

Exploiting Multi-Loop Parallelism on Heterogeneous Microprocessors

Michael Zuzak and Donald Yeung
 Department of Electrical and Computer Engineering
 University of Maryland at College Park
 {mzuzak,yeung}@umd.edu

Abstract

Heterogeneous microprocessors integrate CPUs and GPUs on the same chip, providing fast CPU-GPU communication and enabling cores to compute on data “in place.” These advantages will permit integrated GPUs to exploit a smaller unit of parallelism. But one challenge will be exposing sufficient parallelism to keep all of the on-chip compute resources fully utilized. In this paper, we argue that integrated CPU-GPU chips should exploit parallelism from multiple loops simultaneously. One example of this is nested parallelism in which one or more inner SIMD loops are nested underneath a parallel outer (non-SIMD) loop. By scheduling the parallel outer loop on multiple CPU cores, multiple dynamic instances of the inner SIMD loops can be scheduled on the GPU cores. This boosts GPU utilization and parallelizes the non-SIMD code. Our preliminary results show exploiting such multi-loop parallelism provides a 3.12x performance gain over exploiting parallelism from individual loops one at a time.



1 Introduction

Traditionally, GPUs have been implemented as discrete chips on daughter cards, but recently, processor manufacturers have been producing *heterogeneous microprocessors* in which the CPU and GPU are integrated on the same die [1], [2], [3]. Such heterogeneous chips provide a single image of physical memory to both the CPU and SIMT (single-instruction multiple-thread) cores, resulting in a seamless shared address space. This allows all cores to compute on data “in place,” eliminating copying between separate CPU and GPU memories. In addition, the tight integration permits extremely fast CPU-GPU communication—*e.g.*, through off-chip physical memory.

The vastly improved CPU-GPU communication speeds, along with not having to copy data, imply that integrated GPUs can exploit a much smaller amount of parallelism compared to discrete GPUs which require massive parallelism to amortize their large startup latencies [4]. This will allow integrated GPUs to accelerate a wide range of loops. Furthermore, the support for shared memory between CPU and GPU will also greatly simplify programming—*e.g.*, programmers will be spared the onerous task of having to identify what data to communicate. Together, flexibility to off-load finer-grained loops coupled with reduced programming effort will enable programmers to map more complex programs onto integrated GPUs.

As researchers try to accelerate more complex programs, a major challenge will be parallelizing codes to keep heterogeneous microprocessors fully utilized. The conventional approach is to parallelize loops one at a time, and to schedule each loop separately on the GPU. In this case, only the GPU runs parallel code, and it only runs one parallel loop at a time. Moreover, all unparallelized code regions are scheduled on a single CPU core. This is effective for programs with very large SIMD loops that dominate the program’s execution.

Unfortunately, this approach is *not* effective for more complex programs. While there may still be many GPU-friendly SIMD loops in complex codes, the amount of work in each loop can vary significantly. In many cases, there may not be sufficient parallelism to fully utilize the SIMT cores [5]. Moreover, complex programs also contain code with control divergence and irregular memory access patterns. And, these non-SIMD code regions tend to account for significant portions of execution time. By running them serially, significant performance degradation will occur. It will also underutilize the multiple CPU cores in a heterogeneous microprocessor, which along with the GPU are also a tremendous source of compute power.

We propose to exploit parallelism from multiple loops simultaneously when possible to more effectively utilize heterogeneous microprocessors. Such *multi-loop parallelism* is a generalization of existing parallel idioms, like *nested parallelism* or *pipeline parallelism*. In this paper, we focus on the former in which one or more SIMD loops are nested underneath a parallel outer loop, a common code structure in our benchmarks. We schedule the parallel outer loop on the CPU cores. The resulting multiple CPU threads then spawn multiple instances of the inner SIMD loop(s), often simultaneously, which get scheduled on the SIMT cores. (While the frequency of spawns can be high if the inner loop has smaller amounts of work, the high-speed communication between a CPU and an integrated GPU enables the exploitation of such finer-grained SIMD loops). This exposes more parallelism for the SIMT cores and parallelizes the non-SIMD outer loop. Our preliminary results, which estimate performance by combining separate CPU and GPU simulations, show that multi-loop parallelism outperforms single-loop parallelism by 3.12x across 3 OpenMP benchmarks.

The rest of this paper is organized as follows. Section 2 presents our new parallelization scheme. Then, Section 3 undertakes a quantitative evaluation of its effectiveness. Next, Section 4 discusses related work. Finally, Section 5 concludes the paper.

