

RESEARCH ARTICLE

Efficient and high-quality sparse graph coloring on GPUs

Xuhao Chen | Pingfan Li | Jianbin Fang | Tao Tang | Zhiying Wang | Canqun Yang

College of Computer, National University of Defense Technology, Changsha, 410073, China

Correspondence

Xuhao Chen, College of Computer, National University of Defense Technology, Changsha 410073, China,

Email: cxh@illinois.edu

Funding Information

National Natural Science Foundation of China (NSFC), Grant/Award Number: 61502514, 61602501, 61402488, 61502509; National Key Research and Development Program of China, Grant/Award Number: 2016YFB0200400

Summary

Graph coloring has been broadly used to discover concurrency in parallel computing. To speed up graph coloring for large-scale datasets, parallel algorithms have been proposed to leverage modern GPUs. Existing GPU implementations either have limited performance or yield unsatisfactory coloring quality (too many colors assigned). We present a work-efficient parallel graph coloring implementation on GPUs with good coloring quality. Our approach uses the *speculative greedy* scheme, which inherently yields better quality than the method of finding *maximal independent set*. To achieve high performance on GPUs, we refine the algorithm to leverage efficient operators and alleviate conflicts. We also incorporate common optimization techniques to further improve performance. Our method is evaluated with both synthetic and real-world sparse graphs on the NVIDIA GPU. Experimental results show that our proposed implementation achieves averaged $4.1 \times$ (up to $8.9 \times$) speedup over the serial implementation. It also outperforms the existing GPU implementation from the NVIDIA CUSPARSE library ($2.2 \times$ average speedup), while yielding much better coloring quality than CUSPARSE.

KEYWORDS

GPU, graph coloring, speculative greedy

1 | INTRODUCTION

Graph processing algorithms are getting a growing research interest in the past decade. They are pervasively used in many application domains, such as scientific computing, social networks, simulations and bioinformatics. Parallelizing graph algorithms is challenging because of their inherent irregularity. To leverage modern massively parallel processors, eg GPUs, make the problem even harder because of the difficulty of managing massive hardware resources and sophisticated memory hierarchies. In this paper, we investigate the problem of graph coloring, which assigns colors to all the vertices of a graph such that no neighboring vertices have the same color. Graph coloring is a fundamental graph algorithm that has been used in many applications¹⁻⁵ and is also intensively used by scientific computing to discover concurrency, eg, high performance conjugate gradient⁶ and incomplete-LU factorization,⁷ where coloring is used to identify subtasks that can be performed or data elements that can be updated simultaneously.

To deal with large-scale datasets, parallel graph coloring algorithms^{8,9} have been proposed to leverage the massive hardware resources on modern multicore CPUs or GPUs. Existing parallel implementations of graph coloring can be classified into 2 categories: (1) speculative greedy (SGR) scheme based¹⁰ and (2) maximal independent set (MIS) based¹¹. There are existing GPU implementations

of both categories. With different algorithms, they exhibit different characteristics of performance and coloring quality. The MIS implementations⁷ are usually fast because multiple threads can find MIS in parallel independently, and more importantly, they can substantially reduce the total number of memory accesses. But they inherently yield too many colors. On the other hand, SGR implementations¹² generally use fewer colors than MIS ones, but without careful mapping and optimizations, they spend much more time to complete coloring.

To overcome the limitations of existing approaches, we propose a high-performance GPU graph coloring implementation, which can produce high-quality coloring. Our method is built on the basis of the SGR scheme so that good coloring quality is guaranteed. It is then optimized specifically for the GPU architecture to improve performance. We choose data-driven instead of topology-driven mapping strategy for better work efficiency and make algorithm trade-offs to leverage efficient operators and alleviate the side effects of massive parallelism on GPUs. Meanwhile, we incorporate common optimization techniques, eg, kernel fusion, to further improve performance. The major insight of this work is that *algorithm-specific optimizations* are as important as common optimization techniques for high performance graph algorithm on GPUs. The main contributions of this paper are:

1. We present a work-efficient GPU graph coloring algorithm on the basis of the SGR scheme. The algorithm is carefully refined to better leverage GPU's bulk-synchronous model. It shows the importance of algorithm refinement to achieve high performance on GPUs.
2. We use optimization techniques specifically for the GPU architecture to take advantage of GPU's computation resources and memory hierarchies. Our practice further demonstrates GPU's capability on accelerating graph algorithms.

The rest of the paper is organized as follows: the existing serial and parallel algorithms as well as the state-of-the-art GPU implementations are introduced in Section 2. Our proposed design is presented in Section 3. We present the experimental results in Section 4. Section 5 discusses related work, and Section 6 concludes.

2 | BACKGROUND AND MOTIVATION

The graph coloring problem refers to the assignment of colors to elements (vertices or edges) of a graph subject to certain constraints. In this paper, we focus on *vertex coloring*, which assigns colors to vertices so that no 2 neighboring (connected) vertices are assigned the same color. There are several known applications of graph coloring, such as time-tabling and scheduling,¹⁻³ register allocation,⁴ high-dimensional nearest-neighbor search,⁵ sparse-matrix computation,^{6,7} and assigning frequencies to wireless access points¹³.

Graph coloring that minimizes the number of colors is an NP-complete problem and is known to be NP-hard even solved approximately.¹⁴ In this paper, we focus on *approximate graph coloring*, which yields near-optimal coloring quality. Many heuristics have been developed for approximate solutions, including first fit and largest degree first (LF). These heuristics make trade-offs between minimizing the number of colors and execution time, but generally faster algorithms have poor coloring quality while slower ones tend to yield fewer colors. In the following, we introduce some existing sequential and parallel algorithms.

2.1 | Sequential graph coloring

A sequential algorithm^{10,15} on the basis of the greedy scheme is shown in Algorithm 1. In all the algorithms specified in this paper, we use similar data structures to those introduced to the work of Çatalyürek et al.¹⁰ $adj(v)$ denotes the set of vertices adjacent to the vertex v , $color$ is a vertex-indexed array that stores the color of each vertex, and $colorMask$ is a color-indexed mask array used to mark the colors that are impermissible to a particular vertex v . At the beginning of the procedure, the array $color$ is initialized with each entry $color[w]$ set to zero to indicate that vertex w is not yet colored, and each entry of the array $colorMask$ is initialized with some value $a \notin V$. When processing the vertex v , the algorithm scans all its neighbors (line 3), and their colors are forbidden to be assigned to the vertex v (line 4). By the end of the inner for-loop, all of the colors that are impermissible to the vertex v are recorded in the array $colorMask$. It is then scanned from left to right to search the lowest positive index i at which a value different from the current vertex v is encountered; this index corresponds to the smallest permissible

color c to the vertex v (line 6). The color c is then assigned to the vertex v (line 7).

Algorithm 1 Sequential Greedy Algorithm¹⁰

```

1: procedure GREEDY( $G(V, E)$ )
2:   for each vertex  $v \in V$  do
3:     for each vertex  $w \in adj(v)$  do
4:        $colorMask[color[w]] \leftarrow v$ 
5:     end for
6:      $c \leftarrow \min \{i > 0 : colorMask[i] \neq v\}$ 
7:      $color[v] \leftarrow c$ 
8:   end for
9: end procedure

```

Algorithm 2 Parallel GM Algorithm¹⁰

```

1: procedure GM( $G(V, E)$ )
2:    $W \leftarrow V$  ▷ Initialize the worklist
3:   while  $W \neq \emptyset$  do
4:     for each vertex  $v \in W$  in parallel do
5:       for each vertex  $w \in adj(v)$  do
6:          $colorMask[color[w]] \leftarrow v$ 
7:       end for
8:        $c \leftarrow \min \{i > 0 : colorMask[i] \neq v\}$ 
9:        $color[v] \leftarrow c$ 
10:    end for
11:     $R \leftarrow \emptyset$  ▷ Initialize the remaining worklist
12:    for each vertex  $v \in V$  in parallel do
13:      for each vertex  $w \in adj(v)$  do
14:        if  $color[v] = color[w]$  and  $v < w$  then
15:           $R \leftarrow R \cup \{v\}$ 
16:        end if
17:      end for
18:    end for
19:     $W \leftarrow R$  ▷ Update the worklist
20:  end while
21: end procedure

```

2.2 | Parallel graph coloring

Parallel graph coloring has been applied to large-scale problems, such as sparse-matrix computation^{6,7} and chromatic scheduling³ to meet the performance requirement. Because of its sequential nature, the greedy scheme is challenging to parallelize. Basically, 2 classes of approaches have been proposed in the past to tackle this issue.

