

Characterising Optimizations to Memory Access Patterns using Architecture-Independent Program Features

Aditya Chilukuri
aditya.chilukuri@anu.edu.au
Australian National University

Josh Milthorpe
josh.milthorpe@anu.edu.au
Australian National University

Beau Johnston
beau.johnston@anu.edu.au
Australian National University

ABSTRACT

HPC developers are faced with the challenge of optimizing OpenCL workloads for high performance on novel architectures. The Architecture-Independent Workload Characterization (AIWC) tool is a plugin for the Oclgrind OpenCL simulator that gathers metrics of OpenCL programs used to understand and predict program performance on an arbitrary given hardware architecture. However, AIWC metrics are not always easily interpreted and do not reflect some important memory access patterns affecting efficiency across architectures. We propose the concept of parallel spatial locality: the locality of memory accesses by OpenCL work-items (threads) executing simultaneously. We implement the parallel spatial locality metric in the AIWC framework and analyse gathered results on matrix multiply and the Extended OpenDwarfs OpenCL benchmarks. The differences in the parallel spatial locality metric between implementations of matrix multiply reflect the optimizations performed. The new metric showcases the memory access patterns in the OpenDwarfs benchmarks that affect their performance on various architectures. The improvements suggested to AIWC will help HPC developers better understand memory access patterns of complex codes, a crucial step in guiding optimization of codes for arbitrary hardware targets.

KEYWORDS

Architecture Independent Analysis, Heterogenous Computing, Workload Characterization, Memory Access Patterns

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

IWOCL'20, Munich, Germany

© 2020 ACM. 978-x-xxxx-xxxx-x/YY/MM...\$15.00

DOI: 10.1145/nnnnnnnn.nnnnnnn

ACM Reference format:

Aditya Chilukuri, Josh Milthorpe, and Beau Johnston. 2020. Characterising Optimizations to Memory Access Patterns using Architecture-Independent Program Features. In *Proceedings of International Workshop on OpenCL, Munich, Germany, April 27–29, 2020 (IWOCL'20)*, 11 pages.

DOI: 10.1145/nnnnnnnn.nnnnnnn

1 INTRODUCTION

High-performance computing (HPC) systems are increasingly heterogeneous. A single node on a modern supercomputer may combine traditional CPUs with one or more accelerators such as GPUs, Field Programmable Gate Arrays (FPGAs), or many integrated core devices (MICs). High bandwidth interconnects support tight integration between multiple devices of different types on a single compute node.

The OpenCL programming language is designed to support portability for heterogeneous systems by allowing a single code to execute on multiple hardware targets. The language provides support for individual *kernels*, portions of parallelizable codes that make up a HPC code base, to be executed independently on diverse hardware targets such as traditional CPUs, GPUs and FPGAs [11]. This gives HPC software developers greater flexibility by allowing codes aimed at a range of hardware targets to be written in a single programming language environment.

Application codes differ in resource requirements, control structure and available parallelism. Similarly, compute devices differ in number and capabilities of execution units, processing model, and available resources. These heterogeneous computing environments present opportunities for software engineers and HPC integrators to design highly optimized systems with multiple kernels executing on hardware targets best suited for the diverse computational tasks performed by each kernel [16]. However, this opportunity also presents a challenge in optimizing code to run on diverse architectures. The work reported here aims to help HPC developers understand and optimize memory access patterns to improve performance on heterogeneous systems.

In the evolution of computer architectures over the last few decades, the exponential growth in computational capability

has not been matched by proportional increases in memory speeds[6]. Memory accesses pose larger bottlenecks to performance as application demand for main memory scales with the arithmetic capability of computer systems. To mitigate this latency, modern CPU designs have employed a wide range of cache technologies to reduce main memory accesses using the principle of *spatial locality*: the observation that data that a program accesses close together in time tend to also be close together in memory. On the other hand, GPUs rely on hardware multithreading to hide memory latency, and their architectures favour ALU capability over sophisticated logic to manage a cache hierarchy and out-of-order execution. As a result, the performance of kernels on CPUs and GPUs alike is strongly dependent on memory access patterns intrinsic to the code.

Our aim in this paper is to develop a framework to guide HPC software engineers in hardware-dependent code optimization – specifically by guiding the improvement of memory access patterns. We provide examples of manufacturer-recommended code optimizations to improve memory access patterns on the target architecture, and examine these using the architecture-independent workload characterization (AIWC) [8] plugin for Oclgrind. Our work highlights the benefits and challenges arising from an architecture-independent analysis of memory-based program characteristics. We propose two new metrics to AIWC to better characterize these memory-based optimizations. We measure the presented codes using the new metrics and demonstrate the metrics' effectiveness in capturing the essence of the optimizations performed.

The structure of this paper is as follows. In section 2 we use the example of matrix multiplication on a GPU to showcase the specific vendor-recommended optimizations our work aims to measure. In section 3 we discuss how AIWC and its precursors profiled memory access patterns in an architecture-independent fashion, and relevant architecture dependent approaches to memory access profiling. In section 4 we present a methodology to critically examine the architecture-independent workload characterization tool, which we use in section 5 to measure the accuracy of AIWC in profiling the impact of performance optimizations to various OpenCL codes. In section 6 we suggest new methods to accurately capture code patterns correlating to hardware-specific memory optimizations. Finally, in section 8 we discuss further avenues to extend our work and conclude.