2 Multi-Loop Parallelism

Traditional GPU workloads [6], [7] consist of massively parallel SIMD loops within which programs spend the vast majority of their time. (The massive loops have been necessary to amortize the large kernel initiation costs associated with discrete GPUs).

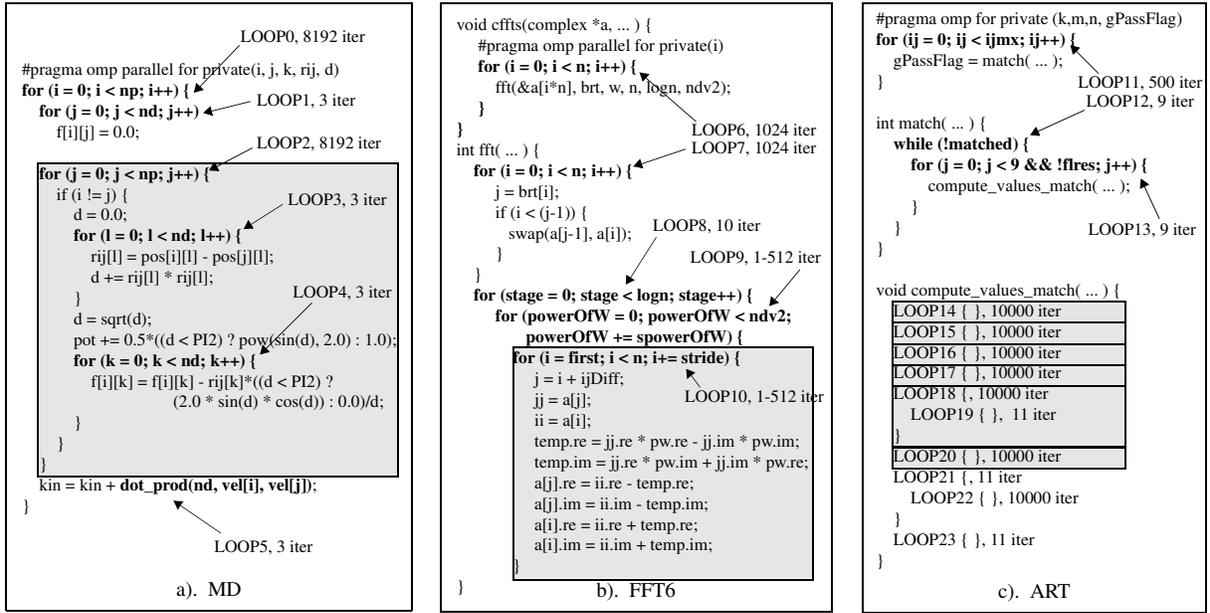


Fig. 2. Code examples from the a). MD, b). FFT6, and c). ART benchmarks exhibiting multi-loop parallelism. Our technique schedules the parallel outer loops (OpenMP pragmas) on CPUs, and the parallel inner loops (shaded gray) on GPUs.

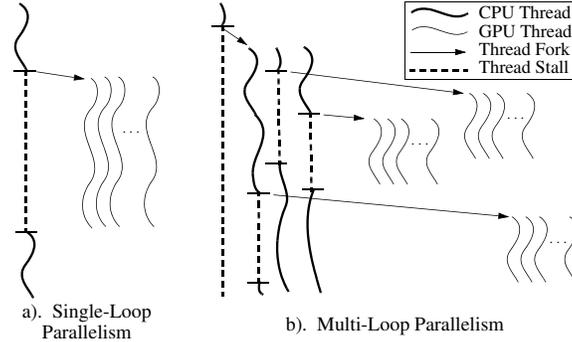


Fig. 1. CPU and GPU threads in a). single-loop and b). multi-loop parallelism.

In such workloads, it is sufficient to parallelize the SIMD loops individually, and to schedule them one at a time on the GPU. Figure 1a illustrates this *single-loop parallelism* case. The program starts out running non-SIMD code serially in a single CPU thread. When the CPU thread reaches a parallel SIMD loop, it spawns the corresponding GPU kernel which gets scheduled on the GPU. The GPU's SIMT cores then execute the SIMD loop in parallel as the CPU thread stalls. When GPU execution completes, the CPU thread is notified and continues executing.

