

MIPP: A Microbenchmark Suite for Performance, Power, and Energy Consumption Characterization of GPU architectures

Nicola Bombieri
Dept. Computer Science
University of Verona
nicola.bombieri@univr.it

Federico Busato
Dept. Computer Science
University of Verona
federico.busato@univr.it

Franco Fummi
Dept. Computer Science
University of Verona
franco.fummi@univr.it

Michele Scala
Dept. Computer Science
University of Verona
michele.scala@univr.it

Abstract—GPU-accelerated applications are becoming increasingly common in high-performance computing as well as in low-power heterogeneous embedded systems. Nevertheless, GPU programming is a challenging task, especially if a GPU application has to be tuned to fully take advantage of the GPU architectural configuration. Even more challenging is the application tuning by considering power and energy consumption, which have emerged as first-order design constraints in addition to performance. Solving bottlenecks of a GPU application such as high thread divergence or poor memory coalescing have a different impact on the overall performance, power and energy consumption. Such an impact also depends on the GPU device on which the application is run. This paper presents a suite of microbenchmarks, which are specialized chunks of GPU code that exercise specific device components (e.g., arithmetic instruction units, shared memory, cache, DRAM, etc.) and that provide the actual characteristics of such components in terms of throughput, power, and energy consumption. The suite aims at enriching standard profiler information and guiding the GPU application tuning on a specific GPU architecture by considering all three design constraints (i.e., power, performance, energy consumption). The paper presents the results obtained by applying the proposed suite to characterize two different GPU devices and to understand how application tuning may impact differently on them.

I. INTRODUCTION

With the growth of computational power and programmability, Graphic Processing Units (GPUs) have become increasingly used as general-purpose accelerators. They not only provide high peak performance, but also excellent energy efficiency [10]. As a consequence, besides supercomputers, GPUs are quickly spreading in low-power and mobile devices like smartphones. NVIDIA Tegra X1 [1] and Qualcomm Snapdragon [2] are some among the several system-on-chip examples available in the mobile market that integrate GPUs with other processing units (i.e., CPUs, FPGAs, DSPs).

On the other hand, the large number of operating hardware resources (e.g., cores and register files) employed in GPUs to support the massive parallelism leads to a significant power consumption. The elevated levels of power consumption have a sensible impact on such many-core device reliability, aging, economic feasibility, performance scaling and deployment into a wide range of application domains. Different techniques have been proposed to manage such high levels of power dissipation and to continue scaling performance and energy. They include approaches based on dynamic voltage/frequency scaling (DVFS) [7], CPU-GPU work division [8], architecture-level/runtime adaptations [15], dynamic resource allocation [6], and application-specific (i.e., programming-level) optimizations [18]. Particularly in this last category, it has been observed that source-code-level transformations and application specific optimizations can significantly affect the GPU resource utilization, performance, and energy efficiency [14].

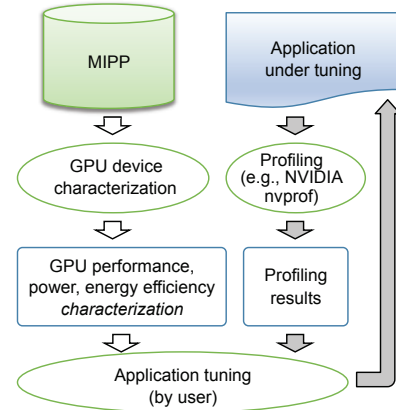


FIG. 1: Overview of application tuning through MIPP.

In this context, even though profiling tools (e.g., CUDA nvprof, AMD APP) exist to help programmers in the application analysis and optimization, they do not provide a complete view of the GPU features (especially on power consumption and energy efficiency) neither they provide a correlation among these design constraints.

This paper presents *MIPP*, a suite of microbenchmarks that aims at characterizing a GPU device in terms of performance, power, and energy consumption. In particular, it aims at understanding how application bottlenecks involving selected functional components or underutilization of them can affect the code performance, power consumption, and energy efficiency on the given device. The functional components include arithmetic instruction units, memories (shared, cache, DRAM, constant), scheduling and synchronization components. Fig. 1 shows how the suite can be applied during an application tuning. First, the microbenchmarks are run on the GPU device to *characterize* the device in terms of performance, power and energy over its main functional components. Then, the application under tuning is profiled by using a standard profiling tool. The profiling results provide information on bottlenecks and underutilization of functional components. The microbenchmark results extend such an information by quantitatively showing how such bottlenecks and underutilization affect performance, power and energy. The proposed model allows the flow (underlined by grey arrows in Fig. 1) to be iterated for incremental tuning of the application.