2 MOTIVATING EXAMPLE

We aim to capture the essence of vendor-recommended optimizations for target architectures in our metrics. We demonstrate the optimization strategies using an OpenCL kernel

which multiplies square matrices of order N . We start with a simple unoptimized kernel, and then improve the memory access patterns of this kernel by incrementally performing the optimizations recommended in the CUDA Optimization Handbook [3] for NVIDIA GPUs. These incrementally optimized OpenCL codes are then analysed using AIWC to determine the accuracy in profiling favourable memory access patterns of the current framework and the proposed extensions to AIWC.

Simple Unoptimized Matrix Multiplication

The unoptimized matrix multiplication kernel presented in Appendix A is used as a baseline to validate the performance improvements measured from each following optimization. Each thread of the kernel updates the (*globalRow*, *globalCol*) element of matrix C by computing the dot product of the corresponding row of A and column of B .

Using Shared Memory to Coalesce Global Memory Access to Matrix A

We first notice that the number of global memory accesses to matrices A and B increases as $O(N^3)$ with respect to matrix size N . Global memory is typically located off-chip and accesses induce large delays. NVIDIA GPUs coalesce global memory loads and stores issued within thread-groups into as few DRAM transactions as possible. Multiple global memory loads and stores are coalesced into a single transaction when certain device-specific conditions are met. On most NVIDIA GPUs, data accesses are coalesced when multiple requests are made for memory locations from the same cache line in global memory [4]. Appendix B contains the `coalescedAMultiply` kernel, which coalesces accesses to matrix A by storing *tiles*, or square blocks, of A 's values into shared memory.

Using Shared Memory to store Tiles of Matrix B

The code is further optimized by improving how the previous `coalescedAMultiply` kernel reads from Matrix B . While each thread in `coalescedAMultiply` reads only a single element of matrix A from global memory, each thread reads a full column of matrix B . The repeated reads of elements of matrix B across multiple threads can be alleviated by reading tiles from matrix B into shared memory. Appendix C contains the `coalescedAB` kernel which performs this optimization.

Optimising handling of shared memory

NVIDIA shared memory is divided into multiple banks – stored in independent memory modules – to allow parallel accesses to shared memory. Bank conflicts occur when shared memory in the same bank is accessed concurrently. The code

in coalescedAB kernel is further optimized by implicitly transposing tiles of Matrix *A* while loading from global memory. This improves memory bank utilization during reads to shared memory tiles are performed during the necessary computation. The coalescedABT kernel demonstrates this optimization by modifying lines 20,21 of coalescedAB to:

```
20   ASub[localCol][localRow] = A[tiledRow];
21   BSub[localRow][localCol] = B[tiledCol];
```

Alignment of Memory Allocation

Memory access alignment is important to best utilize all parts of GPU memory architecture. Global memory buffer alignment can allow threads to access blocks of global memory aligned to the nearest cache line. This enables coalescing of memory accesses. If buffers are misaligned, parallel memory requests may cross over cache lines, and may double the number of slow global memory accesses needed as demonstrated in figures 1 and 2.

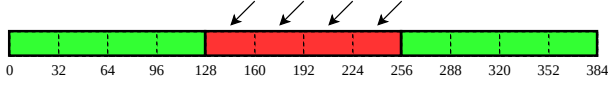


Figure 1: All threads access memory aligned to nearest cache line in parallel [3]

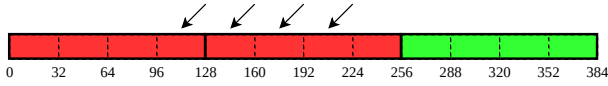


Figure 2: Unaligned sequential memory addresses fit in two cache lines [3]

A similar principle applies when using shared memory. Bank conflicts may be reduced by aligning allocations of shared memory buffers. The alignedABT kernel improves the alignment of shared memory tiles in coalescedABT. An arbitrary large alignment value of 4096 is chosen as it is larger than the cache line size of modern hardware and attributed to the local array declarations in the existing code as shown below.

```
1   __local float aTile[TILE_DIM][TILE_DIM]
    __attribute__((aligned(4096)));
2   __local float bTile[TILE_DIM][TILE_DIM]
    __attribute__((aligned(4096)));
```

All the examples above are an incomplete listing of optimization strategies developers can apply to code targeted at NVIDIA GPUs. In section 5, we will compare the collected

AIWC metrics with the expected effect of each optimization, to examine how effectively our methodology uncovers underlying bottlenecks and guides optimization efforts.

3 RELATED WORK

Hoste and Eeckhout [7] show that while conventional architecture dependent characteristics are useful in locating performance bottlenecks, they can hide the underlying, inherent program behaviour causing performance bottlenecks. A conceptual understanding of the performance characteristics of complex codes is necessary for the programmer to effectively optimize these codes. Architecture dependent characteristics typically include instructions per cycle (IPC), cache and branch prediction miss rates, page faults and DRAM bus data transfer rates. These are typically collected using hardware performance counters available on most target architectures. These performance counters do not serve to guide optimization beyond highlighting potential bottlenecks [5, 7]. Further, in many cases, architecture dependent characteristics usually cannot be directly correlated to specific code patterns. For example, the causes of high cache miss rates in the execution of a program are too complex and widely variant on microarchitecture specific features such as cache size, prefetch behaviour and cache placement policies. In summary, we argue that a HPC developer tasked with optimizing code for a given hardware target would benefit from architecture-independent metrics of the code that can be conceptually correlated to hardware specific optimization patterns. Such architecture-independent metrics will effectively guide the HPC developer’s workflow in both finding and fixing bottlenecks in code performance.