Gebremedhin and Manne (GM)⁹ used *speculation* to deal with the inherent sequentiality of the greedy scheme. It colors as many vertices as possible in parallel, tentatively tolerating potential conflicts, and resolve conflicts afterwards. Algorithm 2 shows the details of the GM algorithm. It can be divided into 2 parts: the first part (from lines 4 to 10) is the same as the sequential algorithm but done in parallel. The second part (from lines 12 to 18) does the conflict resolve (line 14) and puts the conflicting vertices into the remaining worklist (line 15). On the basis of this SGR algorithm, Çatalyürek et al developed OpenMP implementations for the multicore and massively multithreaded architectures.¹⁰ Rokos et al improved Çatalyürek algorithm and implemented it on the Intel Xeon Phi coprocessor.¹⁶

Another approach relies on iteratively finding an MIS of vertices in a progressively shrinking graph and coloring the vertices in the independent set in parallel. In many of the methods in this class, the independent set is computed in parallel using some variant of Luby algorithm.¹¹ An example is the work of Jones and Plassmann (JP).¹⁷ Algorithm 3 shows the details of the JP algorithm. Gjertsen et al¹⁸ introduced an advanced parallel heuristic, PLF, that consistently generates better colorings than the JP heuristic with slight overhead. Two new parallel color-balancing heuristics, PDR(k) and PLF(k) are also introduced. Hasenplaugh et al¹⁹ further improve the ordering heuristics on the basis of the JP algorithm.

2.3 | CUDA programming and GPU graph coloring

With the success of CUDA²⁰ programming model, general-purpose graphics processing units (GPGPUs)²¹ have been widely used for high performance computing (HPC) and many other application domains during the last decade.

In CUDA, individual functions executed on the GPU device are called *kernel* functions, written in a single program multiple-data form. Each instance of the single program multiple-data function is executed by a GPU *thread*. Groups of such threads, called *thread blocks*, are guaranteed to execute concurrently on the same streaming multiprocessors (SMs). Within each group, subgroups of threads called *warps* are executed in lockstep, evaluating 1 instruction for all threads in the warp at once. One of the major difficulties of CUDA programming is to manage the GPU memory hierarchy. It consists of register files, L1 memories (scratchpad, L1 cache, and read-only data cache), the shared L2 cache, and the off-chip GDDR DRAM.²²

Scratchpad memory (*shared memory* in CUDA terminology) is programmer visible and can be used for explicit intra thread block communication.

The L2 cache works as the central point of coherency and is shared across all threads of the entire kernel.

Algorithm 3 Parallel JP Algorithm⁷

```

1: procedure JP( $G(V, E)$ )
2:    $W \leftarrow V, c \leftarrow 1$ 
3:   while  $W \neq \emptyset$  do
4:      $S \leftarrow \emptyset$  ▷ Initialize the independent set
5:     for each vertex  $v \in W$  in parallel do
6:        $r(v) \leftarrow \text{random}()$ 
7:     end for
8:     for each vertex  $v \in W$  in parallel do
9:        $flag \leftarrow true$ 
10:      for each vertex  $w \in adj(v)$  do
11:        if  $r(v) \leq r(w)$  then
12:           $flag \leftarrow false$ 
13:        end if
14:      end for
15:      if  $flag = true$  then
16:         $S \leftarrow S \cup \{v\}$ 
17:      end if
18:    end for
19:    for each vertex  $v \in S$  in parallel do
20:       $color[v] \leftarrow c$  ▷ Color an independent set
21:    end for
22:     $W \leftarrow W - S, c \leftarrow c + 1$ 
23:  end while
24: end procedure

```

Several GPU graph coloring implementations have been proposed so far using either GM or JP algorithm. Grosset et al¹² implement the GM algorithm using CUDA. They use a 3-step graph coloring framework: (1) *graph partitioning*, which partitions the graph into subgraphs and identifies boundary vertices; (2) *graph coloring and conflicts detection*, which colors the graph using the specified heuristic, eg, first fit, and identifies color conflicts; and (3) *sequential conflicts resolution*, which goes back to CPU and resolves the conflicts. Note that step 2 is performed multiple times on GPU to reduce the number of conflicts before

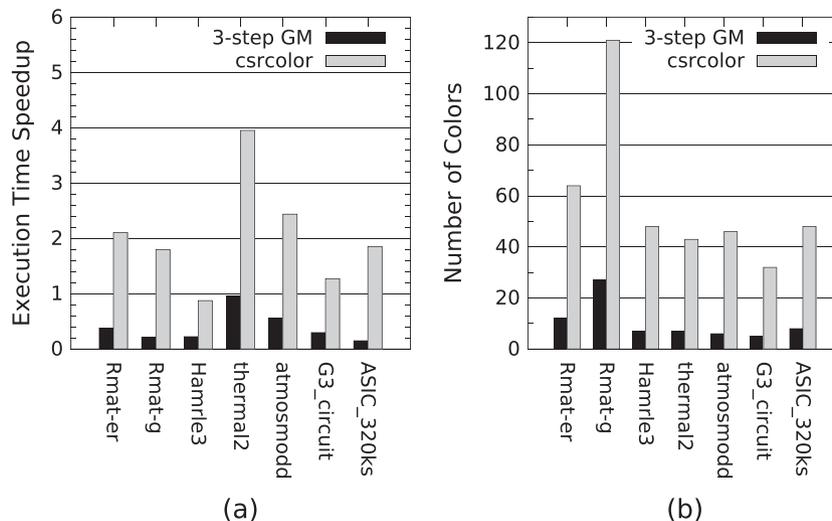


FIGURE 1 Comparison between 2 existing GPU graph coloring implementations: 3-step GM and csrcolor. A, Performance, ie, runtime speedup normalized to the serial implementation (the more the better); B, Coloring quality, ie, the number of colors assigned (the less the better). This figure shows that existing GPU implementations either have poor performance or yield unsatisfactory coloring quality, which motivates our work

going back to CPU. Although this 3-step GM algorithm assigns as few colors as the serial algorithm, its performance is poor, or even worse than the sequential graph coloring for many datasets, meaning the GPU computation horsepower is not leveraged very well.

The CUSPARSE²³ library offered by NVIDIA includes a `csrcolor`⁷ routine, which does graph coloring on a given graph in compressed sparse row (CSR) format.²⁴ The algorithm of `csrcolor` is derived from the JP algorithm, but uses the *multihash* method to find independent sets. Basically, several hash functions (instead of random number generators) are selected and used to generate hash values for each vertex with the vertex number as the input of the hash functions. Given the generated hash values, local maximum and minimum values can be found, and distinct (maximal) independent sets are generated for each of the hash values. Assume N hash values are associated with each vertex and used to create different pairs of (maximal) independent sets, this multihash method can generate $2N$ (maximal) independent sets at once. Compared to the GM algorithm, this method significantly reduces accesses to the *color* array, because it compares the generated hash values (in the registers) instead of the colors of neighbors (in the memory). As reported in the work of Naumov and Cohen,⁷ the `csrcolor` implementation runs pretty fast on modern NVIDIA GPUs. However, it usually produces several times more colors than the sequential algorithm, which is not satisfactory for many applications. For example, when applied to exploiting concurrency in parallel computing, more colors means less parallelism, because tasks (vertices) with the same color can be processed concurrently.

We evaluate the 2 existing GPU implementations of graph coloring on the NVIDIA K40c GPU. Figure 1 shows the performance and coloring quality of both implementations. As illustrated, 3-step GM yields much better coloring quality than `csrcolor`, but its performance is even worse than the sequential implementation, meaning it does not exploit GPU hardware very well. On the other hand, `csrcolor` runs much faster than 3-step GM and gains a certain degree of speedup over the sequential implementation. However, this good performance comes at the expense of much worse coloring quality: it yields several times more colors than the sequential implementation and 3-step GM. The limitations of `csrcolor` and 3-step GM motivate us to design a better implementation of parallel graph coloring for GPUs to achieve both high performance and good coloring quality.

3 | DESIGN

Graph algorithms are typical irregular algorithms²⁵ that are considered to be difficult to parallelize on GPUs. However, recent works²⁶⁻³¹ show that GPUs are capable to substantially accelerate graph algorithms if they are carefully designed and optimized for the GPU architecture. Although the previously proposed optimization techniques for other graph algorithms can be applied to graph coloring, we show that refining the algorithm for GPUs is essential for our case.

As mentioned in Algorithm 2, the graph coloring workload is composed of 2 major components: assign the first permissible color (Algorithm 4. `FirstFit`) and resolve conflicting vertices (Algorithm 5. `ConflictResolve`). The operations are trivial, but GPU's massively parallel model makes it challenging to efficiently parallelize these

workloads. We investigate the 2 activities in the following analyses using NVIDIA Tesla K40c GPUs.

Algorithm 4 FirstFit routine

```

1: function FIRSTFIT( $v$ )
2:   for each vertex  $w \in adj(v)$  do
3:      $colorMask[color[w]] \leftarrow v$ 
4:   end for
5:    $c \leftarrow \min \{i > 0 : colorMask[i] \neq v\}$ 
6:    $color[v] \leftarrow c$ 
7: end function

```

Note that we use the well-known CSR²⁴ sparse matrix format to store the graph in memory consisting of 2 arrays. Figure 2 provides a simple example. The column-indices array C is formed from the set of the adjacency lists concatenated into a single array of m (m is the number of edges) integers. The row-offsets R array contains $n + 1$ (n is the number of vertices) integers, and entry $R[i]$ is the index in C of the adjacency list of the vertex v_i .