The suite has been applied to characterize two different GPU devices (i.e., an NVIDIA Kepler GTX660 and a low power embedded system NVIDIA Jetson TK1). The results show how the same code optimizations have a different impact on the design constraints on the two GPU architectures.

The paper is organized as follows. Section II presents

the analysis of the related work. Section III presents the microbenchmark suite. Section IV reports the experimental results, while Section V is devoted to concluding remarks.

II. RELATED WORK

Different papers have been presented to evaluate functional and architectural characteristics of GPUs through microbenchmarking. The work in [16] proposes a microbenchmark suite to measure the CUDA-visible architectural characteristics of the Nvidia GT200 (GTX280) GPU. Such a measure includes various undisclosed characteristics of the processing elements and the memory hierarchies. The analysis exposes undocumented features that impact program performance and correctness.

In [11], the authors evaluated throughput and the power efficiency of three 128-bit block ciphers (AES, Camellia, and SC2000) on Nvidia Geforce GTX 680 with Kepler architecture and on AMD Radeon HD 7970 with GCN architecture. For the comparison, the authors used Nvidia Geforce GTX 580 and AMD Radeon HD 6770 with architecture of one generation earlier. The authors developed a microbenchmark suite that allows understanding that arithmetic logical instructions are required by encryption processing but are eliminated from some of the processing cores in NVIDIA Kepler architecture, unlike AMD graphics core next (GCN) architectures.

In [17], the authors propose an OpenCL microbenchmark suite for GPUs and CPUs. They present the performance results of hardware and software features such as bus bandwidth, memory architectures, branch architectures and thread hierarchy, etc., evaluated through the proposed microbenchmarks on multi-core X86 CPU and NVIDIA GPUs.

In [4], the authors propose a microbenchmarking methodology based on short elapsed-time events (SETEs) to obtain comprehensive memory microarchitectural details in multi- and many-core processors. This approach requires detailed analysis of potential interfering factors that could affect the intended behavior of such memory systems. They lay out guidelines to control and mitigate those interfering factors.

In [9], the authors propose a fine-grained benchmarking approach and apply it on two popular GPUs (i.e., Fermi and Kepler), to expose the previously unknown characteristics of their memory hierarchies. The authors investigate the structures of different cache systems, such as data cache, texture cache, and the translation lookaside buffer (TLB). They also investigate the impact of bank conflict on shared memory access latency, the GPU memory hierarchy, which can help in the software optimization and the modelling of GPU architectures.

Each of these contributions either presents microbenchmarks for characterizing a GPU device from a specific design constraint point of view (performance or power) or presents and analysis of GPU performance and power of a specific application. In contrast, we propose a comprehensive microbenchmark suite that characterizes all the functional and architectural units of a GPU architecture from both the performance and power/energy consumption point of view. Such a characterization allows understanding how application bottlenecks involving selected functional components or under-utilization of them can affect the code quality by considering all the three design constraints and the given GPU device.

III. THE MICROBENCHMARK SUITE

We developed a suite of microbenchmarks to selectively study the behaviour of a wide range of GPU functional components. Fig. 2 gives an overview of such a microbenchmark suite by reporting, for each microbenchmark, the exercised GPU component, the involved specific instructions, and the considered features.

ARITHMETIC PROCESSING	MEMORY
32-bit Integer (ALU) <ul style="list-style-type: none"> Simple (add, subtract, multiply-add or mad) Complex (multiply) Bit operations (clz, msb, bit-reverse, bit-insert) Shift Compare (<, >, ≤, ≥, min, max) Bit-Counting (population count) 	DRAM Memory <ul style="list-style-type: none"> Throughput Coalescence* Access size*
64-bit Integer (ALU) <ul style="list-style-type: none"> Simple (add, subtract) 	Shared Memory <ul style="list-style-type: none"> Throughput Conflict*
32-bit Floating-point (FPU) <ul style="list-style-type: none"> Simple (add, subtract, multiply, fused multiply-add or fma) Complex (division, division FTZ) 	L2 Cache <ul style="list-style-type: none"> Throughput
32-bit FP Special functions (SFU) <ul style="list-style-type: none"> Transcendental functions (sin, cos, exp, rsqrt, reciprocal, log) 	Constant Memory <ul style="list-style-type: none"> Throughput
64-bit Floating-point (DFU) <ul style="list-style-type: none"> Simple (add, subtract) 	

FIG. 2: Microbenchmark Classes

A microbenchmark consists of a GPU kernel code that exercises a specific functional component of the architecture and whose instructions can be evaluated at a clock-cycle accuracy. The generic structure of the microbenchmark main procedure consists of a long sequence of one or more selected instructions (e.g., arithmetic instructions, memory accesses) that executes without any interference deriving from other instructions. The microbenchmarks have been implemented to stress only a specific functional component at a time, while affecting the others as little as possible to obtain reliable and accurate feedback.