The Architecture-Independent Workload Characterization (AIWC) tool [8] relies on a set of instruction set architecture (ISA)-independent features determined by Shao and Brooks [15]. AIWC measures memory based metrics to characterize the target code’s memory behaviours. AIWC firstly collects aggregate simple counts such as the (i) total memory footprint, the total number of unique addresses accessed; and (ii) 90% memory footprint, the number of unique addresses covering 90% of memory accesses. A small ratio of 90% memory footprint to total memory footprint correlates to programs accessing the same small subset of memory addresses repeatedly – which can be linked to increased performance in a cached memory hierarchy.

AIWC also records the global memory address entropy (GMAE), a positive real number corresponding to the randomness of the memory access distribution of a program. To measure locality of memory accesses, AIWC collects the local memory address entropy (LMAE) of memory addresses accessed after dropping n least significant bits of all memory addresses accessed by the program. To calculate this, AIWC

collects a frequency distribution of all non-register memory accesses by all threads in the target kernel. Using the collected frequency distribution, AIWC calculates 10 separate local memory address entropy (LMAE) values according to increasing number of least significant bits (LSB) skipped using the explicit formula for the n -bits skipped LMAE:

$$LMAE_{n-bits} = \sum_{a \in A_n} p_a \log_2(p_a^{-1}) \quad (1)$$

- A_n is the set of all addresses accessed after skipping n LSBs of each address.
- $p_a := \frac{\#access_a}{\#access_{total}}$ is the probability (calculated as relative frequency) at which each memory address is accessed.

LMAE measures the locality of memory accesses performed over the full execution of a program. A steeper drop in entropy with increasing number of bits may correlate to more localized memory accesses over the program's execution.

Limitations of the memory footprint and local memory address entropy values as memory access pattern metrics are further discussed by this work.

Advancements in Memory Access Pattern Characterization

We provide a brief summary of past work in architecture independent program characterization and Memory Access Pattern analysis:

- The Microarchitecture Independent Workload Characterization Tool (MICA) [7] was the first microarchitecture independent workload characterization tool for single threaded programs.
- The Workload ISA-Independent Characterization for Applications (WIICA) [15] extends MICA to present a framework to analyse single threaded programs independent of the ISA.
- CuMAPz: CuMAPz [12] is a CUDA memory access profiling tool and the first to present metrics for CUDA codes which can be generalized to any given NVIDIA GPU.

MICA

The MICA framework was developed by Hoste and Eeckhout [7] in response to their findings on the limitations of conventional microarchitecture dependent characteristics. Hoste and Eeckhout discovered that performance counter based approaches to profiling codes often struggled in finding underlying program features that map to improved or worsened usage of performance critical hardware features of the target architecture. The MICA framework is a holistic characterization tool, and thus collects features including

instruction mix, instruction-level parallelism, register traffic, data stream strides and branch predictability.

Of these metrics, data stream strides are of particular interest in memory access pattern profiling. MICA's stride length metric measures the distance between consecutive memory accesses in a single threaded application. For CPU architectures running single threaded applications, this metric correlates to the spatial locality of memory accesses – a measure of how closely bunched are memory access in nearby times. This is directly correlated to cache reuse rates, critical to code performance on CPUs [2].

The MICA approach was tailored for single-threaded applications as the metrics collected rely heavily on Pin instrumentation [14]. As such, MICA is unsuited to analysing conventional supercomputing workloads with heavy use of parallelism.

CuMAPz

Kim and Shrivastava [12] present CuMAPz, a CUDA memory access pattern analysis tool to guide NVIDIA GPU optimizations. CuMAPz focuses on the problem of improving CUDA application performance using NVIDIA memory-based optimizations. CuMAPz analyses CUDA codes structurally and simulates code execution on the memory hierarchy of specific NVIDIA GPU models.

During the simulated run of the target code, CuMAPz records all memory accesses in various buffers to simulate the global and shared memories on NVIDIA GPUs. Using this detailed simulation data, CuMAPz can estimate the performance critical (i) shared memory data reuse profit, (ii) profit from coalesced access (iii) memory channel skew cost and (iv) bank conflict cost characteristics of the target code.

The simulation environment described in CuMAPz and the attached analysis framework is highly specific to CUDA enabled GPUs. Replicating the CuMAPz framework for all target architectures is challenging. However, CuMAPz is an interesting simulator from the standpoint of this paper's work since it adheres to the core parts of NVIDIA GPUs' memory models in its analysis, while allowing the user to specify their GPU model specific hardware information.

Our approach is an architecture independent analysis of memory access patterns that provides metrics correlating to similar performance critical memory access optimizations as CuMAPz. We aim to further the state of the art by providing a framework to guide developers in optimizing OpenCL codes for any given target architecture. To the best of our knowledge, none of the previous works present a set of performance metrics that accurately characterize memory access patterns of parallel applications independent of the target architecture.

4 METHODOLOGY

To provide memory access pattern metrics that represent vendor recommended optimizations on GPUs, we improved AIWC by adding a new metric: relative local memory usage. This measures the proportion of all memory accesses from the symbolic execution of the kernel that occurred to memory allocated as `_local`. On GPUs, this memory address space is mapped to fast on-chip shared memory.

Relative local memory usage is an example of a metric that is useful to measure performance critical access patterns on some architectures such as GPUs, and not others, such as CPUs. CPUs do not typically have a notion of user controlled on-chip memory shared between hardware threads such as NVIDIA GPUs' shared memory. This is a natural consequence of programming for a heterogeneous system. Specific code patterns may translate to performance improvements only on certain hardware.