Algorithm 5 ConflictResolve routine

```

1: function CONFLICTRESOLVE( $v$ )
2:   for each vertex  $w \in adj(v)$  do
3:     if  $color[v] = color[w]$  and  $v < w$  then
4:        $color[v] \leftarrow 0$ 
5:     end if
6:   end for
7: end function

```

3.1 | The baseline design

In the previous evaluation we find that SGR (ie, GM) algorithm inherently yields better coloring quality than the MIS (ie, JP) method. Thus, we choose to use the SGR scheme and design our baseline algorithm on top of it. Compared to the 3-step GM algorithm, our proposed GPU implementation maps the entire coloring work onto the GPU; consequently, removing the data transfer between the CPU and the GPU while the CPU is only responsible for controlling the progress. The rationale behind this change of mapping is that throughput-oriented processors are good at exploiting data-level parallelism, and thus, recomputing the conflicting vertices rather than serializing it onto the CPU would be more straightforward and efficient.

Nasre et al³² introduced the concept of *topology-driven* and *data-driven* implementations of irregular applications on GPUs. For graph algorithms, the topology-driven implementation simply maps each vertex to a thread, and in each iteration, the thread stays idle or is responsible to process the vertex depending on whether the corresponding vertex has been processed or not. The topology-driven implementation is straightforward, and since GPUs are suitable for accelerating data-parallel applications, it is easy to map onto the GPU hardware and possibly get speedup. By contrast, the data-driven implementation maintains a worklist that holds the remaining vertices to be processed. In each iteration, threads are created in proportion to the size of the worklist (ie, the number of vertices in the worklist). Each thread is responsible for processing a certain amount of vertices in the worklist, and no thread is idle. Therefore, the data-driven implementation is generally more work efficient than the topology-driven

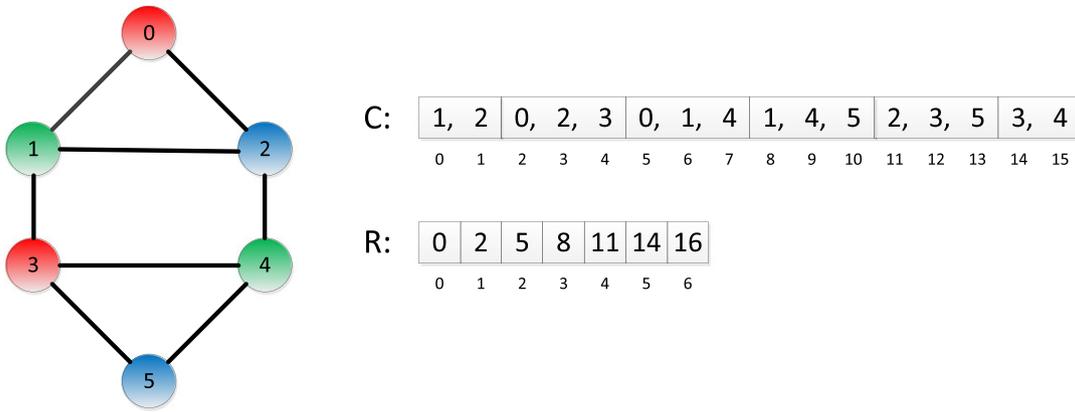


FIGURE 2 An example of the compressed sparse row (CSR) format. For this graph, at least 3 colors (red, green, and blue) are needed.

one, but it needs extra overhead to maintain the worklist. Note that the data-driven implementation still suffers from load imbalance problem, because vertices may have different amount of edges to be processed by the corresponding threads.

Algorithm 6 Topology-driven Parallel Graph Coloring

```

1: procedure TOPO-GC( $G(V, E)$ )
2:   do
3:      $changed \leftarrow false$ 
4:     for each vertex  $v \in V$  in parallel do
5:       if  $color[v] = 0$  then ▷ Not colored yet
6:         FIRSTFIT( $v$ )
7:          $changed \leftarrow true$ 
8:       end if
9:     end for
10:    for each vertex  $v \in V$  in parallel do
11:      if  $colored[v] = false$  then ▷ Not Colored yet
12:        CONFLICTRESOLVE( $v$ )
13:        if  $v$  is not conflicting then
14:           $colored[v] = true$ ;
15:        end if
16:      end if
17:    end for
18:    while  $changed = true$ 
19:  end procedure

```

We implement graph coloring in these 2 fashions. Algorithm 6 shows the topology-driven graph coloring algorithm. In this topology-driven algorithm, a flag *changed* is used to indicate whether all the vertices are colored or not. It is cleared at the beginning of each iteration, and set by 1 or more threads if any vertex is colored. Once all the vertices have been colored, the flag remains *false*, and the algorithm finally terminates. Both *FirstFit* and *ConflictResolve* are similar to those in the GM algorithm, but in *ConflictResolve* a bitmask *colored* is used to avoid recomputation. Algorithm 7 shows the data-driven graph coloring algorithm. It is almost the same as the GM algorithm except that Algorithm 7 uses *double buffering*³² to avoid copying the worklist. The 2 worklists W_{in} and W_{out} are referenced by pointers, and they are swapped at the end of each iteration. Since they are operated using pointers instead of data values, no copy operation is required between the 2 worklists.

Atomic operation reduction. In Algorithm 7, since the *out* worklist is a shared data structure, pushing elements into the worklist (line

11) requires atomic operations to ensure correctness. Although GPU architects have paid a lot of effort to optimize atomic operation, serialization from atomic synchronization is still expensive for GPUs.²⁶ Merrill et al²⁶ proposed to use software *prefix sum*^{33,34} for updating the shared worklist. Given a list of allocation requirements for each thread, prefix sum computes the offsets for where each thread should start writing its output elements. Fortunately, efficient GPU prefix sums³⁵ have been proposed, and the CUB³⁶ library has already provided standard routines for CUDA users to invoke. Thus, we need only 1 atomic operation for each block.

Algorithm 7 Data-driven Parallel Graph Coloring

```

1: procedure DATA-GC( $G(V, E)$ )
2:    $W_{in} \leftarrow V$  ▷ Initialize the in worklist
3:   while  $W_{in} \neq \emptyset$  do
4:     for each vertex  $v \in W_{in}$  in parallel do
5:       FIRSTFIT( $v$ )
6:     end for
7:      $W_{out} \leftarrow \emptyset$  ▷ Initialize the out worklist
8:     for each vertex  $v \in W_{in}$  in parallel do
9:       CONFLICTRESOLVE( $v$ )
10:      if  $v$  is conflicting then
11:         $W_{out} \leftarrow W_{out} \cup \{v\}$  ▷ Atomic push
12:      end if
13:    end for
14:     $swap(W_{in}, W_{out})$  ▷ Swap the worklists
15:  end while
16: end procedure

```

Color clearing. In Algorithm 7, when a vertex is determined to be conflicting, it is pushed into the worklist. Intuitively, its color should be cleared, and it will be assigned color in the next iteration. However, functionally, it is not a necessary operation. In the CPU parallel algorithm, this is not an issue. But for the GPU implementation, it is important to clear the color, so that when its neighbors check its color, there will be no conflicts. Thus, in Algorithm 5, the color is cleared (line 4). We observe nontrivial performance drop if the operation is removed. In the following sections we will see that the techniques to alleviate conflicts are performance critical to our GPU implementations.

Figure 3 compares their performance. As shown in the figure, the data-driven implementation outperforms the topology-driven one on

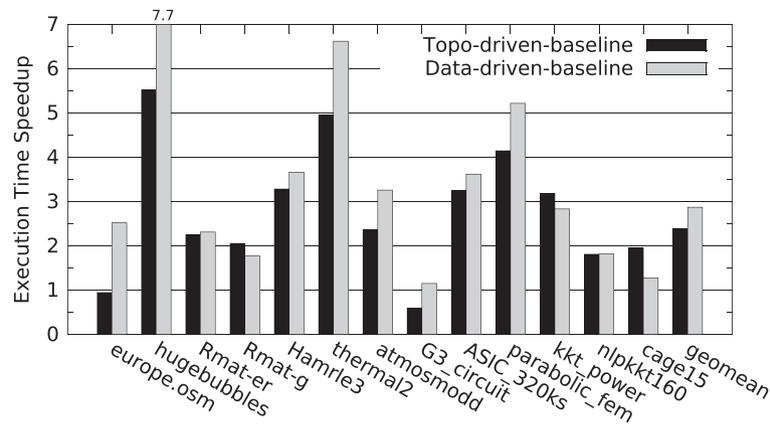


FIGURE 3 Runtime speedup of topology- and data-driven implementations, normalized to the sequential implementation