The microbenchmark code is written by combining the CUDA C/C++ language with inline intermediate assembly to avoid compiler side effects that may elude the target properties. The parallel thread execution (PTX) is a GPU machine-independent language that allows expressing general purpose computation through virtual ISA. We exploited the PTX language to force a specific operation on a data type, to avoid compiler optimizations, to prevent caching/local-storage mechanisms, and to drive the memory accesses. We adopted several arrangements to preserve the code functionality. As an example, registers are initialized with dynamic values to avoid constant propagation in arithmetic benchmarks.

Finally, to perform an extensive computation, we applied template meta-programming and nested loops. Such programming techniques are required to prevent compiler optimizations (e.g., loop collapsing and dead code elimination). The same full control on the compiler cannot be obtained by simply handling the compilation flags, since they apply only to arithmetic instruction optimizations (such as floating-point precision and floating-point multiply-add enabling). The code controls the intensity variability, the amount of computation, and other aspects through parameterized procedures.

Each microbenchmark run returns information like *execution time*, *actual throughput* (to compare with the theoretical throughput from the device specifications), *average* and *max power consumption*, *energy consumption* and *energy efficiency*. Some microbenchmarks (marked with "*" in Fig. 2) are also applied to exercise functional components with different intensity. As an example, a microbenchmark allows analysing the shared memory throughput by generating a different amount of bank conflicts, from zero to the maximum value, and by measuring the corresponding access time.

The microbenchmarks provide a quantitative model of the target GPU architecture based on performance and power, and provide important guidelines for the application optimization. As shown in Fig. 2, the microbenchmarks are grouped into two classes: *Arithmetic processing* and *memory hierarchy*.

```

__global__ ADD_THROUGHPUT()
1: int R1 = clock(); //assign dynamic values to R1,R2 to
2: int R2 = clock(); //avoid constant propagation
3: int startTimer = clock();
4: Computation<N>(R1, R2); //call the function N times
5: int endTimer = clock();

template<int N>() //template metaprogramming
__device__ __forceinline__ COMPUTATION(int R1, int R2)
1: #pragma unroll 4096 //maximum allowed unrolling
2: for (int i = 0; i < 4096; i++) do
3:     asm volatile("add.s32 : "=r"(R1) : "r"(R1), "r"(R2));
4: end //volatile: prevent ptx compiler optimization
5: Computation<N-1>(R1, R2); //recursive call

```

FIG. 3: Example of microbenchmark code. The code aims at measuring the maximum instruction throughput of the add operation.

A. Arithmetic processing benchmarks

This class of microbenchmarks targets the complete set of arithmetic instructions natively supported by the GPU, by distinguishing between integer and floating point over 32 and 64-bit word sizes.

The benchmarks perform a long sequence of instructions to stress the ALU components. All PTX instructions of arithmetic benchmarks have a direct translation into the native ISA, called SASS (Shared ASSEMBLY), except 64-bit integer operations and floating-point divisions that are compiled into multiple instructions. A SM executes native instructions in one clock cycle, providing a throughput (instructions per clock cycles) limited by the concurrency of the exercised ALU component. Depending on the compute capability of the device and on the architecture, the implementation of non-native instructions may correspond to a different number and type of ISA instructions. Arithmetic benchmarks include also four different types of division operations classified by approximation (IEEE754 Compliance and fast hardware approximation) and normalization (normal and de-normal numbers).

As an example, Fig. 3 summarizes the microbenchmark developed to analyse the 32-bit integer arithmetic processing unit (simple add). The code implements dynamic value assignments to registers (see rows 1 and 2 in the upper side of the figure) to avoid the *constant propagation* optimization by the compiler¹. The code also adopts *recursive* and *template-based* metaprogramming. This allows generating an arbitrarily long sequence of arithmetic instructions without any control flow instructions. ($4096 \times N$ add instructions in the example²).

B. Memory benchmarks

This class of benchmarks focuses on the impact of throughput and access patterns on DRAM, shared, constant, and L2 cache memories. The *DRAM throughput* benchmark executes several global accesses to different memory locations with a stride of 128 bytes between grid threads to avoid L1 coalescing. The *L2 benchmark* repeats a compile-time sequence of store instructions on the same memory address. We use cache modifiers [12] to avoid L1 cache hits that can occur in the store operations. *Shared* and *constant memory* benchmarks consist of a sequence of load/store instructions respectively. In the *coalescing* benchmark, we vary the number of threads within a warp that access to continuous locations.