In the original AIWC tool [10], global and local memory address entropy (MAE) were calculated using physical addresses of memory used by the Oclgrind simulator back-end of AIWC. This caused irregularities in entropy calculations across multiple runs of the same simulation. We improved the calculation of memory address entropy by using *virtual addresses* to calculate MAE values using an abstract ideal address space on which all memory accesses by the kernel occur. This allows AIWC to accurately abstract over the hardware and ISA-specific differences in memory layouts across the diverse hardware targets.

5 RESULTS

Figure 3 shows that across varying problem sizes, the vendor recommended optimizations to the matrix multiplication code lead to increased performance. Note that a matrix size of 80 is small – overheads attached to optimization techniques may be causing the observed decrease in performance as some optimizations are performed.

Table 1 summarizes AIWC memory metrics collected for each of the matrix multiply kernels in section 2. Note that only two LMAE values are shown for brevity. We now analyse the effectiveness of AIWC metrics in profiling the NVIDIA recommended memory optimizations applied to the matrix multiplication kernel.

Relative local memory usage: As reliance on NVIDIA GPUs' shared memory increases in each kernel from `simpleMultiply` to `coalescedAB`, we find that the ratio of OpenCL local memory usage increases as expected.

Global and local MAE: Entropy measurements decrease from `simpleMultiply` to `coalescedAB`. This is explained by the optimizations performed incrementally reducing accesses to matrix *A* and *B*, and instead performing the bulk of computations on smaller tiles in on-chip shared memory.

We observe an almost completely uniform distribution of memory accesses in `simpleMultiply`, where the program makes *N* loads to each element of *A* and *B* with *N* the dimension of the matrices. The frequency distribution of memory addresses accessed become increasingly non-uniform as we perform less accesses to matrices *A* and *B* and offset this by increasing the frequency of accesses to smaller, more localized shared memory buffers. The non-uniformity of relative frequency with which each memory cell is accessed by the program translates directly to decreases in local and global memory entropy values as reliance on small shared memory buffers increases from `simpleMultiply` to `coalescedAB`.

Memory Footprint: Similar to trends in global and local MAE, we find that the ratio of 90% memory footprint to total memory footprint decreases from 60.12% for `simpleMultiply` to 0.25% (`coalescedAB`). Increased utilization of shared memory in the optimized kernels means that the shared memory buffers make up a greater proportion of total memory accesses in the program. As the shared memory buffers are small and reused within a workgroup, the memory footprint of the shared memory buffers is also smaller.

AIWC's metrics strongly reward optimizations that tend to localize memory accesses. Shared memory buffers are typically smaller than global memory arrays when programming for GPUs due to hardware limitations on sizes of shared memory buffers [4]. However, the metrics currently measured by AIWC do not have a direct causal relation to code patterns that optimize memory accesses on GPUs. The proposed relative local memory usage metric is the first to correspond to a recommended optimization strategy of using fast on-chip shared memory. Further, we find that all current AIWC metrics do not measure any sizeable difference between `coalescedAB`, `coalescedABT` and `alignedABT` codes. We address the root cause of this – a lack of temporal information about memory accesses, in the next section.

6 DISCUSSION

AIWC metrics are inherently lossy. Given the set of metrics calculated for an unknown kernel, we can not reconstruct the kernel based on these metrics. Specifically to memory access pattern analysis, the calculation of memory address entropy of a kernel relies only on the frequency distribution of memory accesses to all addresses accessed by the kernel. We observe that calculation of memory address entropy does not make use of temporal information. The order of sequential memory accesses performed by each thread, as well as addresses requested in parallel by multiple threads in an OpenCL thread group, are both vital in accurately characterising parallel codes.

We propose a new architecture independent metric to measure memory access patterns of parallel programs. The

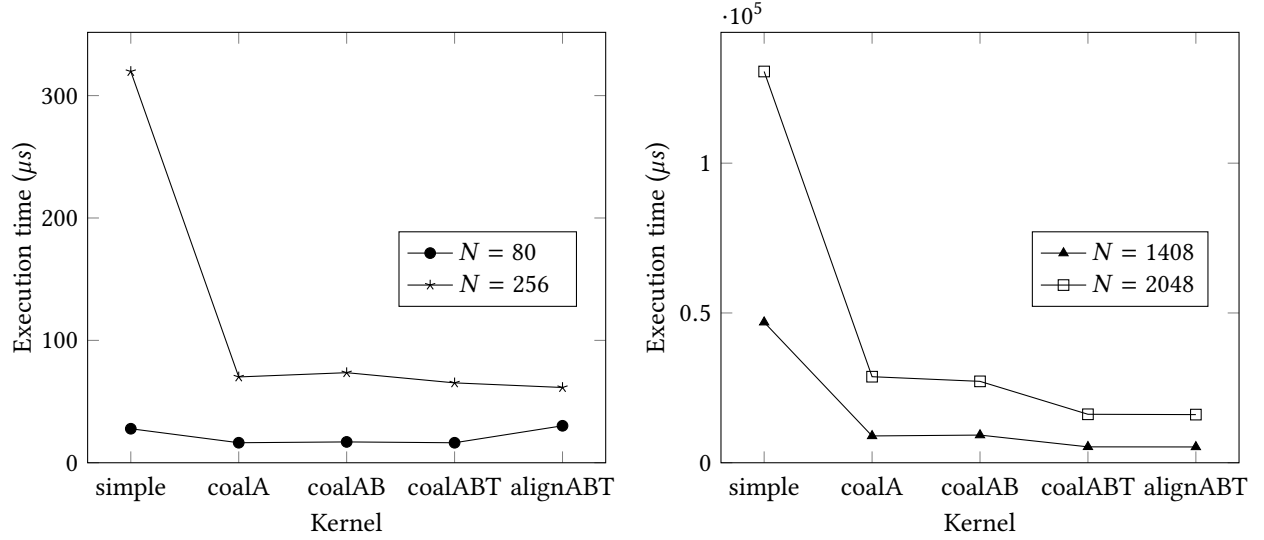


Figure 3: NVIDIA Tesla P100 Runtime performance of section 2 codes averaged over 100 runs