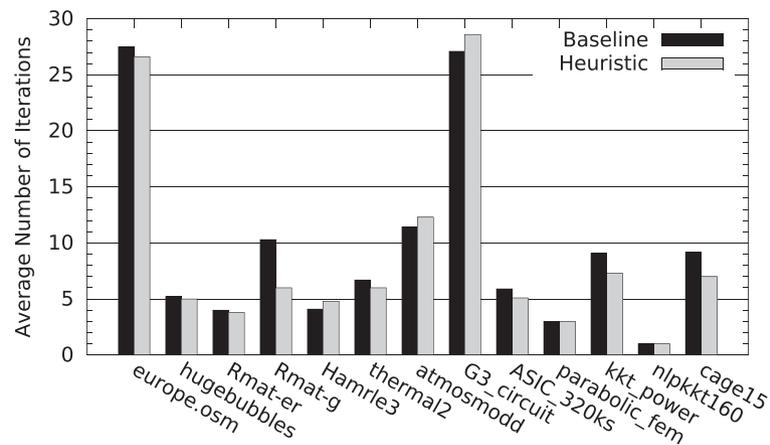


FIGURE 4 Average number of iterations with the baseline and heuristic implementations. Faster convergence leads to an average of 10.4% entire program speedup over the baseline

average, although the latter is more intuitive to implement on the GPU. This is easy to understand because the parallelism decreases in graph coloring as the iteration moves forward, and the topology-driven implementation has plenty of threads with no work to do, while the data-driven implementation is work-efficient although maintaining the worklist costs extra overhead. In the following discussion, we take this data-driven implementation as our baseline implementation. To achieve higher performance, we refine the algorithm to alleviate the side effects of massive parallelism and leverage efficient operators. We call them *algorithm-specific optimizations*. We also use common (nonalgorithm-specific) optimization techniques in Section 3.3.

3.2 | Algorithm-specific optimizations

As most parallel graph processing algorithms, parallel graph coloring is iterative. Therefore, it is important to ensure quick convergence for high performance. In the case of graph coloring, the number of iterations required to complete coloring highly depends on the conflict situation. For dense graphs, conflicts happen so frequently that no parallel algorithm can efficiently solve the problem. Our work thus focuses on sparse graph coloring, which is more common in real-world applications. Even so it is still challenging to parallelize it on GPUs, because the thousands of threads in the massively parallel program-

ming model make the conflicts happen much more frequently. This is not an issue on CPUs because there are only several or dozens of threads running simultaneously. We propose *heuristic conflict resolve* and use *thread coarsening* technique to alleviate this side effect of GPU parallelism.

3.2.1 | Heuristic conflict resolve

To reduce the number of iterations, an important part is to reduce conflicts. Since conflicts happen when 2 adjacent vertices are assigned the same color, deciding which of the 2 conflicting vertices to be reassigned in the next iteration affects the following conflict situation. An intuitive scheme is to pick the one with smaller or larger vertex id, but this is surely far away from optimal. We apply *heuristic conflict resolve* that prioritizes coloring the vertex with larger degree and puts the smaller one into the worklist to be processed in the next iteration. The rationale behind this heuristic is that vertices with larger degrees have more neighbors and thus are more likely to cause conflicts in the future. So it is better to color large-degree vertices first and reduce the possibility of conflicts. When the 2 vertices have the same degree, the one with smaller vertex id is picked. Figure 4 illustrates the average number of iterations required to complete coloring. It is shown that benchmarks, eg, *rmat-g* and *cage15* can have significant iteration reduction using

the heuristic. We also observe 43% and 50% execution time speedup of the entire program for the 2 benchmarks compared to the baseline. On average, the heuristic yields 10.3% speedup over the baseline.

3.2.2 | Thread coarsening

Thread coarsening is a common technique used in CUDA or OpenCL programs. It merges several threads together and thus have each thread do more work. This reduces the total number of threads and directly affects how data parallel work is mapped to the underlying hardware. Usually it is used to reduce the amount of redundant computation and thus can improve performance. For our case, however, it is used to reduce conflicts, because massive amount of threads on the GPU cause severe conflicts, which is not an issue on the CPU. Figure 5 shows the effect of thread coarsening applied to `FirstFit`, `ConflictResolve`, or both. Here we launch $nSM \times max_blocks$ thread blocks, where nSM is the number of SMs on the GPU and max_blocks is the maximum number of thread blocks that is allowed to be launched on each SM. max_blocks depends on how many resources (eg, registers and shared memory) a thread block allocates. Each block has 128 threads. Note that this is not the optimal configuration, which is different for different benchmarks, and some benchmarks, would be faster with even fewer thread blocks. Autotuning techniques would be helpful, but this is out

of the range of this paper. As shown, benchmarks, eg, `G3_circuit` and `cage15` can remarkably benefit from thread coarsening. On average, applying thread coarsening on both kernels can improve performance by 4.4% over the baseline.

3.2.3 | Bitset operation

Another major time-consuming part stems from writes and reads on the `colorMask` data structure in the `FirstFit` kernel. For each vertex, all its neighbors are visited to collect impermissible colors, which are written into `colorMask`. This information is then sequentially checked to find the first permissible color. In the worst case, all the elements in the `colorMask` array are checked, but actually we only need to find 1 permissible color. To reduce the costs of this operation, we propose to use `bitset` operations to implement reads and writes on the `colorMask` array. `bitset` is a standard class template in C++, but no similar support is provided in CUDA yet. Thus, we implement similar operations to mimic the functionality of the `bitset` class.

Fortunately, NVIDIA GPU architecture provides the `__ffs()` intrinsic for our use. Find first set (ffs) or find first one is a bit operation that identifies the least significant index or position of the bit set to one in the word. So our scheme is to initialize the bits as all "1"s and clear the bit if the corresponding color is impermissible. To

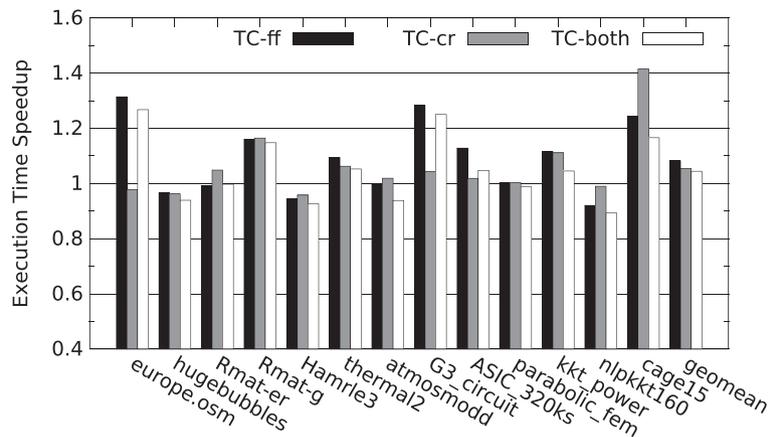


FIGURE 5 Program execution time speedup of thread coarsening on `FirstFit` (TC-ff), `ConflictResolve` (TC-cr), and both kernels (TC-both), all normalized to the baseline

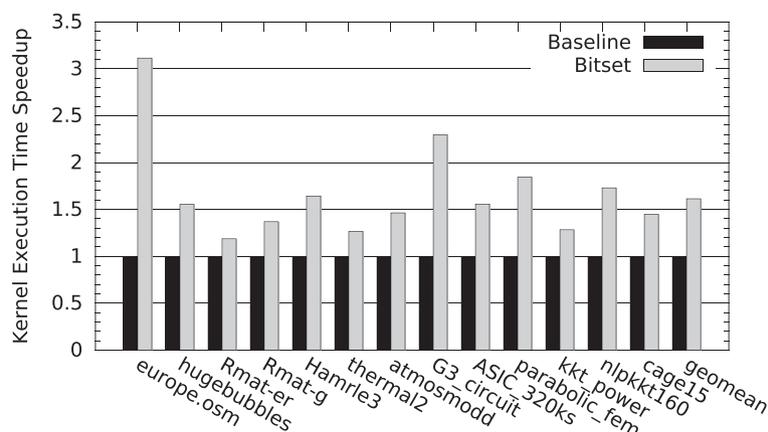


FIGURE 6 `FirstFit` kernel execution time speedup of `bitset` over the baseline. The kernel execution time is obtained by `nvprof`

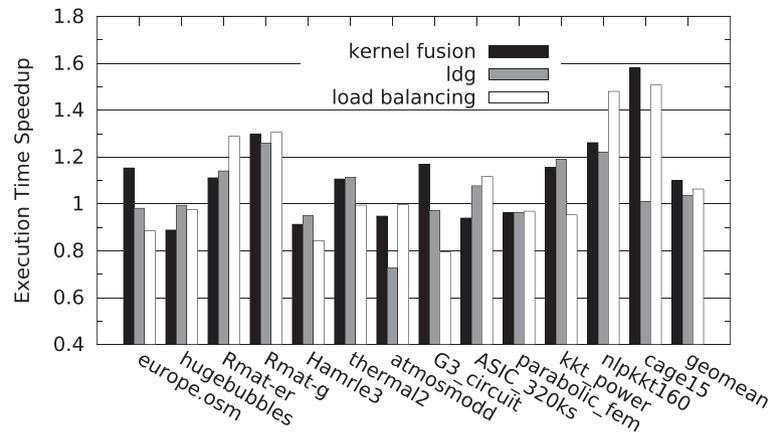


FIGURE 7 Program execution time speedup of *kernel fusion*, *ldg*, and *load balancing* over the baseline

find the first permissible color, we need only to call the `__ffs()` intrinsic. This implementation turns a for-loop into a single instruction and thus significantly reduces the operations required to complete the `FirstFit` kernel. Figure 6 shows a 61% speedup of the `FirstFit` kernel runtime over the baseline on average when `bitset` is applied. We also observe that this kernel improvement leads to an average of 28% speedup of the entire program compared to the baseline.

