

Modeling Emerging Memory-Divergent GPU Applications

Lu Wang¹, Magnus Jahre², Almutaz Adileh³,
Zhiying Wang, and Lieven Eeckhout¹

Abstract—Analytical performance models yield valuable architectural insight without incurring the excessive runtime overheads of simulation. In this work, we study contemporary GPU applications and find that the key performance-related behavior of such applications is distinct from traditional GPU applications. The key issue is that these GPU applications are memory-intensive and have poor spatial locality, which implies that the loads of different threads commonly access different cache blocks. Such memory-divergent applications quickly exhaust the number of misses the L1 cache can process concurrently, and thereby cripple the GPU's ability to use Memory-Level Parallelism (MLP) and Thread-Level Parallelism (TLP) to hide memory latencies. Our Memory Divergence Model (MDM) is able to accurately represent this behavior and thereby reduces average performance prediction error by 14× compared to the state-of-the-art GPUMech approach across our memory-divergent applications.

Index Terms—Analytical performance prediction, memory divergence model, GPU

1 INTRODUCTION

QUANTITATIVE evaluation is an essential part of the computer architect's tool box. Simulation is the most common evaluation tool since it enables detailed, even cycle-accurate, performance analysis. However, simulation is excruciatingly slow and hence parameter sweeps commonly require thousands of CPU hours. An alternative approach is analytical modeling, which captures the key performance-related behavior of the architecture with a set of mathematical equations. Analytical models are much faster than simulation—making them ideally suited for early-stage architectural exploration [1], [2] and helping programmers understand application performance [3], [4].

GPUs are the de facto standard platform for executing performance-critical applications. Their highly parallel execution model and high-performance memory system makes GPUs a popular choice for emerging applications such as data analytics [5], [6]. The diversity of modern-day GPU applications makes them challenging to model. Several contemporary GPU applications differ from traditional GPU-compute workloads because they put a much larger strain on the memory system. More specifically, they are memory-intensive and memory-divergent—i.e., the memory accesses from concurrently executing threads map to multiple cache lines. While simulators account for this behavior by modeling cycle-by-cycle activities, state-of-the-art GPU modeling approaches are unable to predict performance with sufficient accuracy.

Our objective is to provide an analytical performance model for GPUs that is able to accurately predict the performance of the

- L. Wang, A. Adileh, and L. Eeckhout are with the Department of Electronics and Information Systems (ELIS), Ghent University, Gent 9000, Belgium. E-mail: {luluwang.wang, almutaz.adileh, lieven.eeckhout}@ugent.be.
- M. Jahre is with the Department of Computer Science, Norwegian University of Science and Technology, Trondheim 7012, Norway. E-mail: magnus.jahre@ntnu.no.
- Z. Wang is with the School of Computer, National University of Defense Technology, Changsha 410073, P.R. China. E-mail: zyuwang@nudt.edu.cn.

Manuscript received 21 May 2019; revised 14 June 2019; accepted 14 June 2019. Date of publication 17 June 2019; date of current version 10 July 2019.
(Corresponding author: Lu Wang.)

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below.

Digital Object Identifier no. 10.1109/LCA.2019.2923618

various GPU applications, including divergent memory-intensive applications. Our starting point is interval modeling [7], which is a widely used approach for CPU performance evaluation. The key observations are that an application will have a certain steady-state performance in the absence of miss events (e.g., data cache misses), and that miss events are independent of each other. Therefore, performance can be predicted by predicting steady-state performance and subtracting the performance loss due to each miss event. GPUMech [1] applies interval modeling to GPUs. While GPUMech is accurate for traditional GPU-compute workloads, we find that it falls short for memory-divergent applications.