	simple	coalescedA	coalescedAB	coalescedABT	alignedABT
Total memory footprint	196608	196608	196608	196608	196608
90% Memory Footprint	118196	56176	489	489	489
Global MAE	17.02	13.18	9.78	9.78	9.78
LMAE #bits=3	16.02	12.18	8.78	8.78	8.78
LMAE #bits=10	9.02	5.18	1.78	1.78	1.78
Relative Local Memory Usage	0	49.91	94.28	94.28	94.28

Table 1: Selection of AIWC[8] metrics for 256×256 matrix multiplication

proposed metric is inspired by CuMAPz’ direct approaches to measuring optimization specific characteristics of CUDA codes.

AIWC currently collects a list of all memory accessed by each thread of execution. In the OpenCL programming model, threads within thread groups execute the same code and are allowed to cooperate on *local memory* addresses. We can group together threads by their work groups to align histories of memory accesses of threads in a group to analyse memory accesses performed by thread groups at each logical time step in the symbolic execution of the code. We argue that analysing the accesses occurring at the same logical time step is an accurate, conceptually sound method of predicting the extent of memory access coalescing and bank conflicts we are likely to find upon executing the code on a GPU.

The precise definition of the proposed metric is provided in the following section.

A Parallel Spatial Locality Metric

There are three steps involved in generating an AIWC metric: *recording*, *calculating* and *summarising* data collected from the symbolic execution of the kernel under inspection.

Record: we first record memory accesses performed by each thread in an OpenCL thread group as described above to achieve a global ordering of all memory accesses performed in a thread group. This ordering is collected in the form of logical timestamps ($t_0..t_{last}$) at which memory accesses occur.

Calculate: for each timestamp $t = t_0..t_{last}$, calculate the n -bits dropped entropy of memory addresses accessed by all threads in a work group within the timestamp t . Here n can range between 1 and 10 as was the case for *LMAE*.

Summarize: average the collected entropy values across all the timestamps to calculate the parallel spatial locality metrics for one thread group. We then average the n -bits dropped entropy summaries across thread groups to obtain

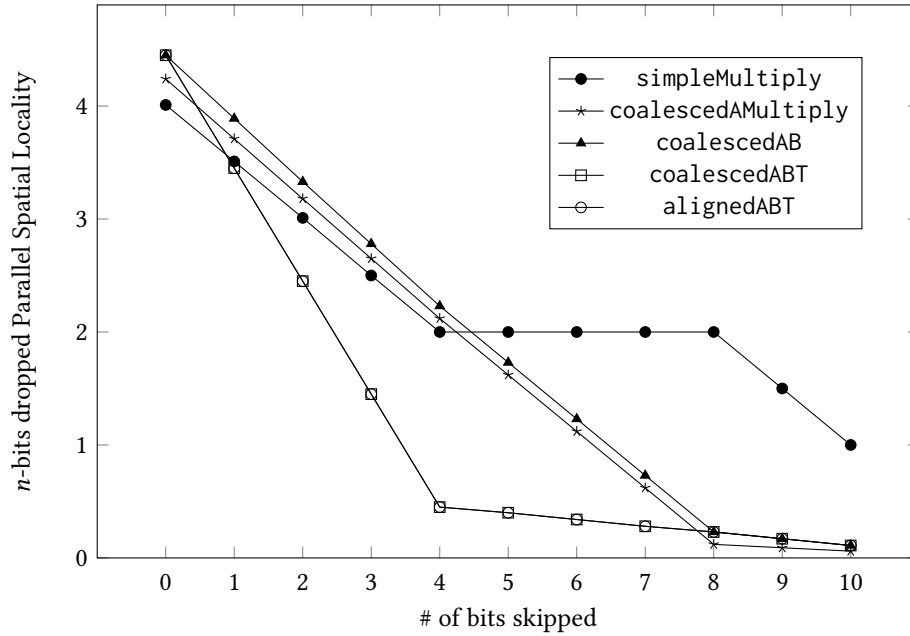


Figure 4: Calculated parallel spatial locality metric for matrix multiply kernels for 256×256 matrix multiplication

the *n-bits dropped parallel spatial locality metrics* for the kernel’s execution.

The proposed metric is a parallel computing analogue for MICA’s data stride metric that measures the distance between consecutive data accesses in a single threaded environment. In parallel programs, to accurately measure spatial locality of accesses, we must consider memory accesses performed by multiple threads in a close temporal scope. The proposed metric calculates the locality of accesses in each time step of the program’s execution and steeper drop-offs in *n-bits dropped parallel spatial locality* scores will be observed in programs that access nearby memory addresses within the same timestamp. Such programs will perform better on GPUs, as they will make better use of both global memory access coalescing and shared memory bank structures. To a lesser extent, the proposed metric captures performance-critical memory access patterns on CPUs, as pulling a single cacheline from global memory services the threads running on all CPU cores.