3.3 | Common optimization techniques

Existing GPU graph processing algorithms have already used many optimization techniques to improve performance. In graph coloring we use some of these optimizations, including *kernel fusion*, *read-only data caching*, and *load balancing* to enhance our implementation.

3.3.1 | Kernel fusion

Previous techniques focus on individual kernels. However, another important optimization technique called *kernel fusion* combines multiple GPU kernels into a single one and thus can keep the entire program on the GPU. Since adjacent kernels in CUDA share no state, this technique can leverage producer-consumer locality between operations and thus save significant memory bandwidth.³⁷ Note that global barrier is required between `FirstFit` and `ConflictResolve` operations. We use the existing method proposed by Xiao et al.³⁸ With this global barrier, kernels can only launch limited number of thread blocks; and thus, thread coarsening is forced to be applied to both kernels. Figure 7 shows an average 10% speedup of *kernel fusion* over the baseline. As shown, benchmarks, eg, `cage15` and `rmat-g` can benefit from better locality brought by kernel fusion because they are relatively denser and more irregular than others. We observe an improved L2 cache hit rate for `cage15`.

3.3.2 | Read-only data caching

In CUDA devices of compute capability 3.5 and higher, data that are read-only for the entire lifetime of the kernel can be kept in the read-only data (unified L1/texture) cache by reading it using the intrinsic `__ldg()`.²⁰ We use the texture cache to hold the read-only data, ie,

the C array and the R array. And then more read-only data are forced to be cached in the L1 read-only cache whose access latency is around 30 cycles, which is much shorter than the DRAM access latency (about 300 cycles). Therefore, `__ldg()` can capture temporal locality and improve the performance because of reduced DRAM accesses. As shown in Figure 7, `__ldg()` can bring 3.6% speedup over the baseline

3.3.3 | Load balancing

Another important issue for graph algorithms is load imbalance. The problem is particularly worse for scale-free (power-law) graphs. Merrill et al²⁶ proposed a hierarchical load balancing strategy that maps the workload of a single vertex to a thread, a warp, or a thread block, according to the size of its neighbor list. At the fine-grained level, all the neighbor list offsets in the same thread block are loaded into shared memory, then the threads in the block cooperatively process per-edge operations iteratively. At the coarse-grained level, per-block and per-warp schemes are used to handle the extreme cases: (1) neighbor lists larger than a thread block and (2) neighbor lists larger than a warp but smaller than a thread block respectively. We implement this strategy on graph coloring. Figure 7 illustrates the effect of load balancing on the benchmarks. Irregular benchmarks with uneven degree distribution, eg, `rmat-g` and `cage15` can substantially benefit from this technique. On average, it achieves 6.4% speedup over the baseline.

4 | EVALUATION

We use the R-MAT³⁹ graph generator to create synthetic graphs. The R-MAT algorithm determines the degree distribution by using 4 non-negative parameters (a; b; c; d) whose sum equals 1. We generated 2 graphs (`Rmat-er` and `Rmat-g`) with 1M vertices size but varying structures by using the following set of parameters: (0:25; 0:25; 0:25; 0:25); (0:45; 0:15; 0:15; 0:25). We also pick real-world sparse graphs from the University of Florida Sparse Matrix Collection.⁴⁰ These benchmarks are also used in previous works.^{7,26} The matrices with the respective number of vertices (ie, rows) and edges (nonzero elements) are shown in Table 1. The graphs vary widely in size, degree distribution, density of local subgraphs, and application domain.

TABLE 1 Suite of benchmark graphs

| Name | $n(10^6)$ | $m(10^6)$ | \bar{d} | σ | Description |
|---------------|-----------|-----------|-----------|----------|-----------------|
| europe.osm | 50.9 | 108.1 | 2.1 | 0.23 | Road network |
| hugebubbles | 21.2 | 63.6 | 3.0 | 0 | Adaptive mesh |
| rmat-er | 1.0 | 10.0 | 10.0 | 10.83 | Synthetic |
| rmat-g | 1.0 | 10.0 | 10.0 | 123.34 | Synthetic |
| Hamrle3 | 1.4 | 11.0 | 7.6 | 7.2 | Circuit sim. |
| thermal2 | 1.2 | 8.6 | 7.0 | 0.7 | Thermal sim. |
| atmosmodd | 1.3 | 8.8 | 6.9 | 0.1 | Atmosphere |
| G3_circuit | 1.6 | 7.7 | 4.8 | 0.4 | Circuit sim. |
| ASIC_320ks | 0.3 | 1.8 | 5.7 | 63.2 | Circuit sim. |
| parabolic_fem | 0.5 | 3.7 | 7.0 | 0.02 | General |
| kkt_power | 2.1 | 14.6 | 7.1 | 54.8 | Optimization |
| nlpkkt160 | 8.3 | 229.5 | 27.5 | 7.3 | Optimization |
| cage15 | 5.2 | 99.2 | 19.2 | 32.9 | Electrophoresis |

n means number of vertices; m means number of edges; \bar{d} means average degree; σ means degree variance.

4.1 | Experiment setup

We compare 6 implementations including (1) *Serial*, the serial implementation in CUSP¹⁵; (2) *OpenMP*, the baseline OpenMP implementation in the work of Çatalyürek et al¹⁰; (3) *3-step GM*, the previously proposed GM GPU implementation¹²; (4) *csrcolor*, the routine provided by NVIDIA CUSPARSE⁷; (5) *Proposed-base*, our proposed baseline data-driven implementation; and (6) *Proposed-opt*, our proposed optimized data-driven implementation. We conduct the experiments on the NVIDIA K40c GPU with CUDA Toolkit 7.5 release. *Serial* is executed on Intel Xeon E5-2690V2 2.30 GHz CPU with 12 cores. All the benchmarks are executed 10 times, and we collect the average execution time to avoid system noise. Timing is only performed on the computation part of each program. For all the GPU implementations, the input/output data transfer time (usually takes 10%-20% of the entire program execution time) is excluded because data is resident on the GPU in real applications⁷.

4.2 | Coloring quality

Figure 8 shows the number of colors needed by different implementations for each graph. It is not surprising that implementations except *csrcolor* need similar amount of colors, because they are all on the basis of the greedy scheme. The slight difference among these 5

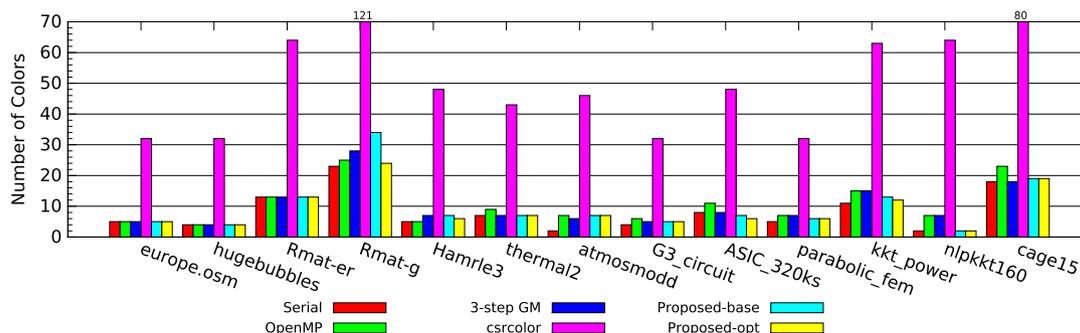
implementations may result from the different orderings that are caused by different thread mapping strategies and so on. *csrcolor*, however, needs $3.9 \times \sim 31 \times$ more colors than *Serial*, making this MIS-based implementation unattractive or even unapplicable in many scenarios. This substantial difference of coloring quality between *csrcolor* and another implementations stems from the inherent algorithm property of the SGR scheme and the MIS scheme. The SGR uses greedy scheme, and for parallel versions it optimistically does coloring in parallel with later conflict resolve. The MIS, however, tries to find independent sets iteratively, which does not cause any conflict, but for performance concern, the methods used to find independent sets should be simple enough and thus generate solutions that are far away from the optimal.

4.3 | Performance

Figure 9 illustrates the execution time speedup normalized to *Serial*. OpenMP on CPU achieves only moderate speedup ($1.54 \times$). As mentioned before, *3-step GM* gets unacceptable performance: 62% average slowdown compared to *Serial*. The slowdown stems from its mapping strategy and different data representation. In contrast, *csrcolor* is a much faster GPU implementation. It achieves an average speedup of $1.84 \times$ over *Serial*. For regular graphs, such as *hugebubbles* and *parabolic_fem*, it performs much better than OpenMP. This shows the high throughput and bandwidth advantages of GPUs over CPUs.