We propose the *Memory Divergence Model (MDM)* which captures the key performance-related behavior of modern, memory-divergent GPU applications. We find that the poor spatial locality of memory-divergent applications leads to inefficient utilization of the Miss Status Holding Registers (MSHRs). The number of MSHRs determines the number of concurrent misses the cache can sustain without blocking (a blocked cache cannot accept any requests). Blocking has a profound performance impact. First, a blocked cache limits the ability of the GPU core to hide memory latencies with Memory and Thread-Level Parallelism (i.e., MLP and TLP). Second, the memory system becomes saturated as the cores issue a large number of requests to fetch all required data. MDM accounts for these effects by accurately modeling MSHR behavior and the Network-on-Chip (NoC) and DRAM queuing latencies. Overall, MDM improves performance prediction accuracy by 14× on average compared to the state-of-the-art GPUMech [1] approach across our memory-divergent applications.

2 UNDERSTANDING EMERGING GPU APPLICATIONS

2.1 The Architectural Effects of Memory Divergence

GPUs use multiple Streaming Multiprocessors (SMs) to execute code. Each SM can run a limited number of software threads concurrently. Thus, software threads are divided into groups, called warps, that match the width of the SM. An SM executes the instructions of all threads within a warp in lock-step. For load instructions, each thread issues a load for a single data element. These per-thread requests are aggregated to cache requests by the coalescer. On a cache hit, the cache line is read by the SM. On a miss, an MSHR is allocated and a memory request is sent to the lower levels of the memory hierarchy. If the access pattern is favorable (e.g., sequential), the coalescer can map the misses of the warp's concurrent threads to a single cache request, consuming a single MSHR. However, a significant fraction of emerging GPU applications are memory-divergent (i.e., the threads of a warp tend to access different cache blocks), exerting significant pressure on the limited number of MSHRs. If the cache runs out of MSHRs, it blocks until an MSHR becomes available. A blocked cache causes SM stalls because no load instructions can be executed.

To understand how the poor spatial locality of memory-divergent applications affects the memory system, Fig. 1 breaks down the average memory latency of GPU applications into the memory unit where it is incurred (see Section 4 for a description of our experimental setup). The key observation is that the benchmarks are clearly partitioned into two categories: The *Non-Memory Divergent (NMD)* benchmarks—where the latency due to insufficient MSHRs is negligible—and the *Memory Divergent (MD)* benchmarks—which on average spend hundreds and even thousands of cycles waiting for MSHRs to become available. Fig. 1 also shows that MD-applications tend to experience significant queuing latencies in the NoC and DRAM subsystems. Thus, an effective performance model for MD-applications needs to accurately model MSHR behavior, NoC queuing and DRAM queuing.

TABLE 1
Simulator Configuration

Parameter	Value
Clock frequency	1.4 GHz
Number of SMs	28
Number of mem. ctrl.	24
Warp schedulers per SM	4 (LRR)
Issue width per sched.	2 warp-instructions/cycle
L1 cache per SM	48 KB, 6-way, LRU, 128 MSHRs
L2 cache per mem. ctrl.	128 KB, 8-way, LRU, 128 MSHRs
NoC bandwidth	1050 GB/s
DRAM bandwidth	480 GB/s
Maximum warps per SM	64
Minimum L2 hit latency	120 cycles
Minimum DRAM latency	220 cycles

We now use the service latency predictions to predict the average queuing latency—and thereby the SM stall cycles caused by queuing latencies. The average queuing latency is determined by the average number of pending requests an arriving request must wait for times the average service latency. We first predict the average number of concurrent L1 misses M

$$M = \min(M^{\text{Read}} \times W, \#\text{MSHRs}) + M^{\text{Write}} \times W. \quad (4)$$

Read misses allocate MSHR entries and are therefore bounded by the number of L1 MSHRs. In other words, the application will either: (1) issue the number of read misses of the current interval of the representative warp times the number of warps; or, (2) as many read misses as there are MSHRs. Since the L1 caches in our GPU models are write-through and no-allocate, write misses effectively bypass the L1 and are independent of the number of MSHRs.