In systems with integrated GPUs, the kernel initiation costs are much lower. This raises the possibility to use GPUs not only for massively parallel loops, but also for smaller loops (*i.e.*, with fewer loop iterations). Despite containing less work, such loops can still be performance critical if they occur within loop nests that execute them a large number of times. Unfortunately, scheduling smaller SIMD loops on the GPU one at a time as is done in Figure 1a will perform poorly because 1). the small amounts of parallel work may not fully utilize the GPU, and 2). the computations in between the SIMD loops—*i.e.*, the rest of the loop nest which would execute serially—may also be important, limiting gains due to Amdahl's Law.

To address this problem, we propose exploiting parallelism from multiple loops simultaneously—*i.e.*, *multi-loop parallelism*. In addition to executing SIMD loops on the GPU, our approach tries to *find additional parallelism outside of the SIMD loops to enable parallel execution on multiple CPU cores as well*. One source of multi-loop parallelism is nested parallelism. In this case, we look for one or more parallel SIMD loops that are nested underneath an outer loop that is also parallel.

Figure 1b illustrates the multi-loop parallelism case. Like Figure 1a, the program starts out running serially in a single CPU thread. The CPU thread again spawns multiple threads, but this time it spawns CPU threads corresponding to a parallel outer (non-SIMD) loop that gets scheduled on multiple CPU cores. Then, each CPU thread reaches a nested parallel SIMD loop, and spawns the corresponding GPU kernel which gets scheduled on the GPU. Exploiting such multi-loop parallelism increases the performance of heterogeneous microprocessors in two ways. First, the multiple dynamic instances of the SIMD loops provide more parallelism to boost the GPU's utilization, especially when each SIMD loop is smaller. And second, portions of the loop nests outside of the SIMD loops that would have otherwise executed serially now execute in parallel on the CPU cores, thus addressing Amdahl's Law.

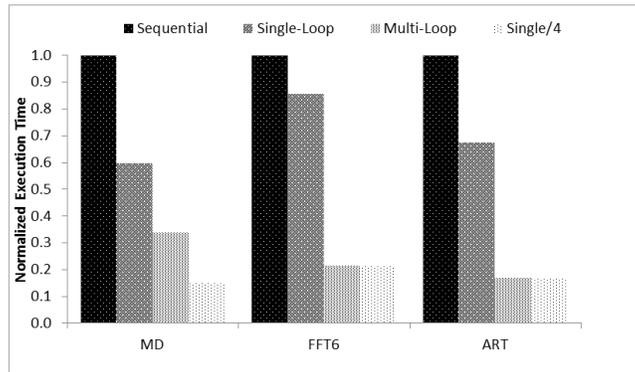


Fig. 3. Normalized execution time for sequential, single-loop and multi-loop parallelism, and multi-loop parallelism without contention.

2.1 Code Examples

We surveyed different programs to look for multi-loop parallelism—specifically nested parallelism—and to understand the nature of these loop structures. We avoided looking at dense numeric codes commonly found in GPU benchmark suites. Instead, we focused on codes with more complex loop nests that one would not normally think of accelerating using GPUs.

Figure 2 presents three representative examples from our code survey: MD and FFT6 from the OpenMP source code repository [8], and ART from the SPEC OMP 2001 suite [9]. All of these benchmarks have been parallelized using OpenMP [10]. In the figure, we show the main parallelized loop from each benchmark—*i.e.*, the loop identified by the “#pragma omp” directive. We treat this as the parallel outer loop for one instance of multi-loop parallelism. We also show all of the loops nested underneath the parallel outer loop, labeling each inner loop and indicating its iteration count. (For the inner-most loops in ART—*i.e.*, in the “compute_values_match()” function—we only show the loop labels in place of the full code due to the large size of this example).

The code examples in Figure 2 illustrate several important features of the programs that our techniques try to target. First, there are no massive loops that dominate the computation. Instead, the codes contain a number of smaller loops. Some loops exhibit a trivial number of iterations (≤ 11) while others contain modest iteration counts (100s to a few 1000 iterations). In addition to being smaller, some loops contain control divergence (*e.g.*, LOOP2 and LOOP7) or strided memory accesses (*e.g.*, LOOP10) that can reduce GPU performance. The codes also exhibit complex loop structure. There is a large number of nesting levels, and in some cases, nesting occurs across multiple function calls. ART in Figure 2c is an extreme example of code distributed across a large number of deeply nested loops.