Our proposed baseline implementation performs even better than *csrcolor*. We observe $2.87 \times$ speedup on average over *Serial*. It is 85.8% and 56.1% faster than OpenMP and *csrcolor* respectively. For some benchmarks, eg, *Hamrle3* and *parabolic_fem*, *Proposed-base* significantly outperforms *csrcolor* ($4.18 \times$ and $2.46 \times$). This performance boost mainly comes from the selection of data-driven algorithm structure and the atomic operation reduction. However, for relatively dense or irregular benchmarks, eg, *cage15*, it performs worse than *csrcolor*, because no specific work is done to handle irregular cases, and *csrcolor* has fewer memory accesses as mentioned before.

With careful algorithm refinement and optimization techniques, we further improve the performance with an average speedup of $4.08 \times$ over *Serial*. It is $2.63 \times$, $2.21 \times$, and $1.42 \times$ speedup over OpenMP, *csrcolor*, and *Proposed-base*, respectively. Generally, for regular benchmarks, it takes advantage of GPU's high throughput as *csrcolor* does and performs even better because of the efficient

**FIGURE 8** Total number of colors assigned with different implementations

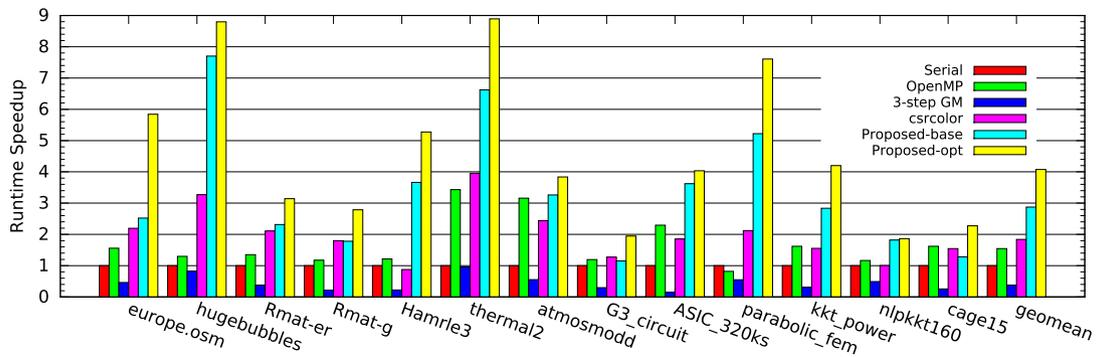


FIGURE 9 Runtime speedup normalized to the serial algorithm

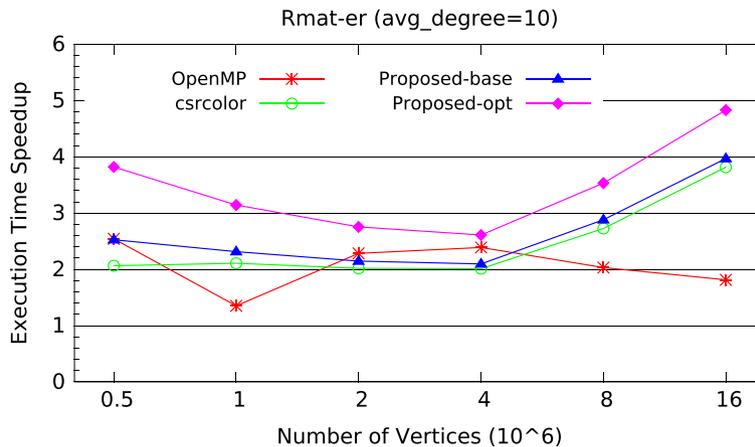


FIGURE 10 Execution time speedup of *Rmat-er* and *Rmat-g* with various graph size (number of vertices), all normalized to *Serial*

bitset operator, fast convergence, and so on, eg, *hugebubbles* ($8.8 \times$) and *thermal2* ($8.9 \times$). For irregular benchmarks, better locality and load balance lead to better performance. Thus, *Proposed-opt* can consistently outperform existing CPU and GPU parallel implementations.

We also notice that for some benchmarks, eg, *G3_circuit* and *nlpkkt160*, *Proposed-opt* gets very limited performance improvement compared to *Serial*. It is clear that the performance of graph coloring highly depends on the graph characteristics (scale, density, degree distribution, and topology). For example, *nlpkkt160* has a relatively large average degree and suffers from conflicts. And some are small in size, which limits the potential of performance improvement using GPUs. But more importantly, because the compute operation is trivial, the performance is likely to be limited by memory operations. For sparse graphs, not much temporal locality exists, and thus, the kernel becomes extremely memory bound with large-scale datasets, which could not be mitigated by the optimizations that we use. We suggest system software or hardware support for efficient memory access to overcome this performance bottleneck.

4.4 | Scalability

To evaluate the scalability of our design on the input size, we vary the graph size (number of vertices) of *Rmat-er* from 500K to 16M with fixed average degree ($\bar{d} = 10$). Figure 10 illustrates that *Proposed-opt* could achieve even more performance speedup given larger input datasets. Our proposed implementation can consistently

gain more than $2.5 \times$ speedup as the graph size changes and always outperforms *csr-color*. After 4M vertices, the speedup increases significantly as the graph size increases, while *OpenMP* changes moderately and even drops at the extremely large size. There is also a slight drop around 3M size for GPU implementations. This drop is related to the graph characteristics on which the performance highly depends on as mentioned. Here the cause is most likely the graph topology and degree distribution. Even so, *Proposed-opt* is still 9.3% faster than *OpenMP* at the 4M size, while it gets 2.66 speedup over *OpenMP* at the 16M size. For *Rmat-g* (not illustrated), we see a similar trend.

4.5 | Sensitivity to density

As mentioned, graph algorithms are highly sensitive to the characteristics of the input datasets. We evaluate sensitivity to the graph density of our proposed graph coloring implementation. In Figure 11, we vary the average degree \bar{d} of *Rmat-er* with fixed graph size (1M vertices). We compare *OpenMP*, *csr-color*, *Proposed-base*, and *Proposed-opt*, all normalized to *Serial*. As shown, *Proposed-opt* significantly outperforms the others when \bar{d} is small. This means our proposal can efficiently handle sparse graphs. However, as the average degree increases, the performance improvement over *Serial* decreases for *csr-color* and our proposals. Their curves drop below *OpenMP* when \bar{d} is larger than 20.

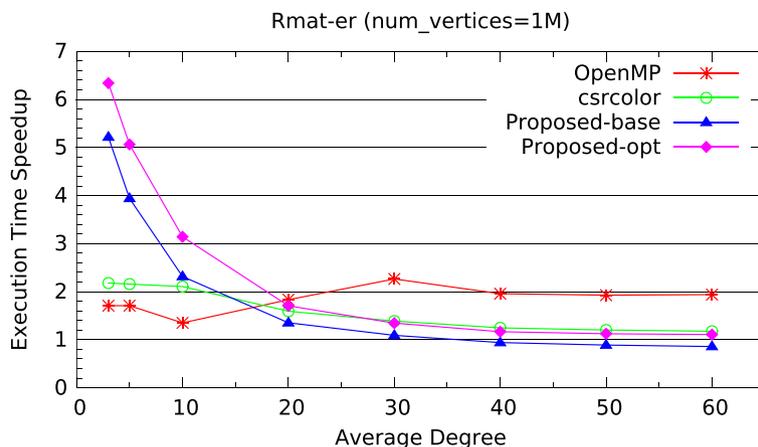


FIGURE 11 Execution time speedup of `Rmat-er` with various average degrees, all normalized to `Serial`

In contrast, `OpenMP` is more stable than GPU implementations. The drop of GPU ones results from the conflicts between neighbors, which is not an issue in CPUs. As the graph becomes denser, the conflicts happen more frequently. In this case, the GPU implementations need much more iterations to complete than `OpenMP`. For dense graphs, thanks to the techniques that alleviate conflicts, `Proposed-opt` still achieves comparable performance to `csrcolor`, while `Proposed-base` becomes worse than another 2 GPU ones and finally becomes slower than the serial implementation (blow 1). Remember that our proposals still consistently yield much better coloring quality than `csrcolor`. Although GPU implementations achieve high performance for sparse graphs, for dense graphs we suggest to use CPUs instead of GPUs to solve the graph coloring problem.

5 | RELATED WORK

Many graph algorithms have been developed on GPUs. Harish et al⁴¹ are the pioneers to implement GPU graph algorithms. They developed topology-driven breadth-first search (BFS) and shortest path algorithms. Hong et al⁴² proposed another topology-driven BFS to map warps rather than threads to vertices. Luo et al⁴³ developed the first work-efficient BFS on GPUs. Merrill et al²⁶ improved Luo's work. They used prefix sum to reduce atomic operations and used dynamic load balancing to deal with scale-free graphs. This implementation thus achieves high throughput and good scalability. The 2 major techniques of their work are also applicable to our implementation, while our work focuses more on the algorithm-specific refinement, eg, the specific strategies to alleviate side effects of GPU's massive parallelism.