The number of queued requests is determined by application behavior while the service latency is an architectural parameter. Thus, we can use the same model to predict both NoC and DRAM stalls by providing $L^{\text{NoCService}}$ ($L^{\text{DRAMService}}$) as input to compute S^{NoC} (S^{DRAM})

$$S^{\text{NoC}} = \begin{cases} \#\text{SMs} \times M \times L^{\text{NoCService}}, & M^{\text{Read}} \times W > \#\text{MSHRs} \\ (1/2) \times \#\text{SMs} \times M \times L^{\text{NoCService}}, & \text{otherwise.} \end{cases} \quad (5)$$

The equation formalizes the key observations of Section 2. For MD-applications, the number of MSHRs is the bottleneck and the high degree of divergence keeps the memory system saturated. Since the memory system is saturated, each request needs to wait for all other requests. For NMD-applications, the memory requests are not sufficient to keep the memory queue saturated. In this case, the first request is serviced directly and the last request needs to wait for all other requests. Thus, a request waits for approximately half the concurrent requests.

3.2 MDM's MSHR Contention Model

The warps of MD-applications send their requests to the memory subsystem over consecutive batches (see Section 2.2). To estimate the length of these batches, we start by determining the memory latency in the absence of contention

$$L^{\text{NoContention}} = L^{\text{MinLLC}} + \text{LLCMissRate} \times L^{\text{MinDRAM}}. \quad (6)$$

Here, L^{MinLLC} is the round-trip latency of an LLC hit without NoC contention. The round-trip latency through the DRAM system is L^{MinDRAM} (again assuming no contention), but only LLC misses incur this latency. We then combine $L^{\text{NoContention}}$ with the average stall cycles due to queuing in the NoC and DRAM subsystems (obtained with Equation (5))

TABLE 2
Benchmarks

Benchmark	Suite	Abbr.	Type
Hotspot	Rodinia [8]	HS	NMD
B+trees	Rodinia	BT	NMD
Back Propagation	Rodinia	BP	NMD
FDTD3d	SDK [9]	FDTD	NMD
Srad	Rodinia	SRAD	NMD
Ray tracing	GPGPUsim [10]	RAY	NMD
2D Convolution	Polybench [11]	2DCONV	NMD
Stencil	Parboil [12]	ST	NMD
CFD solver	Rodinia	CFD	MD
Breadth-first search	Rodinia	BFS	MD
PageView Rank	MARS [13]	PVR	MD
RangeView Count	MARS	PVC	MD
Inverted Index	MARS	IIX	MD
Sparse matrix mult.	Parboil	SPMV	MD
Kmeans clustering	Rodinia	KMEANS	MD

$$S^{\text{Mem}} = L^{\text{NoContention}} + S^{\text{NoC}} + S^{\text{DRAM}}. \quad (7)$$

S^{Mem} is the predicted stall cycles due to L1 misses—considering both NoC and DRAM contention. We then use S^{Mem} to predict the SM stall cycles due to MSHR contention

$$S^{\text{MSHR}} = \begin{cases} (\lceil \frac{M^{\text{Read}} \times W}{\#\text{MSHRs}} \rceil - 1) \times S^{\text{Mem}}, & M^{\text{Read}} \times W > \#\text{MSHRs} \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

Equation (8) checks whether the number of requests of the current warps exceeds the number of MSHRs. If it does, we compute the number of batches needed to issue the memory requests of all warps by dividing the total number of read misses by the number of MSHRs. The latency of the final batch is covered by the queuing model, so we need to subtract one from this quantity to avoid adding this latency twice. Then, we multiply by S^{Mem} to obtain the combined SM stall cycles of these batches. NMD-applications are typically able to issue the requests of all warps in a single batch (see Fig. 1). Therefore, we set S^{MSHR} to zero for non-divergent intervals.

4 EXPERIMENTAL SETUP

We use GPGPU-sim 3.2 [10], a cycle-accurate GPU simulator, to evaluate MDM's prediction accuracy. We model an architecture similar to Nvidia's Pascal [14] as shown in Table 1. We select 15 applications: 8 NMD-applications and 7 MD-applications, from the main GPU benchmark suites. Table 2 provides details on the selected benchmarks. We simulate the benchmarks to completion with the (largest) default input set and report performance prediction error relative to simulated performance.