¹Static value assignments to registers are generally solved and substituted by the compiler optimizations through inlining operations.

²In the example, 4,096 unrolls are a good compromise between loop body replication and template recursion. Over a fixed number of loop unrolling iterations the compiler would insert control statements in the loop to reduce the size of the binary code.

```

__device__ clock_t devClocks[RESIDENT_WARPS];
__device__ int devTMP;

template<int CONFLICTS>
__global__ SHAREDMEMCONFLICTS()
1: __shared__ volatile int SMem[1024];
2: volatile int* Offset = SMem + LANE_ID * (CONFLICTS+1);
3: clock_t startTimer = clock64();
4: Computation<N>(Offset); //call the function N times
5: clock_t endTimer = clock64();
6: if (LANE_ID == 0) then
7:     devClocks[WARP_ID] = endTimer - startTimer;
8: if (THREAD_ID == 1024) then
9:     devTMP = SMem[0]; //never executed

template<int N> //template metaprogramming
__device__ __forceinline__ COMPUTATION(volatile int* Offset)
1: #pragma unroll 4096
2: for (int i = 0; i < 4096; i++) do
3:     asm volatile("st.volatiles32 [%0], %1;" :
4:         "l"(Offset), "r"(i) : "memory" );
5:         //asm volatile: prevent PTX compiler optimization
6: end
7: Computation<N-1>(Offset); //recursive call

```

FIG. 4: Example of the microbenchmark code to measure the impact of shared memory bank conflicts.

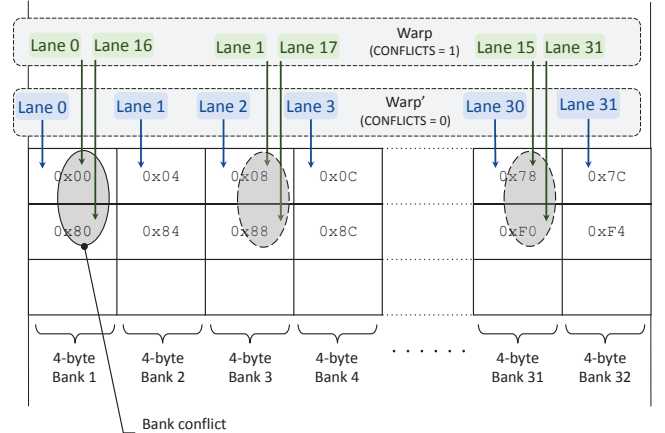


FIG. 5: Example of memory accesses with no and one bank conflict.

For example, to test the impact of the worst memory access pattern (no coalescence) we apply the same stride of the DRAM throughput benchmark, to evaluate 1/16 of coalescence we divide the warp threads into 16 groups of two threads where each group accesses in different addresses. The *access size* benchmark copies one large array into another multiple times and, in each execution, we vary the data type size.

Fig. 4 summarizes the microbenchmark developed to measure the impact of shared memory bank conflicts on the memory access throughput. The procedure first computes, for each thread of a warp, the address offsets of shared memory that can lead to bank conflicts (line 2 in the upper side of Fig. 4, where `LANE_ID` represents the thread id in the warp, and `CONFLICTS` represents the number of conflicts to generate). Fig. 5 shows, for example, the offsets generated to lead to no and to one bank conflict. Then, it performs a long sequence of store operations with no interrupt or intermediate operation. The code implements *volatile* quantifiers (lines 1,2 in the upper side of Fig. 4) to avoid *local-storage* optimization by the compiler. As for the simple add example, the code adopts

recursive and *template-based* metaprogramming to generate an arbitrarily long sequence of arithmetic instructions ($4,096 \times N$ store instructions in the example). In the last step (Line 7 in the upper side of the figure) each warp sends the timing results to the host through global memory. Line 8, 9 ensure that the result is stored in global variable with a fake write instruction (that is never executed since the higher thread id in a block is 1023) to prevent dead code elimination by the compiler.

Similarly, other microbenchmarks of this class exercise their corresponding functional components with different intensity. This allows characterizing the components behaviour under different workloads. For example, the global memory throughput is affected by access pattern that involves a different number of memory transactions or by saturation of the instruction pipeline. Varying the intensity of a microbenchmark allows us to better understand the main factors that affect the performance and the power consumption in real-world applications, where functional components show a wide range of utilization values.

In general, the microbenchmarks have been developed to guarantee enough computation time (i.e., at least of some milliseconds) to overcome the limitation of the sampling frequency in the measurement of the power features (P_{\max} , P_{total} , P_{avg}) and to minimize the RLC effect in the kernel starting/ending phases.

