# Atos: A Task-Parallel GPU Scheduler for Graph Analytics

Yuxin Chen\*
University of California, Davis
Davis, California, USA
yxxchen@ucdavis.edu

Aydın Buluç Lawrence Berkeley National Laboratory Berkeley, California, USA abuluc@lbl.gov Benjamin Brock University of California, Berkeley Berkeley, California, USA brock@cs.berkeley.edu

Katherine Yelick University of California, Berkeley Berkeley, California, USA yelick@berkeley.edu Serban Porumbescu University of California, Davis Davis, California, USA sdporumbescu@ucdavis.edu

John D. Owens University of California, Davis Davis, California, USA jowens@ece.ucdavis.edu

## **ABSTRACT**

We present Atos, a task-parallel GPU dynamic scheduling framework that is especially targeted at dynamic irregular applications. Compared to the dominant Bulk Synchronous Parallel (BSP) frameworks, Atos exposes additional concurrency by supporting task-parallel formulations of applications with relaxed dependencies, achieving higher GPU utilization, which is particularly significant for problems with concurrency bottlenecks. Atos also offers implicit task-parallel load balancing in addition to data-parallel load balancing, providing users the flexibility to balance between them to achieve optimal performance. Finally, Atos allows users to adapt to different use cases by controlling the kernel strategy and task-parallel granularity. We demonstrate that each of these controls is important in practice.

We evaluate and analyze the performance of Atos vs. BSP on three applications: breadth-first search, PageRank, and graph coloring. Atos implementations achieve geomean speedups of 3.44x, 2.1x, and 2.77x and peak speedups of 12.8x, 3.2x, and 9.08x across three case studies, compared to a state-of-the-art BSP GPU implementation. Beyond simply quantifying the speedup, we extensively analyze the reasons behind each speedup. This deeper understanding allows us to derive general guidelines for how to select the optimal Atos configuration for different applications. Finally, our analysis provides insights for future dynamic scheduling framework designs.

### **CCS CONCEPTS**

Computing methodologies → Massively parallel algorithms;
 Theory of computation → Parallel computing models;
 Software and its engineering → Software prototyping.

#### **KEYWORDS**

GPU, irregular workloads, task-parallel, asynchrony, speculation, graph algorithms



This work is licensed under a Creative Commons Attribution International 4.0 License.

ICPP '22, August 29-September 1, 2022, Bordeaux, France © 2022 Copyright held by the owner/author(s). ACM ISBN 978-1-4503-9733-9/22/08. https://doi.org/10.1145/3545008.3545056

#### **ACM Reference Format:**

Yuxin Chen, Benjamin Brock, Serban Porumbescu, Aydın Buluç, Katherine Yelick, and John D. Owens. 2022. Atos: A Task-Parallel GPU Scheduler for Graph Analytics. In 51st International Conference on Parallel Processing (ICPP '22), August 29-September 1, 2022, Bordeaux, France. ACM, New York, NY, USA, 11 pages. https://doi.org/10.1145/3545008.3545056

#### 1 INTRODUCTION

Bulk-synchronous parallel (BSP) programming [27] is the traditional model for GPU applications. It is a natural fit for statically schedulable, regular problems, such as many dense-matrix, image-analysis, and structured-grid computations. Programming environments like NVIDIA's CUDA and Khronos's SYCL support this relatively simple model, which maps efficiently to the massive fine-grained parallelism on GPUs and can deliver near-peak performance.

However, some important problems are instead irregular, with frequent control flow branches, non-unit-stride memory accesses, variable amounts of work across loop iterations, and dynamically varying degrees of parallelism. Algorithms that operate on graphs or trees or those with recursive formulations often exhibit such irregularity. These more naturally use a task-based programming model. Programming systems like Legion [2], PTask [22], and StarPU [1] use tasking on the CPU to feed GPUs with kernels to mask the latency of communication and keep the GPU busy. In contrast, we consider the problem of very fine-grained tasking where (traditionally) a set of similar application-level tasks are aggregated to form a data-parallel GPU task. This idea is used in state-of-the-art GPU graph libraries like Gunrock [28], where the application-level tasks are vertices or edges. In Gunrock and similar frameworks, each frontier in a graph sweep is launched as a separate GPU kernel in the BSP model. In practice, this may result in insufficient parallelism, uneven finish times, and high kernel launch overhead for small frontiers.

To address the above issues, we present of Atos, a task-scheduling framework for GPUs, that is adaptable to different usage scenarios:

- It supports both expensive and inexpensive frontiers by providing persistent and non-persistent task schedulers. The persistent scheduler is a GPU kernel that runs continuously to minimize launch overhead.
- It allows the user to trade off task and data parallelism by selecting the worker size, which is the number of GPU threads within each worker, and the number of items in each task.

It uses a single shared task queue, which balances load more through the shared task queue, which balances load more quickly than a distributed queue, yet is fast enough to keep quickly than a distributed queue, yet is fast enough to keep QU workers occupied.