Figure 2 also shows the loops exhibit nested parallelism. In addition to the parallel loops that are explicitly identified by the OpenMP pragmas, a careful examination of the other loops inside the loop nests reveals a number of them are parallel as well. (Although not shown in Figure 2c, all of the inner loops in ART, labeled LOOP14–LOOP23, are parallel). We find this is fairly common in OpenMP programs because the marked parallel loops often appear at the outer levels of deep loop nests, as in Figure 2. Such coarse-grained parallelism is almost impossible for a compiler to extract automatically, but is natural for a programmer to express given his/her knowledge of the code at the algorithm level. Moreover, programmers are incentivized to express coarse-grained parallelism since it is a good match for CPU cores, the main target for OpenMP programs. Given such explicitly parallel outer loops, nested parallelism arises whenever one or more inner loops are found to be parallel, which is quite likely in OpenMP programs when each parallel region contains many inner loops.

Notice, both CPUs and GPUs are needed to exploit the multiple levels of parallelism in Figure 2. Inner loops with non-trivial iteration counts have a greater chance for GPU acceleration, whereas loops at outer levels are more appropriate for CPUs. Such codes are a good match for heterogeneous microprocessors.

3 Experimental Results

This section conducts a preliminary evaluation of multi-loop parallelism for the codes in Figure 2. We do not yet have these codes running on simulators of heterogeneous microprocessors. Instead, we employ simulators that model CPUs and GPUs *separately*, and then we estimate the integrated CPU-GPU performance by combining the separately acquired results.

Our study uses SimpleScalar [11] to simulate a single CPU core, and GPGPU-Sim [12] to simulate a GPU. We configure SimpleScalar to model a 4-way out-of-order core running at 2.4 GHz with an 128-entry reorder buffer. Our CPU core has a 2-level cache hierarchy with a split 16KB L1 cache and a unified 256KB L2 cache. We configure GPGPU-Sim to model an Nvidia GTX-480 running at 700 MHz. We simulate 8 streaming multiprocessors (SMs), each containing 16 streaming processors (SPs) that support 32 threads at a time (a warp). Each SM also contains a 16KB L1 cache, and all 8 SMs share a 786KB L2 cache. For both simulators, we assume off-chip main memory incurs a latency of 50 ns and provides 88.7 GB/s of memory bandwidth.

First, we estimate the performance of single-loop parallelism—*i.e.*, off-loading loops one at a time as shown in Figure 1a. We timed every loop from Figure 2 on both SimpleScalar and GPGPU-Sim. For SimpleScalar, we simulated each OpenMP region once using a binary in which each loop entry and exit point was marked so that we could break down the execution time per loop. For GPGPU-Sim, we created multiple CUDA kernels to off-load each loop individually onto the GPU, and simulated each kernel separately on GPGPU-Sim. (Although GPGPU-Sim assumes discrete GPUs, we omit the data transfer times since these would not be incurred by integrated GPUs computing on the data “in place”). After acquiring these results, we identified the loops that run faster on the GPU. These are the loops shaded gray in Figure 2. To estimate the performance for off-loading these GPU-friendly loops, we performed a second run of each OpenMP region on SimpleScalar, but this time, we flushed the CPU’s caches at each entry and exit point to those loops. (The original GPU simulations for each loop already start with cold caches and include a flush

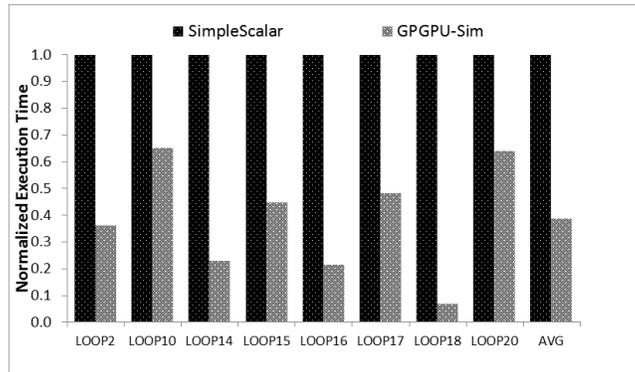


Fig. 4. Per-loop execution time on the GPU and CPU.

at the end). Then, we constructed a trace of execution through each OpenMP region, alternating between the CPU core (with flushing) and the GPU at the off-loaded loop entry and exit points. We also added a fixed 100-cycle delay for each off-loading to account for initiating a new kernel on the GPU. (This is a conservative number considering the CPU and GPU are on the same chip).