IV. EXPERIMENTAL RESULTS

We run the suite to characterize two different GPU devices. The first is an NVIDIA Kepler GeForce GTX 660 with CUDA Toolkit 7.5, AMD Phenom II X6 1055T (3GHz) host processor, and Ubuntu 14.04 OS. The second is a Tegra K1 SoC (Kepler architecture) on an NVIDIA Jetson TK1 embedded system, with CUDA Toolkit 6.5, 4-Plus-1 NVIDIA host multi processor (four ARM Cortex-A15 cores on the general cluster and one ARM Cortex-A7 in the low power cluster), and Ubuntu 14.04 OS.

Performance information has been collected through the `clock64` device instruction for all microbenchmarks to provide high accuracy. The microbenchmark kernels are organized in such a way that each warp stores the starting and the ending clock counters for each execution and sends the difference to the host, which computes the average among all device warps.

Power and energy consumption information has been collected through the *Powermon2* power monitoring device [3]. The device allows measuring the voltage and the current values from different sources at the same time with a frequency of 1024 Hz for every sensors. The GTX 660 requires five 12V pins, three for the pci-express power connectors and two for auxiliary connectors. We used a pci-express *interposer* to isolate the GPU power connectors from the motherboard. The Jetson TK1 requires only a DC barrel connector adapter to enable the power monitoring. We designed specific API and procedures to allow microbenchmarks to communicate and to synchronize with the Powermon device. In particular, we ensured that the kernel calls are synchronized with the first power measurement. Each microbenchmark is repeated ten times with two seconds of idle activity between consecutive executions. The analysis has been performed with the default GPU frequency setting on both devices. To measure, as much accurately as possible, the power consumption of the Jetson TK1 SoC, in which is not possible to physically isolate the GPU (Tegra K1) from the rest of the system, we operated as follows. We implemented the communication with the platform remotely without any additional connected peripherals, we disabled the Linux display manager (lightdm) and the HDMI port from the OS, we physically disabled the general CPU

	DRAM	L2	SHARED	CONSTANT
Execution Time (ms)	1,170	1,013	220.2	60.8
Real Throughput (OPs per Cycle)	1.0	1.1	5.3	19.6
Avg. power (W)	92.1	71.5	59.2	63.3
Max power (W)	102.0	80.0	62.0	68.0
Energy (J)	107.8	72.5	13.0	3.8
Energy efficiency (10^6 Transactions/Watt)	10.0	14.8	82.4	279.1
nano Joule per mem. transaction	100.4	67.5	12.1	3.6

TABLE III: GTX 660 - Characteristics of accesses on DRAM, L2, shared and constant memories.

cluster, and we forced the low power CPU cluster to run with the lowest frequency.

Tables I and II report the results obtained by running the *Arithmetic Processing* benchmarks on the GTX 660 and TK1, respectively. The benchmarks, which are organized over columns, consist each one of 10^9 instructions per SM (i.e., consider that the GTX 660 consists of 5 SMs, while the TK1 consists of 1 SM). For each benchmark, the tables report the execution time, the theoretical peak throughput of the corresponding functional unit provided in the device specifications [13] (*Spec Throughput*) and that measured through the proposed benchmark (*Real throughput*). The device specifications do not include the theoretical peak throughput of the Integer 64-bit simple unit since such an operation has not an embedded hardware implementation (it is performed by combining different hardware units). The tables also show the average power, the peak power, the energy consumption for a single SM, the energy efficiency [5], i.e. performance per watt (Million Instructions Per Second per Watt - MIPS/Watt) and the energy consumption for a single instruction in nanojoules.

For both the GTX 660 and TK1, the benchmarks underline that the theoretical peak throughput of several functional units (e.g., complex multiply, population count, shift Integer 32-bit etc.) can be actually reached. The measured peak values often exceed the theoretical values provided by the device specifications. We assume this is due to the fact that the theoretical values refer to the declared compute capability of the GPU rather than actual GPU manufacturing of the vendor. In any case, the difference between the two value is negligible. In contrast, the peak throughput of selected functional units, such as the simple 32-bit either Integer or Floating Point add cannot be actually fully exploited. We assume this is due to the actual latency of the fetching subunits, which do not support the throughput of the computation subunits. Even though the two devices are fairly different (desktop-oriented GTX 660, and low-power embedded system TK1), the benchmarks underline they rely on equivalent Kepler SMs, whose peak performance are fully comparable.