The original GPUMech proposal does not model NoC queuing delay and does not account for the DRAM and NoC queuing delays when estimating the MSHR stall latencies. GPUMech+ models a NoC queuing delay that resembles GPUMech's DRAM queuing model, whereby each request waits for half the total number of requests on average. GPUMech+ also accounts for the NoC and DRAM queuing delays when estimating the MSHR waiting time. MDM-Queue improves upon GPUMech+ by using MDM's NoC and DRAM queue model. MDM-MSHR improves upon GPUMech+ by using MDM's MSHR model. This enables us to independently evaluate MDM's queue model and MSHR model. MDM incorporates the improved NoC, DRAM and MSHR queuing delays.

5 RESULTS

Fig. 3 reports the relative IPC prediction error for our NMD and MD-benchmarks for all model combinations. MDM reduces

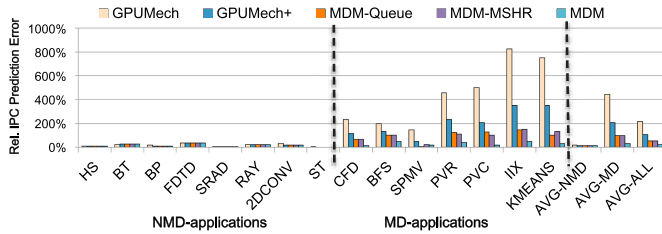


Fig. 3. IPC prediction error for our NMD and MD-benchmarks and the different performance models. The key take-away is that MDM significantly reduces prediction error for MD-applications while minorly reducing error for the NMD-applications.

prediction error by $14\times$ on average compared to GPUMech for the MD-benchmarks, from 444 to 32 percent. For the NMD-benchmarks, MDM reduces prediction error marginally compared to GPUMech, from 19 to 16 percent on average. Across all benchmarks, GPUMech has an average performance prediction error of 217 percent. MDM achieves an average prediction error of 23 percent. The execution times of the MDM and GPUMech models are practically equal.

GPUMech+, MDM-Queue and MDM-MSHR shed light on the relative importance of the different components of MDM for the MD-applications. Although GPUMech+ improves accuracy significantly compared to GPUMech, it still has a high average prediction error of 206 percent. This reinforces that minorly modifying GPUMech is insufficient and that MD-applications need a fundamentally new modeling approach. MDM-Queue improves upon GPUMech+ by applying the saturation model described in Section 3.1 to memory-divergent intervals, thereby reducing the average prediction error to 100 percent. Similarly, MDM-MSHR improves upon GPUMech+ by applying the batching model of Section 3.2 to memory-divergent intervals, which reduces the average prediction error to 98 percent. Neither MDM-Queue nor MDM-MSHR are able to accurately predict MD-application performance in isolation, indicating that modeling both queuing effects and MSHR behavior is critical to achieve low prediction error.

6 RELATED WORK

Prior work uses GPU modeling techniques to guide runtime optimizations (e.g., DVFS configuration [15] and cache miss-related optimizations [16]) or GPU resource scaling analysis [2]. Our work provides an accurate model for fast design space exploration. In general, prior performance modeling efforts make simplifications that lead to inaccuracies when modeling the cache hierarchy [4] and divergent applications [1], [3], or do not provide insight [2]. Volkov [17] studies GPU performance using simple synthetic benchmarks and shows that recent GPU models do not accurately capture the effects of memory bandwidth, non-coalesced accesses, and memory-intensive applications.

7 CONCLUSION

In this paper, we analyze the key performance characteristics of contemporary GPU applications and find that the poor spatial locality of these applications cause them to be memory-divergent. The modeling assumptions made by state-of-the-art GPU performance models such as GPUMech do not capture the characteristics of such applications. Applying GPUMech to memory-divergent applications leads to significant performance prediction errors (444 percent on average). We propose the Memory Divergence Model (MDM), which accurately models the batching and saturation behavior caused by high memory intensity and poor spatial locality. MDM significantly improves performance prediction accuracy compared to GPUMech, by $14\times$ on average across a set of memory-divergent applications.