Figure 3 presents the results. The bars in Figure 3 labeled “Sequential” report the execution time when the entire OpenMP execution trace runs on the CPU core, and the bars labeled “Single-Loop” report the estimated execution time when the execution trace adopts the GPU performance for those loops that run faster on the GPU. All bars are normalized to the sequential bars. As Figure 3 shows, the GPU provides a speedup of 1.67x, 1.17x, and 1.48x for MD, FFT6, and ART, respectively. On average, the speedup is 1.44x.

While these speedups are significant, they are also quite disappointing considering the capability of the GPU. The problem is the complex nature of the codes in Figure 2, as described earlier. In particular, the codes contain a mix of SIMD and non-SIMD loops, so only a portion of each OpenMP region can run on the GPU. As Figure 2 shows, there is a lot of unshaded code that must run on the CPU core. This limits the overall gain from the GPU due to Amdahl’s Law. To make matters worse, the gains for the SIMD loops are mixed. Figure 4 compares the execution time of each shaded (off-loaded) loop from Figure 2 running on SimpleScalar and GPGPU-Sim. In some cases, the gains are large (14.3x speedup for LOOP18), but in other cases the gains are small (1.53x speedup for LOOP10). This is due to the characteristics of our more complex codes: small number of iterations (LOOP10), control flow divergence (LOOP2), and non-unit stride memory accesses (LOOP10, LOOP14, LOOP16, and LOOP18).

Next, we estimate the performance for multi-loop parallelism. As before, we construct a trace of execution through each OpenMP region, but this time we build a parallel trace. We assume the outer-most parallel loops identified by the OpenMP pragmas in Figure 2 run on 4 CPU cores, so we divide all of their iterations into 4 blocks that are executed simultaneously by 4 parallel traces. Within each trace, we again alternate between the CPU core and the GPU (*i.e.*, for the unshaded and shaded codes in Figure 2, respectively), selecting the CPU execution or GPU execution from the SimpleScalar/GPGPU-Sim simulations to estimate the speed of each trace. We also assume the 4 CPU cores *share* the 8-SM GPU, so there is contention whenever 2 or more CPU cores off-load loops onto the GPU. To model this, we schedule loops in the order that they arrive at the GPU, delaying later loops behind earlier loops. Moreover, as a loop’s thread blocks complete and the next loop’s thread blocks begin execution, there may be multiple kernels in the GPU at the same time. We ran additional simulations on GPGPU-Sim to measure the speed of each off-loaded loop when paired with another off-loaded loop in the GPU, trying all possible ways of dividing the 8 SMs between the pair. (For ART, we did this for all pairings of LOOP14, LOOP15, LOOP16, LOOP17, LOOP18, and LOOP20).

In Figure 3, the bars labeled “Multi-Loop” report the estimated execution time for multi-loop parallelism normalized to the “Sequential” bars. As Figure 3 shows, multi-loop parallelism significantly out-performs single-loop parallelism. Compared to the “Sequential” bars, multi-loop parallelism provides an overall speedup of 2.94x, 4.62x, and 5.91x for MD, FFT6, and ART, respectively. On average, the gain is 4.49x which is 3.12x greater than single-loop parallelism’s gain. This is due to the parallel execution of the non-SIMD code in the OpenMP regions, thus addressing the Amdahl’s Law limitation of single-loop parallelism. It is also due to executing multiple loops simultaneously on the GPU, providing more work (especially for smaller loops like LOOP10) to better utilize the SIMT cores. We find single-loop parallelism usually under-utilizes the GPU, leaving ample GPU compute bandwidth for multi-loop parallelism. To show this, the bars labeled “Single/4” in Figure 3 report the single-loop parallelism execution time divided by 4, which represents the ideal speedup of multi-loop parallelism (*i.e.*, in the absence of GPU contention). For FFT6 and ART, there is almost no difference between the “Multi-Loop” and “Single/4” bars, indicating virtually no contention. For MD, the ideal speedup is another 2.27x faster, indicating contention is an issue. But even in this case, multi-loop parallelism is still 1.8x faster than single-loop parallelism.

4 Related Work

Simultaneously utilizing CPU and SIMT cores in a heterogeneous microprocessor is not a new idea. Recent work has already proposed doing this [13], [14], [15]. But these existing techniques still only exploit single-loop parallelism: they schedule iterations from a single SIMD loop across both types of cores. As we have shown, single-loop parallelism cannot keep the GPU (let alone the CPU and GPU) fully utilized for complex codes with small SIMD loops, nor can it parallelize the non-SIMD code. By exploiting parallelism from multiple loops within a loop nest, we expose greater amounts of SIMD parallelism and parallelize the non-SIMD code as well.