Finally, Tables I and II report information about power and energy consumption, which are not provided with the device specifications. Power and energy characteristics refer to the whole GPU device and underline the structural characteristics of the two GPU device architectures (5 SMs vs. 1 SM). The tables also show that FP 64-bit operations have a significant impact on the energy consumption of both the devices (15 times higher than Integer/FP 32-bit simple and Integer compare instructions).

Tables III and IV report the results obtained by running the microbenchmarks of the *memory* class to evaluate the throughput of the different device memories. The results allow

	INTEGER 32-BIT						INTEGER 64-BIT	FP 32-BIT		FP 64-BIT
	SIMPLE	COMPLEX	POP. COUNT	SHIFT	BIT OP.	COMPARE	SIMPLE	SIMPLE	SPECIAL	SIMPLE
Execution Time (ms)	9.6	34.3	34.3	34.3	36.7	9.8	24.5	9.7	36.6	137.0
Spec Throughput (OPs \times Cycle) \times SM	160	32	32	32	32	160	n.a.	192	32	8
Real Throughput (OPs \times Cycle) \times SM	126.1	34.7	34.7	34.7	34.3	122.2	51.2	126.0	33.9	8.8
Avg. power (W)	54.4	54.4	53.7	52.5	56.8	54.3	57.6	55.7	60.6	53.3
Max Power (W)	62	59	56	56	60	59	62	62	65	59
Energy (J)	0.1	0.4	0.4	0.4	0.4	0.1	0.3	0.1	0.4	1.4
Energy efficiency (MIPS \times Watt) \times SM	2,057.1	576.3	583.8	597.0	514.6	2,020.0	760.2	1,990.3	483.9	146.9
nanojoule per instruction	0.1	0.3	0.3	0.3	0.4	0.1	0.2	0.1	0.4	1.3

TABLE I: GTX 660 - Characterization with Arithmetic Processing benchmarks.

	INTEGER 32-BIT						INTEGER 64-BIT	FP 32-BIT		FP 64-BIT
	SIMPLE	COMPLEX	POP. COUNT	SHIFT	BIT OP.	COMPARE	SIMPLE	SIMPLE	SPECIAL	SIMPLE
Execution Time (ms)	140.8	483.0	489.9	486.2	515.5	145.1	399.0	143.3	508.5	1,878.8
Spec Throughput (OPs \times Cycle) \times SM	160	32	32	32	32	160	n.a.	192	32	8
Real Throughput (OPs \times Cycle) \times SM	123.0	32.1	32.1	32.1	32.1	121.0	40.6	123.1	30.7	8.0
Avg. Power (W)	3.2	3.2	3.2	3.2	3.2	3.2	3.2	3.2	3.3	3.2
Max Power (W)	4.0	4.0	4.0	4.0	4.0	4.0	4.0	4.0	4.0	4.0
Energy (J)	0.5	1.5	1.6	1.5	1.6	0.5	1.3	0.5	1.7	6.0
Energy efficiency (MIPS \times Watt) \times SM	2,376.2	703.1	688.3	699.1	658.7	2,318.6	844.5	2,355.2	646.2	180.2
nanojoule per instruction	0.4	1.4	1.5	1.4	1.5	0.4	1.2	0.4	1.5	5.5

TABLE II: Jetson TK1 - Characterization with Arithmetic Processing benchmarks.

	DRAM	L2	SHARED	CONSTANT
Execution Time (ms)	24,703	22,339	14,915	3,905
Real Throughput (OPs per Cycle)	0.6	0.7	1.0	3.8
Avg. power (W)	4.1	3.8	3.3	3.3
Max power (W)	5.0	5.0	5.0	4.0
Energy (J)	100.4	85.4	49.2	12.9
Energy efficiency (10^6 Transactions/Watt)	10.7	12.6	21.8	83.3
nano Joule per mem. transaction	93.5	79.5	45.8	12.0

TABLE IV: Jetson TK1 - Characteristics of accesses on DRAM, L2, shared and constant memories.

understanding how the throughput differs among memories and how it differs between the two devices. As an example, an application running on the TK1 accesses the constant memory 4 times faster than in DRAM. The same application running in the GTX 660 accesses the constant memory 20 times faster than in DRAM. Moreover, the table shows that DRAM accesses strongly affect the average and the max power of GTX 660 and TK1 devices while, the on-chip memories (shared and constant memories) have slightly higher average power than arithmetic instructions.

