

Hierarchically Characterizing CUDA Program Behavior

Zhibin Yu, Hai Jin
Service Computing Technologies and System Lab/
Cluster and Grid Computing Lab,
Huazhong University of Science and Technology
Wuhan, China, 430074

Nilanjan Goswami, Tao Li
Intelligent Design of Efficient Architecture Lab,
University of Florida, Gainesville
Florida, USA

Lizy Kurian John
Laboratory for Computer Architecture,
Department of Electrical and Computer Engineering
University of Texas at Austin
Austin, TX78712 USA

Introduction

Over the last few years, the performance of Graphic Processing Unit (GPUs) has improved more rapidly than that of CPUs [1]. The key to harness the powerful computation power of GPGPUs is an easy and efficient programming model. To this end, NVIDIA created Compute Unified Device Architecture (CUDA) programming mode [2]. It is implemented by extending the standard ANSI C with keywords that designate data-parallel functions called kernels.

CUDA programming mode is very different from sequential programming modes. To characterize CUDA program behavior and understand why and where they can achieve significant speedup comparing to sequential programs, it is important to revisit the basic block level and instruction level properties besides those at the thread level. In this paper, we propose to characterize CUDA program behaviors hierarchically by quantitatively gleaning properties from thread, basic block, and instruction levels.

In addition, previous researchers have demonstrated that basic blocks vectors (BBVs) are one of the most accurate techniques for creating code signatures [3] for sequential programs. In this paper, we firstly employ basic block and basic block vectors to analyze the code signature of CUDA threads. We observed that basic block characteristics of CUDA kernels are very different from those of sequential programs. Based on the basic block vectors, we construct the similarity matrix of threads. We show that the similarity matrix can be a very powerful tool for performance tuning.

Methodology

Metrics

- ❖ Number of instructions per thread
- ❖ Thread performance
- ❖ Number of basic blocks
- ❖ Average basic block size
- ❖ Program footprint
- ❖ Instruction mix
- ❖ Instruction-level parallelism
-

Similarity Matrix

- ❖ Basic block vector per thread
- ❖ Basic block vector per kernel
- ❖ Synchronization vector
- ❖ Similarity matrix based on basic block vectors

Benchmarks

- ❖ CUDA SDK
- ❖ Parboil
- ❖ Rodinia
- ❖ Other programs from recent papers

35 benchmarks in total

Platforms

- ❖ Based on GPGPUsim
- ❖ Extends cuda-sim to support
 - Measure instruction dependency distance
 - Generate basic block vectors per thread and for the whole kernel
 - Generate synchronization vectors
 - Measure instruction mix
 - Measure the instruction count per thread
- ❖ Extends gpu-sim to support
 - Measure the performance of each CUDA thread

Table 1 Hardware Configuration

Number of Shader Cores	28
Warp size	32
SIMD Pipeline Width	8
Number of Threads/Core	1024
Number of CTAs/Core	8
Number of Registers/Core	16384
Shared Memory /Core (KB)	16(16 banks, 1 access/cycle/bank)
Constant Cache Size / Core	8KB (2-way set assoc, 64 lines)
Texture Cache Size / Core	64KB (2-way set assoc, 64 lines)
Number of Memory Channels	8
L1 Cache	None
L2 Cache	None
Bandwidth Per Memory Module	8 (Bytes/Cycle)
DRAM Request Queue Capacity	32
Memory Controller	Out of order (FR-FCFS)
Branch Divergence Method	Immediate Post Dominator
Warp Schedule Policy	Round Robin among read warps

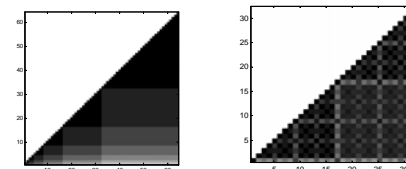
Table 2 Interconnect Configuration

Topology	Mesh
Routing Mechanism	Dimension Order
Routing delay	1
Virtual channels	2
Virtual channel buffers	4
Virtual channel allocator	ISLIP / PIM
Allocation tiers	1
Virtual channel allocation delay	1
Input speedup	2
Flit size (Bytes)	16

