference on High Performance Computing & Simulation (HPCS). pp.  
 43 179–186 (2015)

- 1 6. Chen, X., Chang, L.W., Rodrigues, C.I., Lv, J., Wang, Z., Hwu, W.M.: Adaptive cache man-  
2 agement for energy-efficient gpu computing. In: 2014 47th Annual IEEE/ACM International  
3 Symposium on Microarchitecture. pp. 343–355 (2014)
- 4 7. Collange, C., Dumas, M., Defour, D., Parello, D.: Barra: A parallel functional simulator for  
5 gpgpu. In: 2010 IEEE International Symposium on Modeling, Analysis and Simulation of  
6 Computer and Telecommunication Systems. pp. 351–360 (2010)
- 7 8. Giménez, A., Gamblin, T., Jusufi, I., Bhatele, A., Schulz, M., Bremer, P.T., Hamann, B.:  
8 Memaxes: Visualization and analytics for characterizing complex memory performance be-  
9 haviors. *IEEE Transactions on Visualization and Computer Graphics* 24(7), 2180–2193 (2018)
- 10 9. Guerreiro, J., Ilic, A., Roma, N., Tomás, P.: Dvfs-aware application classifica-  
11 tion to improve gpgpus energy efficiency. *Parallel Computing* 83, 93–117 (2019),  
12 <https://www.sciencedirect.com/science/article/pii/S0167819118300243>
- 13 10. Hong, J., Cho, S., Kim, G.: Overcoming memory capacity wall of gpus with heterogeneous  
14 memory stack. *IEEE Computer Architecture Letters* 21(2), 61–64 (2022)
- 15 11. Islam, T., Ayala, A., Jensen, Q., Ibrahim, K.: Toward a programmable analysis and visualization  
16 framework for interactive performance analytics. In: 2019 IEEE/ACM International Workshop  
17 on Programming and Performance Visualization Tools (ProTools). pp. 70–77 (2019)
- 18 12. Jain, P., Jain, A., Nrusimha, A., Gholami, A., Abbeel, P., Keutzer, K., Stoica, I., Gonzalez,  
19 J.E.: Checkmate: Breaking the memory wall with optimal tensor rematerialization. *CoRR*  
20 [abs/1910.02653](https://arxiv.org/abs/1910.02653) (2019), <http://arxiv.org/abs/1910.02653>
- 21 13. Jog, A., Kayiran, O., Nachiappan, N.C., Mishra, A.K., Kandemir, M.T., Mutlu, O., Iyer, R.R.,  
22 Das, C.R.: OWL: cooperative thread array aware scheduling techniques for improving GPGPU  
23 performance. In: Sarkar, V., Bodík, R. (eds.) *Architectural Support for Programming Lan-  
24 guages and Operating Systems, ASPLOS 2013*, Houston, TX, USA, March 16–20, 2013. pp.  
25 395–406. ACM (2013), <https://doi.org/10.1145/2451116.2451158>
- 26 14. Kandiah, V., Peverelle, S., Khairy, M., Pan, J., Manjunath, A., Rogers, T.G., Aamodt,  