In figure 4 we see theoretical calculations for the proposed parallel spatial locality metric for each of the matrix multiplication kernels from section 2. In particular we observe that the coalescedABT and alignedABT kernels have the steepest drop-offs in entropy as the number of bits skipped is increased, which correlates to better locality of parallel memory accesses. It is expected for these kernels to exhibit better parallel spatial locality than simpleMultiply, as

coalescedABT and alignedABT make use of shared memory, where accesses tend to be localized simply due to the small size of shared on-chip memory typically available on GPUs. Further we find that the proposed metric successfully distinguishes between the coalescedAB and coalescedABT kernels. It accurately depicts a steeper drop-off for the more optimized coalescedABT kernel, where a larger proportion of parallel memory accesses make better use of the memory bank structure of GPU shared memory than all previous kernels. This is a significant improvement over the state-of-the-art AIWC metrics in characterising how codes localize shared memory accesses to better use the hardware provided.

7 RESULTS

The Extended OpenDwarfs benchmark suite is a set of OpenCL workload classified according to the computation patterns or Dwarfs. First presented in 2006, the Dwarfs categorize We present the results of running the AIWC tool on the Extended OpenDwarfs benchmark suite [1, 9].

N-body methods: GEM

The OpenDwarfs benchmark GEM computes the electrostatic potential of a biomolecule by calculating the sum of charges contributed by all atoms in the biomolecule at each specific surface vertex. This is an embarrassingly parallel problem. Each OpenCL work-item operates on a single surface vertex, finding the electrostatic potential generated by

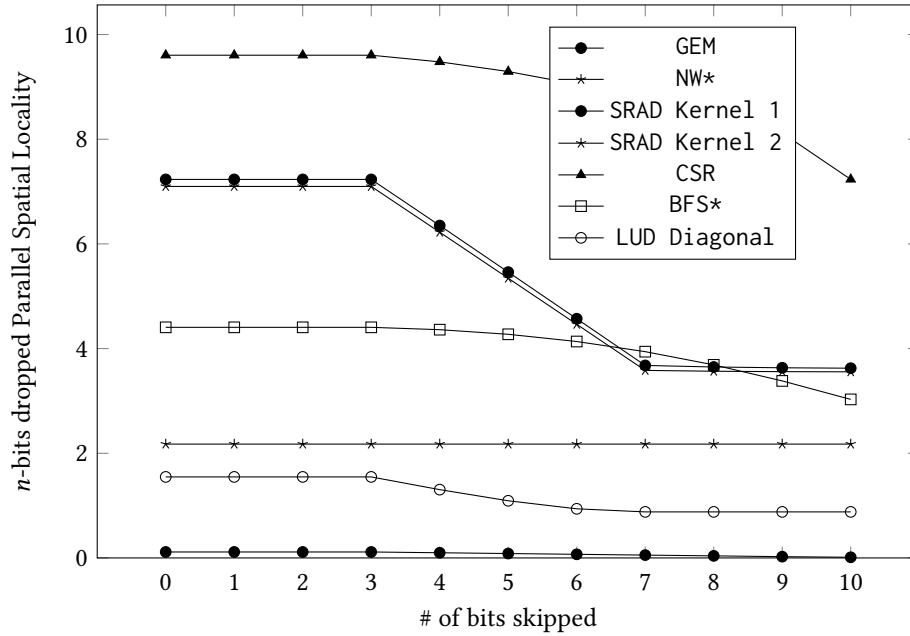


Figure 5: Parallel spatial locality metric for selected OpenDwarfs benchmark kernels

looping through every atom in the biomolecule in a global order. The computation pattern is highly regular and memory accesses are perfectly synchronized. Atom data is accessed consecutively, with all work-items simultaneously accessing each single atom's data. The parallel spatial locality metric reflects this pattern of efficient memory utilization (single loads from global memory servicing all OpenCL work-items). The recorded parallel spatial locality approaches the theoretical limit of 0 (0.0124 at 10 bits dropped). This indicates that almost all memory accesses made by the OpenCL kernel are perfectly synchronized between OpenCL threads.

Dense Linear Algebra: Lower-upper decomposition (LUD)

LUD in OpenDwarfs is a program to decompose an input $N \times N$ matrix as the product of one lower and one upper diagonal matrix. LUD is a Memory access patterns in a dense linear algebra workload such as LUD are typically highly regular and deterministic for each OpenCL work item, based on the the matrix dimension and offset parameters. The OpenDwarfs implementation of LUD [1] splits the LUD computation into three kernels: LUD perimeter and diagonal kernels spawn work-items in a single work dimension while LUD internal performs work in both dimensions of the input matrix.

LUD perimeter and diagonal partially benefit from memory access coalescing, for the lines of code where all the work-items access contiguous memory along a row of the

input matrix. However, when threads parallelly access a column of the input matrix, multiple memory requests are made as addresses accessed are too distant to be coalesced into a single memory transaction on GPUs. This is reflected in the slow drop in parallel spatial locality for LUD perimeter and diagonal kernels.

LUD internal

...

Dynamic Programming: Needleman-Wunsch (NW)

Needleman-Wunsch is a dynamic programming algorithm used to perform protein sequence alignment by identifying the similarity level between two strings of Amino acids. Elements in a matrix are computed, with each element dependent on its west, north and north-west neighbours. This dependency enforces a wavefront computation pattern, travelling along the main diagonal of the matrix. Each iteration of the kernels computes over an antidiagonal of the matrix, starting at its top left corner and finishing at its bottom right corner.