GPU workers occupied.

It supports asynchronous execution across frontiers to maximize available parallelism, while mostly preserving crossimize, available parallelism, while mostly preserving crossimize available parallelism, while mostly preserving crossimize of the property of graph datasets a study three graph algorithms on a variety of graph datasets.

trontier ordering and thus minimizing overwork.

We study three graph algorithms on a variety of graph datasets that stress the importance of this adaptability. These algorithms that stress the importance of this adaptability. These algorithms have nested parallelism with outer-loop dependencies, and the nave nested parallelism with outer-loop dependencies; and the ability to relax those dependencies comes at the cost of possible overwork. Thus a second major theme in this paper is this tradeoff between increased parallelism and overwork. Each of our three algorithms explores this tradeoff in a somewhat different manner, which we discuss and analyze in Sections 5 and 6.

Our contributions include:

• Developing a generalized CDUL.

- Developing a generalized GPU task-parallel framework that explores a broad design space of possible task-parallel im-
- A demonstration of the benefits of mixed task and data parallelism for fine-grained parallel applications;
- Identifying relaxed-synchronization applications as a strong candidate for acceleration with a GPU task-parallel framework; and
- A detailed analysis of application performance that highlights the impact of design decisions both within the task scheduler and at the application level.

## 2 A DYNAMIC. IRREGULAR APPLICATION **PATTERN**

Atos handles a broad set of applications with fine-grained task and data parallelism, but we choose to focus here on a particularly challenging class of irregular nested loops with the following form:

Listing 1 A program with nested loops, expressed with a frontier abstraction.

```
in_frontier = initialize()
while (stop condition not met): // outer loop
               --CUDA-KERNEL--
  for (i in in_frontier): // inner loops
   for (j = 0 to workload(i).size()):
     out_frontier.append(f(in_frontier[i], j))
  cudaDeviceSynchronize()
  in_frontier = out_frontier
```

The inner loops produce data parallelism that may be flattened as The Nieger 1600s produced the parallelism that many his flattened an inement a class - Garager to 1820 the many microwant left or example, imilement a data-narallel load-balancing technique. For example, if the inner loops iterate overly rich evertises und outgoing edges some fire-analysis may be used to evenly distribute the edges rather. than the viertices. Our work also relaxes the outer loop iterations,

which will expand spraturities to find parallelism wing forms of dyname, irregular parallelism or more of the following forms of dynamic, irregular parallelism:

The number of tasks varies across outer loop frontiers: Work
 The number of tasks varies across outer loop frontiers: Work is generated dynamically and the number of output tasks is generated dynamically and the number of output tasks produced from each input task is not fixed.
 produced from each input task is not fixed.

- The cost of each inner loop task (lines 4–5 in Listing 1) may the cost of each inner loop task (lines 4–5 in Listing 1) may vary. The loop bound (work load 1) . \$12e() in line 5 is not pixed.

  In local work may vary: The outer loop has a loop-carried deford work may vary: The outer loop has a loop-carried dependence, so while dynamic generation of tasks will enforce some dependencies, parallel execution of tasks will enforce some dependencies, parallel execution of tasks across frontiers can change program behavior, including the total number of tasks processed.

The pattern in Listing 1 is not specific to graph algorithms. Atos co Il dealer traddrier Listing & of purble wifither glap malgorithms at them, inalldding raddressing awheref perblidteds thractled photolochaise partieted idechedingally techning helperolitidle or ithrofipated photolescare and extend dipelinically chlas Rayes while evil a sibject ageomedure dybarnic pipodėnoing, aluditkeytis evalgreritems sivitlata gienitequi omputatinio protessing PageRitekativehals of idented ithan significant in a protessing pageRitekativehals of identificant in the p pattern to PageRank, such as federated learning algorithms.

## Performance Challenges of Dynamic,

# Feregular Problemsenges of Dynamic,

Traditional RSPa implementations of such applications usually

2.1 Persylhance Chillienges of Dynamic,

Traditional girls protogrations of such applications usually launch a series of kernels, with each kernel corresponding to a fractional BSF implementations of such applications traully faunch a series of kernels, with each kernel corresponding to a bulk-synchronous step lone iteration of the outer loop), between a series of kernels, with each kernel corresponding to abulk-synchronous each, step is a global synchronization burner. The kernels them step lone iteration of the outer loop), between each step is a global synchronization burner. The kernels them step lone iterations for the outer loop), between each step is a global series of the outer loop. The later with the step is a global step of the control of the outer loop. The later is the step of the later is the ste

## **OUR TASK-PARALLEL PROGRAMMING MODEL AND DESIGN SPACE**

Our programming model allows implementations of task-parallel formulations of the workloads we discussed in Section 2, with a focus on providing solutions to BSP's three challenges: smallfrontier, load-imbalance, and loss of concurrency opportunities. We use the following terminology:

worker: one or a group of GPU threads.

task: one or more pieces of work that are scheduled as a single umit im our system.