27 T.M., Hardavellas, N.: Accelwattch: A power modeling framework for modern gpus. In:  
28 MICRO-54: 54th Annual IEEE/ACM International Symposium on Microarchitecture. p.  
29 738–753. MICRO '21, Association for Computing Machinery, New York, NY, USA (2021),  
30 <https://doi.org/10.1145/3466752.3480063>
- 31 15. Khairy, M., Shen, Z., Aamodt, T.M., Rogers, T.G.: Accel-sim: An extensible simulation frame-  
32 work for validated gpu modeling. In: 2020 ACM/IEEE 47th Annual International Symposium  
33 on Computer Architecture (ISCA). pp. 473–486 (2020)
- 34 16. Koo, G., Oh, Y., Ro, W.W., Annavaram, M.: Access pattern-aware cache management for im-  
35 proving data utilization in gpu. In: 2017 ACM/IEEE 44th Annual International Symposium on  
36 Computer Architecture (ISCA). pp. 307–319 (2017)
- 37 17. Krzywaniak, A., Czarnul, P., Proficz, J.: Gpu power capping for energy-performance trade-offs  
38 in training of deep convolutional neural networks for image recognition. In: *Computational  
39 Science – ICCS 2022*. pp. 667–681. Springer International Publishing, Cham (2022)
- 40 18. Krzywaniak, A., Czarnul, P., Proficz, J.: Dynamic gpu power capping  
41 with online performance tracing for energy efficient gpu computing us-  
42 ing depo tool. *Future Generation Computer Systems* 145, 396–414 (2023),  
43 <https://www.sciencedirect.com/science/article/pii/S0167739X23001267>
- 44 19. Leskovec, J., Lang, K., Dasgupta, A., Mahoney, M.: Community structure in large networks:  
45 Natural cluster sizes and the absence of large well-defined clusters. *Internet Mathematics* 6 (11  
46 2008)
- 47 20. Lew, J., Shah, D.A., Pati, S., Cattell, S., Zhang, M., Sandhupatla, A., Ng, C., Goli, N., Sinclair,  
48 M.D., Rogers, T.G., Aamodt, T.M.: Analyzing machine learning workloads using a detailed  
49 gpu simulator. In: 2019 IEEE International Symposium on Performance Analysis of Systems  
50 and Software (ISPASS). pp. 151–152 (2019)

- 1 21. Li, S., Ahn, J.H., Strong, R.D., Brockman, J.B., Tullsen, D.M., Jouppi, N.P.: Mcpat: An integrated power, area, and timing modeling framework for multicore and manycore architectures. In: 2009 42nd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO). pp. 469–480 (2009)
- 2 22. Mittal, S., Vetter, J.S.: A survey of methods for analyzing and improving gpu energy efficiency. *ACM Comput. Surv.* 47(2) (aug 2014), <https://doi.org/10.1145/2636342>
- 3 23. NVIDIA: Quadro gv100 data sheet (2018), <https://www.nvidia.com/content/dam/en-zz/Solutions/design-visualization/productspage/quadro/quadro-desktop/quadro-volta-gv100-a4-nvidia-704619-r3-web.pdf>
- 4 24. NVIDIA: Volta architecture white paper (March 2018), <https://images.nvidia.com/content/volta-architecture/pdf/volta-architecture-whitepaper.pdf>
- 5 25. NVIDIA: Jetson agx xavier and the new era of autonomous machines (2019), [https://info.nvidia.com/rs/156-OFN-742/images/Jetson\\_AGX\\_Xavier\\_New\\_Era\\_Autonomous\\_Machines.pdf](https://info.nvidia.com/rs/156-OFN-742/images/Jetson_AGX_Xavier_New_Era_Autonomous_Machines.pdf)
- 6 26. NVIDIA: Cuda toolkit documentation (Jan 2023), <https://docs.nvidia.com/cuda>
- 7 27. NVIDIA: Nsight compute kernel profiling guide (2024), <https://docs.nvidia.com/nsight-compute/ProfilingGuide/index.html>
- 8 28. NVIDIA: Nsight systems profiling guide (2024), <https://developer.nvidia.com/nsight-systems>
- 9 29. NVIDIA: Nvidia cuda compiler driver nvcc (2024), <https://docs.nvidia.com/cuda/cuda-compiler-driver-nvcc/>
- 10 30. NVIDIA: Parallel thread execution isa version 8.4 (2024), <https://docs.nvidia.com/cuda/parallel-thread-execution/index.html>
- 11 31. NVIDIA: System management interface (2024), <https://developer.nvidia.com/system-management-interface>
- 12 32. O’Neil, M.A., Burtscher, M.: Microarchitectural performance characterization of irregular gpu kernels. In: 2014 IEEE International Symposium on Workload Characterization (IISWC). pp. 130–139 (2014)
- 13 33. Pentecost, L., Gupta, U., Ngan, E., Beyer, J., Wei, G.Y., Brooks, D., Behrisch, M.: Champvis: Comparative hierarchical analysis of microarchitectural performance. In: 2019 IEEE/ACM International Workshop on Programming and Performance Visualization Tools (ProTools). pp. 55–61 (2019)