In particular, the wavefront pattern of computations causes each work-item to request distant memory addresses in a parallel fashion. Thus, parallel memory requests into the input matrix prohibit locality. We note the lack of any dropoff in parallel spatial locality as the number of bits dropped increases. This rightly indicates that memory addresses accessed at each logical timestamp are very distant. On GPUs this will translate to poor utilization of both memory access

coalescing and caching. This trend in the parallel spatial locality metric suggests that a possible improvement for GPU performance would be to load blocks of the input matrix into on-chip local memory to reduce the number of global memory requests.

Structured Grids: Speckle Reducing Anisotropic Diffusion (SRAD)

SRAD removes locally correlated noise from images by following a repeated grid update computation pattern. Conditional statements in the code cause thread divergence, potentially within a work-group, to handle boundary conditions. These boundaries constitute a small portion of the executed work-items, especially on larger data-sets and so the effective thread divergence is minimal. Thus, memory access patterns in SRAD are statically determined and relatively localized. The observed rapid drop-off in parallel spatial locality of the SRAD kernels (figure 5) is due to the highly localized nature of parallel memory accesses. At each point in the kernels' memory access profile, each work-item (i, j) in an OpenCL work-group accesses the $(i, j)^{th}$ elements of various matrices. The sequential memory access pattern is non-linear since different matrices are accessed consecutively by the kernel, prohibiting ideal caching. However, memory requests made parallelly are highly localized, allowing memory access coalescing on GPUs, greatly reducing the number of memory requests made on each line of the kernel code on GPUs. Thus while cache hit rates are typically low on this benchmark, particularly on GPUs [13], GPUs can hide the latency of global memory accesses through coalescing.

Graph Traversal: BFS

Graph traversal algorithms require pointer chasing operations to traverse nodes of a graph and perform calculations. The Breadth First Search implemented in OpenDwarfs calls OpenCL kernels to traverse nodes immediately connected to the list of nodes at each depth starting from the root node. BFS is characterized by an imbalanced workload per kernel launch depending on the degree of the nodes being operated on – with only a proportion of launched nodes performing meaningful work. As such the AIWC features collected for each kernel invocation vary significantly. The parallel spatial locality metric collects the entropy of memory accesses at each time step, and workload imbalances are dealt with by averaging entropy scores across all the time steps in the execution of an OpenCL work-group.

The memory access patterns in BFS are irregular. Work-items fetch discontinuous memory locations attributed to any particular node in the graph, depending on the connectivity of the node they operate on. The program structure

involves multiple levels of pointer-chasing and thus the precise memory address accessed are runtime dependent. This leads to poor parallel spatial locality. Thus, we see a slow dropoff in parallel spatial locality for BFS.

Sparse Linear Algebra: CSR

Compressed Sparse Row Matrix-Vector Multiplication (CSR) computes the sum of each of a matrix's row's elements after it is multiplied by a given vector. The matrix is stored in a compressed row storage sparse matrix format, which is very efficient for storage when the number of zero elements is far greater than the number of non-zero elements. Only non-zero elements of each row of a matrix are stored in Ax , along with separate arrays Aj , Ap indicating the position of each non-zero element in the matrix. Figure ?? depicts the sparse matrix format of storage for a 4×4 matrix.

In this particular implementation, a single OpenCL work-item is assigned to compute over each row of the input sparse matrix. Similar to BFS, memory access patterns are runtime-dependent, since the locations of matrix data read by each thread are dependent on the number and position of non-zero values in the sparse-matrix, decided by the values in Aj and Ap . This pattern of indirect addressing is typical of applications in the Sparse Linear Algebra Dwarf, which severely hinders locality of memory accesses performed. The parallel spatial locality metric reflects this trend. Figure 5 shows the very gradual decline in parallel spatial locality as the number of bits dropped is increased.

Finite State Machine: Temporal Data Mining (TDM)

8 CONCLUSIONS AND FUTURE WORK

To the best of our knowledge this work is the first to critically examine architecture independent methods to measure memory access patterns and guide hardware-specific optimization. We implemented and justified the addition of relative local memory usage to AIWC in our steps to characterize memory-based GPU code optimizations. We have also proposed the methodology to implement a new parallel spatial locality metric to AIWC. Upon implementation, our work would be the first to measure memory access patterns of parallel codes in an architecture-independent, easily interpreted manner.

Our proposed parallel spatial locality metric may also correlate to some memory-based optimization strategies for CPUs. Future work would apply the approach followed by this paper to CPU and FPGA architectures to evaluate the viability of an architecture independent characterization of memory access patterns in codes targeted to these architectures.

Future work will provide a concrete implementation of this new parallel spatial locality metric in the AIWC framework.

It is then feasible to test these new metrics across a larger set of GPU targeted OpenCL codes. This will allow a critical examination of the effectiveness of both proposed metrics in accurately profiling memory access patterns of these diverse codes.