Results

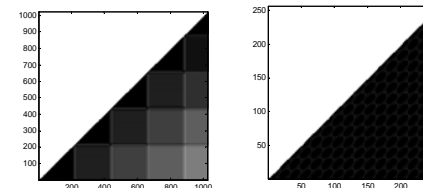
Table 3 Basic Block Properties of CUDA Programs
NBB for % means Number of Basic Blocks account for % of program execution
*: The BS is a modified version of BlockSched

Benchmark	Number of Basic Blocks	NBB for 90%	NBB for 95%	Average Basic Block Size	Average Number of Successive Basic Blocks
CGT_13	11	5	4	7.56	1.628
BS	11	5	4	7.56	1.628
BlockSched	21	4	4	26.4	1.5
BP_k1	10	4	4	9.61	1.5
BSV	4	4	4	26.4	1.5
CL	110	4	4	5.10	1.33
CP	2	5	3	10.43	1.33
CS_1	2	4	4	33.75	1.39
FW_16	10	10	10	12.75	1.41
HS	26	9	11	9.10	1.45
KM_12	16	7	10	4.76	1.44
LTB_1	40	11	13	8.72	1.6
LS	30	12	14	7.61	1.57
LT_13	11	4	4	7.5	1.55
LS*	21	5	5	8	1.428
MC	4	2	2	10.3	1.5
MM	4	2	2	18.57	1.5
MHP_11	15	2	2	12	1.56
MT_11	2	2	2	9.5	1.3
NS	21	12	14	5.98	1.33
NS*	6	4	4	15.14	1.33
NS_11	27	6	10	15.53	1.41
NS_11*	16	4	8	8.51	1.44
PP	103	59	100	8.47	1.49
PR_1	10	7	8	7.16	1.4
RAY	29	10	10	10	1.66
RFS_11	29	6	6	14.43	1.48
SD	19	12	13	12.48	1.42
SEA_11	11	6	9	6.79	1.54
SP	11	8	11	5.68	1.56
SRAD_11	15	12	12	4.52	1.41
SS_12	21	9	14	7.43	1.48
STB	21	8	11	23.78	1.47
STD	19	12	13	12.48	1.42
TPAE*	61	15	21	8	1.55



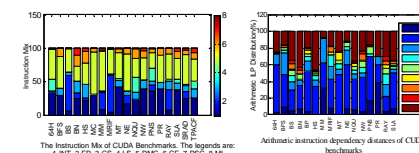
64H_k2,1-64

BP_k1,1-32



LT_k3,1-1024

NW_k1,1-256



The instruction mix of CUDA benchmarks. The legends are: 1-INT, 2-FF, 3-CS, 4-LS, 5-DMC, 6-CF, 7-PSC, 8-MI

Arithmetic instruction dependency distances of CUDA benchmarks

INT --- Integer arithmetic
FP --- Floating point
CS --- Comparison and selection
LS --- Logic and shift
DMC --- Data movement and conversion
CF --- Control flow
PSC --- Parallel Synchronization and Communication
MI --- Miscellaneous

Conclusion

We present a hierarchical methodology to quantitatively characterize CUDA program behavior at thread, basic block and instruction level. We summarize the main findings here. First, the IPC of CUDA thread is only about 1/40~1/100 of the average IPC of CPUs. Second, the average number of basic blocks of CUDA programs is 1/11~1/25 of that of sequential programs. Finally, the data movement and conversion instructions (*mov*, *cvt*) of CUDA programs hold a high percentage (37.8%). There are also a lot of other findings such as ILP of CUDA kernels in the paper. To our best knowledge, we are the first to do such characterization for CUDA programs. The outcome of our work can be used to optimize GPGPU architectures and CUDA compilers.

The CUDA programming model derives from the more general Single-Program Multiple-Data (SPMD) model which is widely available other parallel processing systems. Therefore, the proposed hierarchical characterization methodology, especially the basic block vectors and similarity matrix, can also be used to characterize other SPMD parallel programs.

Acknowledgements

This work is supported by NSF China under Grant No. 60973036

References

- [1] <http://www.nvidia.com/>
- [2] NVIDIA CORPORATION. NVIDIA CUDA Programming Guide, version 3.0.
- [3] T. Sherwood, E. Perelman, G. Hamerly, and B. Calder, "Automatically Characterizing Large Scale Program Behavior", Proceedings of the 10th International Conference on Architectural Support for Programming Languages and Operating Systems, ACM Press, October 5-9, 2002, San, Jose, CA, pp. 45-57