- 14 34. Power, J., Hestness, J., Orr, M.S., Hill, M.D., Wood, D.A.: gem5-gpu: A heterogeneous cpu-gpu simulator. *IEEE Computer Architecture Letters* 14(1), 34–36 (2015)
- 15 35. PowerUP, T.: V100 tech powerup (2018), <https://www.techpowerup.com/gpu-specs/tesla-v100-pcie-16-gb.c2957>
- 16 36. Rhu, M., Sullivan, M., Leng, J., Erez, M.: A locality-aware memory hierarchy for energy-efficient gpu architectures. In: Proceedings of the 46th Annual IEEE/ACM International Symposium on Microarchitecture. p. 86–98. MICRO-46, Association for Computing Machinery, New York, NY, USA (2013), <https://doi.org/10.1145/2540708.2540717>
- 17 37. Rodriguez, I., Kosmidis, L., Lachaize, J., Notebaert, O., Steenari, D.: Gpu4s bench: Design and implementation of an open gpu benchmarking suite for space on-board processing. *Universitat Politècnica de Catalunya* (2019)
- 18 38. Sezgin, Y., Öz, I.: Performance-reliability tradeoff analysis for safety-critical embedded systems with gpus. In: Ulusal Yüksek Başarımlı Hesaplama Konferansı (BAŞARIM) (2024), <https://indico.truba.gov.tr/event/140/attachments/310/642/BASARIM2024-BildiriKitabi.pdf>
- 19 39. Shende, S.S., Malony, A.D.: The tau parallel performance system. *Int. J. High Perform. Comput. Appl.* 20(2), 287–311 (may 2006), <https://doi.org/10.1177/1094342006064482>
- 20 40. Shi, X., Zheng, Z., Zhou, Y., Jin, H., He, L., Liu, B., Hua, Q.S.: Graph processing on gpus: A survey. *ACM Comput. Surv.* 50(6) (jan 2018), <https://doi.org/10.1145/3128571>
- 21 41. Strohmaier, E., Dongarra, J., Simon, H., Meuer, M.: Highlights - november 2022 (2022), <https://www.top500.org/lists/top500/2022/11/highs/>

- 1 42. Sun, Y., Mukherjee, S., Baruah, T., Dong, S., Gutierrez, J., Mohan, P., Kaeli, D.: Evaluat-  
2 ing performance tradeoffs on the radeon open compute platform. In: 2018 IEEE International  
3 Symposium on Performance Analysis of Systems and Software (ISPASS). pp. 209–218 (2018)
- 4 43. Sun, Y., Zhang, Y., Mosallaei, A., Shah, M.D., Dunne, C., Kaeli, D.R.: Daisen: A frame-  
5 work for visualizing detailed GPU execution. *Comput. Graph. Forum* 40(3), 239–250 (2021),  
6 <https://doi.org/10.1111/cgf.14303>
- 7 44. Topçu, B., Öz, I.: Gpprmon: Gpu runtime memory performance and power monitoring tool. In:  
8 Euro-Par 2023: Parallel Processing Workshops. pp. 17–29. Springer Nature Switzerland (2024)
- 9 45. Ubal, R., Jang, B., Mistry, P., Schaa, D., Kaeli, D.: Multi2sim: A simulation framework for  
10 cpu-gpu computing. In: 2012 21st International Conference on Parallel Architectures and Com-  
11 pilation Techniques (PACT). pp. 335–344 (2012)
- 12 46. Vijaykumar, N., Ebrahimi, E., Hsieh, K., Gibbons, P.B., Mutlu, O.: The locality descriptor:  
13 A holistic cross-layer abstraction to express data locality in gpus. In: 2018 ACM/IEEE 45th  
14 Annual International Symposium on Computer Architecture (ISCA). pp. 829–842 (2018)
- 15 47. Wang, Y., Pan, Y., Davidson, A., Wu, Y., Yang, C., Wang, L., Osama, M., Yuan, C., Liu, W.,  
16 Riffel, A.T., Owens, J.D.: Gunrock: Gpu graph analytics. *ACM Trans. Parallel Comput.* 4(1)  
17 (aug 2017), <https://doi.org/10.1145/3108140>
- 18 48. Xu, Z., Chen, X., Shen, J., Zhang, Y., Chen, C., Yang, C.: Gardenia: A graph processing bench-  
19 mark suite for next-generation accelerators. *ACM Journal on Emerging Technologies in Com-  
20 puting Systems* 15(1) (jan 2019), <https://doi.org/10.1145/3283450>
- 21 49. Zhao, D., Samsi, S., McDonald, J., Li, B., Bestor, D., Jones, M., Tiwari, D., Gadepally, V.:  
22 Sustainable supercomputing for ai: Gpu power capping at hpc scale. In: Proceedings of the 2023  
23 ACM Symposium on Cloud Computing. p. 588–596. SoCC '23, Association for Computing  
24 Machinery, New York, NY, USA (2023), <https://doi.org/10.1145/3620678.3624793>
- 25 50. Zhao, X., Adileh, A., Yu, Z., Wang, Z., Jaleel, A., Eeckhout, L.: Adaptive memory-side last-  
26 level gpu caching. In: 2019 ACM/IEEE 46th Annual International Symposium on Computer  
27 Architecture (ISCA). pp. 411–423 (2019)