ACKNOWLEDGMENTS

This work is supported through the European Research Council (ERC) Advanced Grant agreement No. 741097, Research Foundation Flanders (FWO) grants No. G.0434.16N and G.0144.17N, the National Natural Science Foundation of China through grants No. 61572508 and 61672526, NUDT Research Project No. ZK17-03-06. Lu Wang is supported through a CSC scholarship and UGent-BOF co-funding.

REFERENCES

- [1] J.-C. Huang, et al., "GPUMech: GPU performance modeling technique based on interval analysis," in *Proc. 47th Annu. IEEE/ACM Int. Symp. Microarchit.*, 2014, pp. 268–279.
- [2] G. Wu, et al., "GPGPU performance and power estimation using machine learning," in *Proc. IEEE 21st Int. Symp. High-Perform. Comput. Archit.*, 2015, pp. 564–576.
- [3] S. Hong and H. Kim, "An analytical model for a GPU architecture with memory-level and thread-level parallelism awareness," in *Proc. 36th Annu. Int. Symp. Comput. Archit.*, 2009, pp. 152–163.
- [4] Y. Zhang and J. D. Owens, "A quantitative performance analysis model for GPU architectures," in *Proc. IEEE 17th Int. Symp. High-Perform. Comput. Archit.*, 2011, pp. 382–393.
- [5] M. Burtscher, et al., "A quantitative study of irregular programs on GPUs," in *Proc. IEEE Int. Symp. Workload Characterization*, 2012, pp. 141–151.
- [6] N. Chatterjee, et al., "Managing DRAM latency divergence in irregular GPGPU applications," in *Proc. Int. Conf. High Perform. Comput. Netw. Storage Anal.*, 2014, pp. 128–139.
- [7] T. S. Karkhanis and J. E. Smith, "A first-order superscalar processor model," in *Proc. 31st Annu. Int. Symp. Comput. Archit.*, 2004, pp. 338–349.
- [8] S. Che, et al., "Rodinia: A benchmark suite for heterogeneous computing," in *Proc. IEEE Int. Symp. Workload Characterization*, 2009, pp. 44–54.
- [9] NVIDIA CUDA SDK Code Samples. NVIDIA Corp. 2013. [Online]. Available: <https://developer.nvidia.com/cuda-downloads>
- [10] A. Bakhoda, et al., "Analyzing CUDA workloads using a detailed GPU simulator," in *Proc. IEEE Int. Symp. Perform. Anal. Syst. Softw.*, 2009, pp. 163–174.
- [11] S. Grauer-Gray, et al., "Auto-tuning a high-level language targeted to GPU codes," in *Proc. Innovative Parallel Comput.*, 2012, pp. 1–10.
- [12] J. A. Stratton, et al., "Parboil: A revised benchmark suite for scientific and commercial throughput computing," University of Illinois, 2012. [Online]. Available: <http://impact.crhc.illinois.edu/parboil/parboil.aspx>
- [13] B. He, et al., "Mars: A MapReduce framework on graphics processors," in *Proc. Int. Conf. Parallel Archit. Compilation Techn.*, 2008, pp. 260–269.
- [14] NVIDIA GP100 Pascal Architecture. NVIDIA Corp. 2016. [Online]. Available: <https://www.nvidia.com/object/pascal-architecture-whitepaper.html>
- [15] R. Nath and D. Tullsen, "The CRISP performance model for dynamic voltage and frequency scaling in a GPGPU," in *Proc. 48th Annu. IEEE/ACM Int. Symp. Microarchit.*, 2015, pp. 281–293.
- [16] H. Dai, C. Li, H. Zhou, S. Gupta, C. Kartsaklis, and M. Mantor, "A model-driven approach to warp/thread-block level GPU cache bypassing," in *Proc. 53rd ACM/EDAC/IEEE Des. Autom. Conf.*, 2016, pp. 1–6.
- [17] V. Volkov, "Understanding latency hiding on GPUs," PhD dissertation, Dept. Comput. Data Sci., Univ. California, Berkeley, Berkeley, CA, USA, 2016.