Davidson et al³¹ developed a work-efficient single-source shortest path algorithm on the GPU. They used another load balancing strategy, which partitions the work into chunks and assigns each chunk to a block. Researchers also proposed GPU implementations of betweenness centrality,²⁷ minimum spanning tree,^{30,44} strongly connected components,²⁸ and so on. These work together demonstrated that with careful mapping and optimizations graph algorithms can get substantial performance boost on the GPU. Our work further enhances the conclusion of previous practices, while we show the importance of algorithm

refinement and architecture-specific optimizations for the problem of graph coloring.

Researchers have proposed many optimization techniques for graph algorithms, or more generally, for irregular algorithms on GPUs. The LAYER⁴⁵ is a locality-aware vertex scheduling scheme, which reorders the vertex queue to improve temporal locality of vertex data stored in on-chip caches. Nasre⁴⁶ proposed high-level methods to eliminate atomics in irregular programs, eg, BFS and single-source shortest path, on GPUs. Gunrock³⁷ absorbs previous knowledge and provides a library solution for GPU graph processing. It provides a load balancing framework on the basis of Merrill and Davidson strategies and integrates a set of common optimization techniques. A huge amount of efforts⁴⁷⁻⁵⁴ have been made by researchers to generalize graph processing computation and reduce programmer's burden.

Although generalized method can improve programmability, we argue that optimizations customized for the specific algorithm (which is difficult to generalize) is also important.

Che et al⁵⁵ characterize a suite of GPU graph applications and suggest architectural support. Xu et al⁵⁶ evaluate existing GPU graph algorithms on both a GPU simulator and a real GPU card and also suggest GPU hardware support. Wu et al⁵⁷ characterize 3 GPU graph frameworks and suggest to focus on constructing efficient operators. Beamer et al⁵⁸ also measure 3 graph libraries and propose processor architecture change. Green-Marl⁵⁹ is a domain specific language for graph processing. Chen et al⁶⁰ proposed compiler optimization methodology for graph and other irregular applications on Intel Xeon Phi coprocessors. Ahn et al⁶¹ developed a customized processing-in-memory (PIM) accelerator for large-scale graph processing. We believe that language, compiler, runtime, and architecture support is necessary for large-scale graph processing.

6 | CONCLUSION AND FUTURE WORK

Graph coloring is an important graph algorithm that has been applied in many application domains. To process large-scale graphs, parallel graph coloring has been intensively studied in the past. Meanwhile, GPUs have been broadly used to speedup compute intensive kernels of HPC applications in the past decade. In this paper, we explore

parallel graph coloring on the GPU. Existing implementations either achieve limited performance or yield unsatisfactory coloring quality. We present a high performance graph coloring implementation for GPUs with good coloring quality. We use the SGR scheme that guarantees coloring quality and improve performance with algorithm refinement and common optimization techniques. Experimental results show that our proposed implementation outperforms existing GPU implementations in both performance and coloring quality. This work helps us further understand graph algorithms on modern massively parallel processors and gives insight on the importance of both algorithm-specific and nonalgorithm-specific (common) optimizations. We also show the necessity of lower level support from system software and architecture.

In the future, we will further investigate the effect of conflict-resolution heuristics on performance and coloring quality and possibly propose even better heuristic. We will also try to implement our proposal on Intel Xeon Phi coprocessors and try to optimize it for the MIC architecture. Besides, it would be interesting to implement it on a GPU or MIC cluster to evaluate the scalability of our work.

ACKNOWLEDGMENT

We thank the anonymous reviewers for their insightful comments and suggestions, and A.V. Pascal Grosset from University of Utah for generously sharing his source code. This work is partly supported by the National Natural Science Foundation of China (NSFC) under grant numbers 61502514, 61602501, 61402488, and 61502509, and the National Key Research and Development Program of China under grant number 2016YFB0200400.

REFERENCES

- Welsh DJA, Powell MB. An upper bound for the chromatic number of a graph and its application to timetabling problems. *Comput J*. 1967;10(1):85–86.
- Lotfi V, Sarin S. A graph coloring algorithm for large scale scheduling problems. *Comput Oper Res*. January 1986;13(1):27–32.
- Kaler T, Hasenplaugh W, Scharld TB, Leiserson CE. Executing dynamic data-graph computations deterministically using chromatic scheduling. *Proceedings of the 26th ACM Symposium on Parallelism in Algorithms and Architectures, SPAA '14*. ACM, New York, NY, USA; 2014:154–165.
- Chaitin GJ. Register allocation & spilling via graph coloring. *Proceeding of the SIGPLAN Symposium on Compiler Construction*, Boston, Massachusetts; 1982 June:98–101.
- Berchtold S, Böhm C, Braunmüller B, Keim DA. Fast parallel similarity search in multimedia databases. *Proceedings of the acm sigmod international conference on management of data*, Tucson, Arizona; June 1997:1–12.
- Phillips E, Fatica M. A cuda implementation of the high performance conjugate gradient benchmark. In: Jarvis SA, Wright SA, Hammond SD, eds. *High Performance Computing Systems. Performance Modeling, Benchmarking, and Simulation*, Lecture Notes in Computer Science, vol. 8966: Springer International Publishing, Cham, Switzerland; 2015:68–84.
- Naumov PCM, Cohen J. Parallel graph coloring with applications to the incomplete-lu factorization on the GPU. Technical Report, NVIDIA Research, Santa Clara, CA; 2015.
- Allwright JR, Bordawekar R, Coddington PD, Dincer K, Martin CL. A comparison of parallel graph coloring algorithms. Technical Report, Syracuse University, Northeast Parallel Architecture Center, Syracuse, NY, USA; 1995.
- Gebremedhin AH, Manne F. Scalable parallel graph coloring algorithms. *Concurrency-Pract Ex*. 2000:1131–1146.
- Çatalyürek ÜV, Feo J, Gebremedhin AH, Halappanavar M, Pothan A. Graph coloring algorithms for multi-core and massively multi-threaded architectures. *Parallel Comput*. October 2012;38(10-11): 576–594.
- Luby M. A simple parallel algorithm for the maximal independent set problem. *Proceedings of the 17th Annual ACM Symposium on Theory of Computing, STOC '85*. ACM, New York, NY, USA; 1985:1–10.
- Grosset AVP, Zhu P, Liu S, Venkatasubramanian S, Hall M. Evaluating graph coloring on GPUs. *Proceedings of the 16th ACM Symposium on Principles and Practice of Parallel Programming, PPOPP '11*. ACM, New York, NY, USA; 2011:297–298.
- Riihijarvi J, Petrova M, Mahonen P. Frequency allocation for WLANs using graph colouring techniques. *Second Annual Conference on Wireless on-Demand Network Systems and Services*, St. Moritz, Switzerland; 2005 January:216–222.
- Zuckerman D. Linear degree extractors and the inapproximability of max clique and chromatic number. *Proceedings of the 38th Annual ACM Symposium on Theory of Computing, STOC '06*. ACM, New York, NY, USA; 2006:681–690.
- Dalton S, Bell N, Olson L, Garland M. Cusp: generic parallel algorithms for sparse matrix and graph computations. Version 0.5.0. Available: <http://cusplibrary.github.io/>; 2014. Accessed [September 2015].
- Rokos G, Gorman G, Kelly P. A fast and scalable graph coloring algorithm for multi-core and many-core architectures.. In: Triff JL, Hunold S, Versaci F, eds. *Euro-Par 2015: Parallel processing*, Lecture Notes in Computer Science, vol. 9233. Berlin Heidelberg: Springer; 2015:414–425.
- Jones MT, Plassmann PE. A parallel graph coloring heuristic. *SIAM J Sci Comput*. May 1993;14(3):654–669.
- Gjertsen RK, Jr., Jones MT, Plassmann PE. Parallel heuristics for improved, balanced graph colorings. *J Parallel Distr Com*. 1996;37:171–186.
- Hasenplaugh W, Kaler T, Scharld TB, Leiserson CE. Ordering heuristics for parallel graph coloring. *Proceedings of the 26th ACM Symposium on Parallelism in Algorithms and Architectures, SPAA '14*. ACM, New York, NY, USA; 2014:166–177.
- CUDA C Programming Guide v7.0. NVIDIA; March 2015.
- Keckler SW, Dally WJ, Khailany B, Garland M, Glasco D. GPUs and the future of parallel computing. *IEEE Micro*. 2011 Sept;31(5):7–17.
- Nvidia's next generation cuda™compute architecture: Kepler™gk110. NVIDIA; 2012.
- NVIDIA. CUSPARSE Library. Available: <http://docs.nvidia.com/cuda/cusparses/>; 2015. Accessed [September 2015].
- Dongarra J. Compressed row storage. Available: <http://web.eecs.utk.edu/dongarra/etemplates/node373.html>. Accessed [September 2015].
- Burtscher M, Nasre R, Pingali K. A quantitative study of irregular programs on GPUs. *Proceedings of the IEEE International Symposium on Workload Characterization, IISWC '12*, La Jolla, CA, USA; 2012 Nov:141–151.
- Merrill D, Garland M, Grimshaw A. Scalable GPU graph traversal. *Proceedings of the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP '12*. ACM, New York, NY, USA; 2012:117–128.
- McLaughlin A, Bader DA. Scalable and high performance betweenness centrality on the GPU. *Proceedings of the International Conference for High performance computing, networking, storage and analysis, SC '14*. IEEE Press, Piscataway, NJ, USA; 2014:572–583.
- Barnat J, Bauch P, Brim L, Ceska M. Computing strongly connected components in parallel on CUDA. *Proceedings of the 25th IEEE International Parallel Distributed Processing Symposium, IPDPS '11*, Anchorage (Alaska) USA; 2011 May:544–555.
- Nasre R, Burtscher M, Pingali K. Morph algorithms on GPUs. *Proceedings of the 18th ACM SIGPLAN Symposium on Principles and Practice*