Figure 6(a) shows the impact of thread *coalescence* in DRAM memory accesses on the GTX 660 performance, power

and energy. The figure shows the effect starting from no coalescence (one memory transaction per warp thread access), 1/16 coalescence (one transaction per two warp threads), until FULL coalescence (one transaction per a whole 32-threads warp). The figure shows how performance and energy are proportional to the reached coalescence. In contrast, max and average power reach the highest values at 1/8 coalescence, and they decrease until FULL coalescence. This is due to the fact that from NO to 1/8, the coalescence is incrementally supported by the 32-Byte L2 memory banks, which saturate at 1/4 coalescence (i.e., each set of 4 transactions per warp, each one 32-Byte large, saturate a 32-Byte L2 bank). From 1/4 on, the coalescence relies on the 128-Byte L1 memory banks. Figure 6(b) reports the same analysis on the TK1, for which the decreasing of the max power can be observed at the FULL coalescence state only.

Figures 7(a) and 7(b) quantify the impact of bank conflicts in shared memory on power, performance, and energy consumption. They underline that the bank conflicts similarly impact on performance and energy on the two devices. In contrast the analysis underlines that up to 7 conflicts do not affect the max power on the GTX 660, while up to 7 conflicts strongly affect the max power on the TK1.

Overall, the results obtained by running the proposed suite on a given GPU device allows understanding the specific impact of a tuning step on the design constraints (performance, power, and energy consumption). Improving the code performance that affect a specific functional component (e.g., coalescence of memory accesses) may violate a design constraint

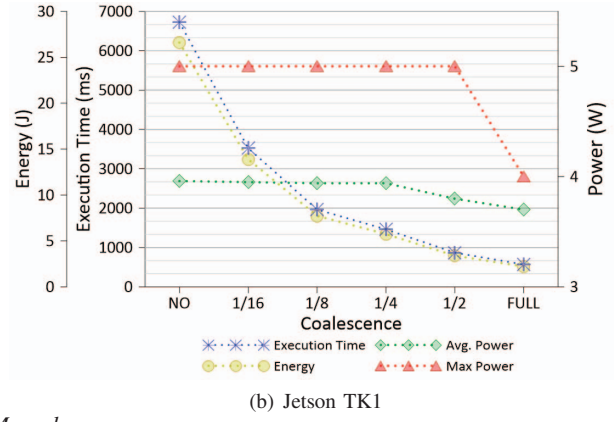
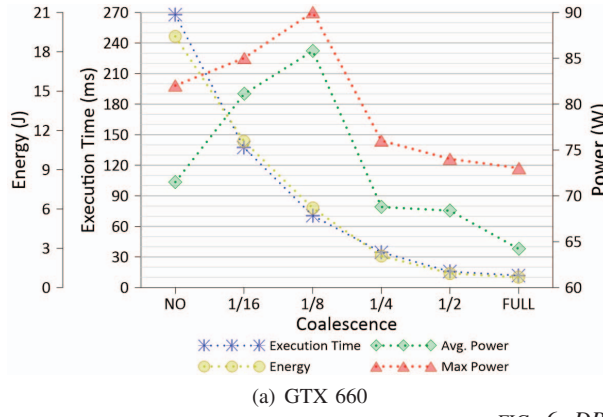


FIG. 6: DRAM coalescence

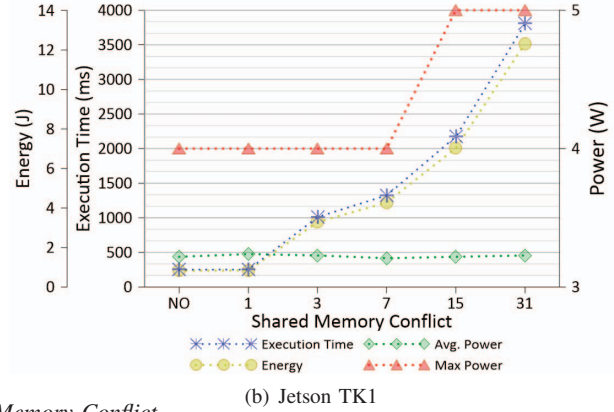
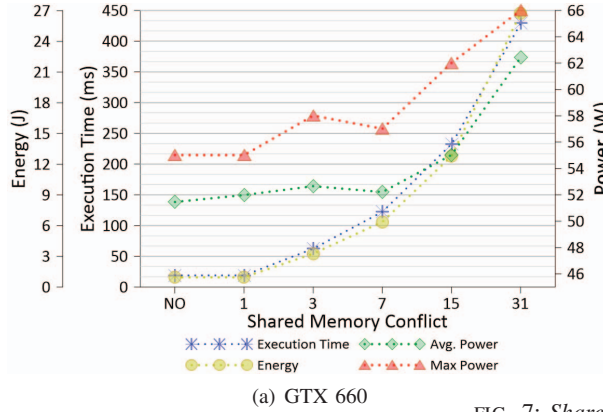


FIG. 7: Shared Memory Conflict

on a device, while it may not on a different device (see for instance the different effect on peak power on GTX660 and TK1 by increasing the memory coalescence). Combined with the standard profiler information, the proposed microbenchmark suite can efficiently guide developers in choosing among the possible optimizations during the whole iterative tuning flow.