REFERENCES

- [1] OpenDwarfs (base version). <https://github.com/vtsynergy/OpenDwarfs/commit/31c099aff5343e93ba9e8c3cd42bee5ec536aa93>, Feb 2017.
- [2] Intel Corporation. Intel 64 and ia-32 architectures optimization reference manual. 2016.
- [3] NVIDIA Corporation. Cuda c best practices guide. 2019.
- [4] NVIDIA Corporation. Cuda c programming guide. 2019.
- [5] Karthik Ganesan, Lizy John, Valentina Salapura, and James Sexton. A performance counter based workload characterization on Blue Gene/P. In *International Conference on Parallel Processing (ICPP)*, pages 330–337. IEEE, 2008.
- [6] John L. Hennessy and David A. Patterson. *Computer Architecture - A Quantitative Approach, 5th Edition*. Morgan Kaufmann, 2012.
- [7] Kenneth Hoste and Lieven Eeckhout. Microarchitecture-independent workload characterization. *IEEE Micro*, 27(3), 2007.
- [8] Beau Johnston and Josh Milthorpe. AIWC: opencl based architecture independent workload characterisation. *CoRR*, abs/1805.04207, 2018.
- [9] Beau Johnston and Josh Milthorpe. Dwarfs on accelerators: Enhancing opencl benchmarking for heterogeneous computing architectures. In *The 47th International Conference on Parallel Processing, ICPP 2018, Workshop Proceedings, Eugene, OR, USA, August 13-16, 2018*, pages 4:1–4:10. ACM, 2018.
- [10] Beau Johnston, James Price, Moritz Pflanzner, Petros Kalos, Tom Deakin, Nido Media, and Daniel Saier. BeauJoh/Oclgrind: Adding AIWC – An Architecture Independent Workload Characterisation Plugin. <https://doi.org/10.5281/zenodo.1134175>, December 2017.
- [11] Khronos OpenCL Working Group. *The OpenCL Specification, Version 1.1*, 2011.
- [12] Yooseong Kim and Aviral Shrivastava. Cumapz: a tool to analyze memory access patterns in CUDA. In Leon Stok, Nikil D. Dutt, and Soha Hassoun, editors, *Proceedings of the 48th Design Automation Conference, DAC 2011, San Diego, California, USA, June 5-10, 2011*, pages 128–133. ACM, 2011.
- [13] Konstantinos Krommydas, Wu-chun Feng, Christos D Antonopoulos, and Nikolaos Bellas. Opencdwarfs: Characterization of dwarf-based benchmarks on fixed and reconfigurable architectures. *Journal of Signal Processing Systems*, 85(3):373–392, 2016.
- [14] Chi-Keung Luk, Robert Cohn, Robert Muth, Harish Patil, Artur Klauser, Geoff Lowney, Steven Wallace, Vijay Janapa Reddi, and Kim Hazelwood. Pin: building customized program analysis tools with dynamic instrumentation. In *ACM SIGPLAN notices*, volume 40, pages 190–200. ACM, 2005.
- [15] Yakun Sophia Shao and David Brooks. Isa-independent workload characterization and its implications for specialized architectures. In *IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, pages 245–255. IEEE, 2013.
- [16] Kyle Spafford, Jeremy Meredith, and Jeffrey Vetter. Maestro: data orchestration and tuning for OpenCL devices. *Euro-Par 2010-Parallel Processing*, pages 275–286, 2010.

A SIMPLEMULTIPLY KERNEL

```

1  __kernel void simpleMultiply(__global float *A,
2      __global float *B,
3      __global float *C,
4      int N)
5  {
6      const int globalRow = get_global_id(0); //
          Row ID of C (0..N-1)
7      const int globalCol = get_global_id(1); //
          Col ID of C (0..N-1)
          // Compute a single element of C (loop over
          K)
8      float acc = 0.0f;
9      for (int k = 0; k < N; ++k) {
10         acc += B[k * N + globalCol] * A[globalRow
11             * N + k];
12     }
13     // Store the result
14     C[globalRow * N + globalCol] = acc;
15 }

```

B COALESCEDMULTIPLY KERNEL

```

1  __kernel void coalescedAMultiply(
2      const __global float* A,
3      const __global float* B,
4      __global float* C,
5      const int N)
6  {
7      __local float aTile[TILE_DIM][TILE_DIM];
8
9      const int localRow = get_local_id(0);
10     const int localCol = get_local_id(1);
11
12     const int globalRow = get_global_id(0);
13     const int globalCol = get_global_id(1);
14
15     __private float sum = 0.0f;
16
17     const int numTiles = N / TILE_DIM;
18     __private const int tiledRow = globalRow*N+
19         localCol;
20     for (int i = 0; i < numTiles; i++) {
21         aTile[localRow][localCol] =
22             A[tiledRow+i*TILE_DIM];
23         barrier(CLK_LOCAL_MEM_FENCE);
24         for (int k = 0; k < TILE_DIM; k++) {
25             sum += aTile[localRow][k] *
26                 B[(i*TILE_DIM+k)*N+globalCol];
27         }
28         barrier(CLK_LOCAL_MEM_FENCE);
29     }
30     C[globalRow*N+globalCol] = sum;
31 }

```

C COALESCEDAB KERNEL

```

1  __kernel void coalescedAB(
2      const __global float* A,
3      const __global float* B,
4      __global float* C,
5      const int N) {
6      __local float ASub[TILE_DIM][TILE_DIM];
7      __local float BSub[TILE_DIM][TILE_DIM];
8
9      const int localRow = get_local_id(0);
10     const int localCol = get_local_id(1);
11     const int globalRow = get_global_id(0);
12     const int globalCol = get_global_id(1);
13
14     float acc = 0.0f;
15     const int numTiles = N/TILE_DIM;
16
17     for (int i = 0; i < numTiles; i++) {
18         const int tiledRow =
19             globalRow*N+i*TILE_DIM + localCol;
20         const int tiledCol = globalCol +
21             (TILE_DIM*i + localRow)*N;
22         ASub[localRow][localCol] = A[tiledRow];
23         BSub[localRow][localCol] = B[tiledCol];
24
25         barrier(CLK_LOCAL_MEM_FENCE);
26
27         for (int k=0; k<TILE_DIM; k++) {
28             acc += ASub[localRow][k] *
29                 BSub[k][localCol];
30         }
31         barrier(CLK_LOCAL_MEM_FENCE);
32     }
33     C[globalRow*N + globalCol] = acc;
34 }

```