- of *Parallel Programming*, PPOPP '13. ACM, New York, NY, USA; 2013:147–156.
30. Nobari S, Cao T-T, Karras P, Bressan S. Scalable parallel minimum spanning forest computation. *Proceedings of the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPOPP '12. ACM, New York, NY, USA; 2012:205–214.
 31. Davidson A, Baxter S, Garland M, Owens JD. Work-efficient parallel GPU methods for single-source shortest paths. *Proceedings of the IEEE 28th International Parallel and Distributed Processing Symposium*, Phoenix, AZ, USA; 2014 May:349–359.
 32. Nasre R, Burtcher M, Pingali K. Data-driven versus topology-driven irregular computations on GPUs. *Proceedings of the 27th IEEE International Parallel Distributed Processing Symposium*, IPDPS '13, Boston, Massachusetts, USA; 2013 May:463–474.
 33. Blelloch GE. Scans as primitive parallel operations. *IEEE T Comput*. 1989 Nov;38(11):1526–1538.
 34. Sengupta S, Harris M, Garland M. Efficient parallel scan algorithms for GPUs. Technical Report NVR-2008-003, NVIDIA, Santa Clara, CA; 2008.
 35. Yan S, Long G, Zhang Y. Streamscan: fast scan algorithms for GPUs without global barrier synchronization. *Proceedings of the 18th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPOPP-18. ACM, New York, NY, USA; 2013:229–238.
 36. Merrill D. CUB. NVIDIA Research. Available: <http://nvlabs.github.io/cub/>; 2015. Accessed [September 2015].
 37. Wang Y, Davidson A, Pan Y, Wu Y, Riffel A, Owens JD. Gunrock: a high-performance graph processing library on the GPU. *Proceedings of the 21st ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPOPP 2016, Barcelona, Spain; March 2016.
 38. Xiao S, Feng W. Inter-block gpu communication via fast barrier synchronization. *Proceedings of the IEEE 24th International Parallel and Distributed Processing Symposium*, Atlanta, GA, USA; 2010 May:1–12.
 39. Chakrabarti D, Zhan Y, Faloutsos C. R-MAT: a recursive model for graph mining. *SDM*. SIAM, Lake Buena Vista, Florida; 2004.
 40. The university of florida sparse matrix collection. Available: <http://www.cise.ufl.edu/research/sparse/matrices/>. Accessed [September 2015].
 41. Harish P, Narayanan PJ. In: Aluru S, Parashar M, Badrinath R, Prasanna VK, eds. *Proceedings of the 14th international conference high performance computing (HIPC)*, ch. Accelerating Large Graph Algorithms on the GPU Using CUDA. Berlin, Heidelberg: Springer Berlin Heidelberg; December 2007:197–208.
 42. Hong S, Kim SK, Oguntebi T, Olukotun K. Accelerating CUDA graph algorithms at maximum warp. *Proceedings of the 16th ACM Symposium on Principles and Practice of Parallel Programming*, PPOPP '11. ACM, New York, NY, USA; 2011:267–276.
 43. Luo L, Wong M, Hwu W-m. An effective gpu implementation of breadth-first search. *Proceedings of the 47th Design Automation Conference*, DAC '10. ACM, New York, NY, USA; 2010:52–55.
 44. Vineet V, Harish P, Patidar S, Narayanan PJ. Fast minimum spanning tree for large graphs on the GPU. *Proceedings of the Conference on High Performance Graphics*, ACM, New Orleans, Louisiana; 2009:167–171.
 45. Park H, Ahn J, Park E, Yoo S. Locality-aware vertex scheduling for GPU-based graph computation. *Proceedings of the IFIP/IEEE International Conference on Very Large Scale Integration (VLSI-SOC)*, Daejeon, Korea; 2015 Oct:195–200.
 46. Nasre R, Burtcher M, Pingali K. Atomic-free irregular computations on GPUs. *Proceedings of the 6th Workshop on General Purpose Processor using Graphics Processing Units*, GPGPU-6. ACM, New York, NY, USA; 2013:96–107.
 47. Malewicz G, Austern MH, Bik AJC, et al. Pregel: a system for large-scale graph processing. *Proceedings of the ACM SIGMOD International Conference on Management of Data*, SIGMOD '10. ACM, New York, NY, USA; 2010:135–146.
 48. Nguyen D, Lenharth A, Pingali K. A lightweight infrastructure for graph analytics. *Proceedings of the 24th ACM Symposium on Operating Systems Principles*, SOSP '13. ACM, New York, NY, USA; 2013:456–471.
 49. Low Y, Gonzalez J, Kyrola A, Bickson D, Guestrin C, Hellerstein JM. Graphlab: a new parallel framework for machine learning. *Proceedings of the UAI*, Catalina Island, California; 2010:340–349.
 50. Gonzalez JE, Low Y, Gu H, Bickson D, Guestrin C. Powergraph: distributed graph-parallel computation on natural graphs. *Proceedings of the 10th USENIX Conference on Operating Systems Design and Implementation*, OSDI'12. USENIX Association, Berkeley, CA, USA; 2012:17–30.
 51. Shun J, Blelloch GE. Ligma: a lightweight graph processing framework for shared memory. *Proceedings of the 18th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPOPP '13. ACM, New York, NY, USA; 2013:135–146.
 52. Zhong J, He B. Medusa: simplified graph processing on GPUs. *IEEE T Parall Distr*. 2014 June;25(6):1543–1552.
 53. Khorasani F, Vora K, Gupta R, Bhuyan LN. Cushman: vertex-centric graph processing on GPUs. *Proceedings of the 23rd International Symposium on High-Performance Parallel and Distributed Computing*, HPDC '14. ACM, New York, NY, USA; 2014:239–252.
 54. Fu Z, Personick M, Thompson B. Mapgraph: a high level api for fast development of high performance graph analytics on GPUs. *Proceedings of Workshop on Graph Data Management Experiences and Systems*, GRADES'14. ACM, New York, NY, USA; 2014:2:1–2:6.
 55. Che S, Beckmann BM, Reinhardt SK, Skadron K. Pannotia: understanding irregular GPGPU graph applications. *Proceedings of the IEEE International Symposium on Workload Characterization (IISWC)*, Portland, Oregon; 2013 Sept:185–195.
 56. Xu Q, Jeon H, Annavaram M. Graph processing on GPUs: Where are the bottlenecks?. *Proceedings of the IEEE International Symposium on Workload Characterization (IISWC)*, Raleigh, North Carolina, USA; 2014 Oct:140–149.
 57. Wu Y, Wang Y, Pan Y, Yang C, Owens JD. Performance characterization of high-level programming models for GPU graph analytics. *Proceedings of the IEEE International Symposium on Workload Characterization (IISWC)*, Atlanta, Georgia; 2015 Oct:66–75.
 58. Beamer S, Asanovic K, Patterson D. Locality exists in graph processing: workload characterization on an ivy bridge server. *Proceedings of the IEEE International Symposium on Workload Characterization (IISWC)*, Atlanta, Georgia; 2015 Oct:56–65.
 59. Hong S, Chafi H, Sedlar E, Olukotun K. Green-marl: a DSL for easy and efficient graph analysis. *Proceedings of the 17th International Conference on Architectural Support for Programming Languages and Operating Systems*, ASPLOS XVII. ACM, New York, NY, USA; 2012:349–362.
 60. Chen L, Jiang P, Agrawal G. Exploiting recent SIMD architectural advances for irregular applications. *Proceedings of the International Symposium on Code Generation and Optimization*, CGO 2016. ACM, New York, NY, USA; 2016:47–58.
 61. Ahn J, Hong S, Yoo S, Mutlu O, Choi K. A scalable processing-in-memory accelerator for parallel graph processing. *Proceedings of the 42nd Annual International Symposium on Computer Architecture*, ISCA '15. ACM, New York, NY, USA; 2015:105–117.

How to cite this article: Chen X, Li P, Fang J, Tang T, Wang Z, Yang C. Efficient and high-quality sparse graph coloring on GPUs. *Concurrency Computat: Pract Exper*. 2017;29:e4064. <https://doi.org/10.1002/cpe.4064>