V. CONCLUSIONS

This paper presented *MIPP*, a suite of microbenchmarks that aims at characterizing a GPU device in terms of performance, power, and energy consumption. *MIPP* aims at understanding how application bottlenecks involving selected functional components or underutilization of them can affect device performance, power consumption, and energy efficiency on a given device. The paper presented the results obtained by applying the microbenchmark suite to characterize two different GPU devices, i.e., an NVIDIA Kepler GTX660 and a low power embedded system NVIDIA Jetson TK1. The results showed how the same code optimizations have a different impact on the design constraints on the two GPU architectures.

REFERENCES

- [1] NVIDIA Tegra X1. <http://www.nvidia.com/object/tegra.html>.
- [2] Qualcomm Snapdragon. <http://www.qualcomm.com/products/snapdragon>.
- [3] D. Bedard, M. Y. Lim, R. Fowler, and A. Porterfield. Powermon: Fine-grained and integrated power monitoring for commodity computer systems. In *Proc. of IEEE SoutheastCon*, pages 479–484, 2010.
- [4] Z. Fang, S. Mehta, P.-C. Yew, A. Zhai, J. Greensky, G. Beeraka, and B. Zang. Measuring microarchitectural details of multi- and many-core memory systems through microbenchmarking. *ACM Transactions on Architecture and Code Optimization*, 11(4):art. n.5, 2015.
- [5] W.-c. Feng and K. W. Cameron. The green500 list: Encouraging sustainable supercomputing. *Computer*, 40(12):50–55, 2007.
- [6] S. Hong and H. Kim. An integrated gpu power and performance model. In *Proc. of ACM/IEEE ISCA*, pages 280–289, 2010.
- [7] T. Komoda, S. Hayashi, T. Nakada, S. Miwa, and H. Nakamura. Power capping of cpu-gpu heterogeneous systems through coordinating dvfs and task mapping. In *Proc. of IEEE ICCD*, pages 349–356, 2013.
- [8] K. Ma, X. Li, W. Chen, C. Zhang, and X. Wang. Greengpu: A holistic approach to energy efficiency in gpu-cpu heterogeneous architectures. In *Proc. of IEEE ICPP*, pages 48–57, 2012.
- [9] X. Mei, K. Zhao, C. Liu, and X. Chu. Benchmarking the memory hierarchy of modern gpus. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 8707 LNCS:144–156, 2014.
- [10] S. Mittal and J. S. Vetter. A survey of methods for analyzing and improving gpu energy efficiency. *ACM Comput. Surv.*, 47(2):19:1–19:23, Aug. 2014.
- [11] N. Nishikawa, K. Iwai, H. Tanaka, and T. Kurokawa. Throughput and power efficiency evaluation of block ciphers on kepler and gcn gpus using micro-benchmark analysis. *IEICE Transactions on Information and Systems*, E97-D(6):1506–1515, 2014.
- [12] NVIDIA. PTX: Parallel Thread Execution ISA, 2015. <http://docs.nvidia.com/cuda/parallel-thread-execution/>.
- [13] Nvidia CUDA. Programming guide, 2015. <http://docs.nvidia.com/cuda/cuda-c-programming-guide>.
- [14] J. Sim, A. Dasgupta, H. Kim, and R. Vuduc. A performance analysis framework for identifying potential benefits in gpgpu applications. In *Proceedings of ACM SIGPLAN PPoPP*, pages 11–22, 2012.
- [15] Y. Wang, S. Roy, and N. Ranganathan. Run-time power-gating in caches of gpus for leakage energy savings. In *Proceedings of ACM/IEEE Design, Automation and Test in Europe, DATE*, pages 300–303, 2012.
- [16] H. Wong, M.-M. Papadopoulos, M. Sadooghi, and A. Moshovos. Demystifying GPU microarchitecture through microbenchmarking. In *Proc. of IEEE ISPASS*, pages 235–246, 2010.
- [17] X. Yan, X. Shi, L. Wang, and H. Yang. An opencl micro-benchmark suite for gpus and cpus. *Journal of Supercomputing*, 69(2):693–713, 2014.
- [18] Y. Yang, P. Xiang, M. Mantor, and H. Zhou. Fixing performance bugs: An empirical study of open-source gpgpu programs. In *Proc. of IEEE ICPP*, pages 329–339, 2012.