Multi-loop parallelism is also related to dynamic thread-block launch (DTBL) [5]. In DTBL, a loop off-loaded by the CPU onto

the GPU can initiate other instances of itself from the GPU. (This supports certain forms of dynamic parallelism—for example, vertex expansion during recursive graph search). Like multi-loop parallelism, DTBL enables multiple dynamic instances of a loop to execute simultaneously on the GPU. However, DTBL only exploits the GPU. It does not try to schedule loops onto the CPU and GPU simultaneously, and hence, is not designed for heterogeneous microprocessors. Moreover, DTBL only exposes SIMD parallelism, and thus cannot parallelize non-SIMD loops.

5 Conclusion

Given the wide availability of heterogeneous microprocessors, we believe an exciting research direction is to find new computations that can benefit from integrated GPUs. We argue this will require new parallelization methods that can make effective use of both CPU and SIMT cores simultaneously. In this paper, we propose one such novel method, multi-loop parallelization, and analyze its benefits for nested parallel loops. Our results show exploiting multi-loop parallelism can outperform the conventional approach of parallelizing loops one at a time by 3.12x for programs with complex loop nests. In the future, we plan to investigate different parallelization strategies that can support other types of looping structures.

References

- [1] "Intel Corporation. Intel Sandy Bridge Microarchitecture. <http://www.intel.com>."
- [2] N. Brookwood, "AMD Fusion Family of APUs: Enabling a Superior, Immersive PC Experience. AMD White Paper." 2010.
- [3] M. Wilkins, "NVIDIA Jumps on Graphics-Enabled Microprocessor Bandwagon," 2011.
- [4] C. Gregg and K. Hazelwood, "Where is the Data? Why You Cannot Debate CPU vs. GPU Performance Without the Answer," in *Proceedings of the International Symposium on Performance Analysis of systems and Software*, 2011.
- [5] J. Wang, N. Rubin, A. Sidelnik, and S. Yalamanchili, "Dynamic Thread Block Launch: A Lightweight Execution Mechanism to Support Irregular Applications on GPUs," in *Proceedings of the International Symposium on Computer Architecture*, Portland, OR, June 2015.
- [6] S. Che, J. W. Sheaffer, M. Boyer, L. G. Szafaryn, L. Wang, and K. Skadron, "A Characterization of the Rodinia Benchmark Suite with Comparison to Contemporary CMP Workloads," in *Proceedings of the International Symposium on Workload Characterization*, Atlanta, GA, December 2010.
- [7] J. A. Stratton, C. Rodrigues, I.-J. Sung, N. Obeid, L.-W. Chang, N. Anssari, G. D. Liu, and W. mei Hwu, "The Parboil Technical Report," March 2012.
- [8] "OpenMP Source Code Repository. <http://www.pcg.ull.es/omp/scr/>." 2004.
- [9] "SPEC OMP 2001. <https://www.spec.org/omp2001/>." 2001.
- [10] "The OpenMP API Specification for Parallel Programming. Intel Corporation. <http://www.openmp.org/wp/>." 2014.
- [11] D. Burger and T. M. Austin, "The SimpleScalar Tool Set, Version 2.0," University of Wisconsin-Madison, CS TR 1342, June 1997.
- [12] A. Bakhoda, G. L. Yuan, W. W. L. Fung, H. Wong, and T. M. Aamodt, "Analyzing CUDA Workloads Using a Detailed GPU Simulator," in *Proceedings of the International Symposium on Performance Analysis of Systems and Software*, Boston, MA, April 2009.
- [13] R. Kaleem, R. Barik, T. Shpeisman, B. T. Lewis, C. Hu, and K. Pingali, "Adaptive Heterogeneous Scheduling for Integrated GPUs," in *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques*, Edmonton, AB, Canada, August 2014.
- [14] V. T. Ravi and G. Agrawal, "A Dynamic Scheduling Framework for Emerging Heterogeneous Systems," in *Proceedings of the 18th International Conference on High Performance Computing*, December 2011.
- [15] V. T. Ravi, W. Ma, D. Chiu, and G. Agrawal, "Compiler and Runtime Support for Enabling Generalized Reduction Computations on Heterogeneous Parallel Configurations," in *Proceedings of the International Conference on Supercomputing*, Tsukuba, Ibaraki, Japan, June 2010.