**application function**  $f(\cdot)$ : the code that processes each task.

Listing 2 SPMD code of each thread worker in Atos.

```
-----CUDA KERNEL-----
for each worker:
  while not queue.empty():
   task = queue.concurrent_pop(task.size())
   new_tasks = f(task)
   queue.concurrent_push(new_tasks)
```

At a high level, our model maintains a queue of tasks. Work-At a high level, our model maintains a queue of tasks. Work-ers fetch a task from the queue, process the task, and add newly ers fetch a task from the queue. Process the task, and add newly generated tasks (if any) to the queue. The program runs until e-generated tasks (if any) to the queue. The program runs until re-ther a stop condition is met or the queue is empty (Listing 2). This programming model addresses the performance challenges from section 2.1: Section 2.1:

We can implement Listing 2 with a single kernel invocation.
We can implement Listing 2 with a single kernel invocation, avoiding multiple launches of small kernels.
avoiding multiple launches of small kernels.
Task parallelism is implicitly load-balanced because workers task parallelism is implicitly load-balanced because workers can run independent fasks and stay busy even if tasks require different amounts of work.
Listing 2 has no global synchronization; instead the program—Listing 2 has no global synchronization; instead the programmer controls the scheduling of work, allowing more flexible mer controls the scheduling of work, allowing more flexible dependencies and thus more opportunities for concurrency.

Managing task dependencies In Atos's programming model; tasks themselves (run inside a worker) are executed synchronously, but different tasks (across workers) are executed asynchronously: For this execution model to work; tasks are only added to the task queue when their dependencies are satisfied. Our current implemenfation of Atos supports tree-structured task dependency graphs; which is sufficient for the graph applications considered in this paper (and we believe most graph algorithms). Atos can be extended in a straightforward way to BAGs by adding (atomic) counters for each join; the last worker to reach the join would continue the computation beyond the join. This strategy does put a burden on the programmer; but the overall advantage of this approach is to support task-based computations that are generated dynamically (in contrast to Legion [2], StarPU [1], Juggler [4], etc. that must build a static task-dependency graph in advance):

The primary focus of this paper is how to best implement Listing 2 on the GPU: We identify below the four key design decisions in such an implementation.

**Relaxing barriers** A BSP implementation separates work in each iteration with a global barrier. Can we benefit from relaxing this global barrier constraint?

Worker size We can choose what GPU resources we assign to each worker. What is the worker size that yields the best performance?

Data vs. Task Parallelism We expect to leverage the parallelism between tasks. We can also choose the size of our tasks, and within each task, potentially leverage data parallelism. What is the right balance between task and data parallelism?

Kernel strategy Listing 2 is written in a "persistent" style, implementable with a single kernel call. We could alternatively interchange its outer and inner loops, making one "discrete" kernel call per iteration. Which is best for optimal performance?

## 3.1 Relaxing Barriers

As discussed in Section 2, many applications have the nested loop structure from Listing 1 and their BSP implementations may lose concurrency opportunities because of the global barrier between the outer loops. In many cases, we can remove the barrier and relax the outer loop dependency while still computing the correct result. How? Consider two tasks A and B that, in a BSP implementation, are ordered: A is in an iteration that precedes B and thus must run before B.

- One possibility is to speculate that we can compute A and B at the same time, or even in the order B then A, without changing the correctness of the computation. If our speculation is correct, then we expose more concurrency. If our speculation is incorrect, then we must fix it. This fix might be cheap or costly.
- · A second possibility is a problem formulation that is robust to computing items out of order. This is also called Dijkstra's don't care non-determinism [11].

In either case, we can relax global barriers and expose additional concurrency; depending on the problem and dataset, we may see significant performance gains. However, relaxing barriers may incur additional costs: the cost of performing incorrectly speculated work, the cost of repairing incorrectly speculated work, and less predictable convergence rates when compared to the BSP counterpart. Overall, if the performance improvements from increased concurrency outweigh these costs, we can deliver performance

Related work: Hassaan et al. [16] studied unordered and ordered versions of several algorithms, demonstrating a tradeoff between parallelism and work efficiency. However, the relaxed barrier formulations we study differ from unordered ones. Consider breadth-first search (BFS). Both Hassaan et al. and we begin with a work-efficient Dijkstra BFS, but they compare to a work-inefficient Bellman-Ford BFS, while we compare to a relaxed (speculative) Dijkstra BFS (Section 5.1). The speculative Dijkstra BFS is more work-efficient than Bellman-Ford BFS. Empirically, speculative Dijikstra's workload is within a small constant factor of that of BSP Dijkstra, which is #edges (see Table 4). This is much smaller than Bellman-Ford's workload of diameter × #edges. Kulkarni et al. [19] studied an abstraction and runtime scheme for workloads with optimistic parallelism, which differ from the relaxing barriers we study in this paper. Their notion of optimistic parallelism assumes many

tasks can run in parallel and stops a task the moment it violates a dependency. In contrast, we allow the computation to commit, even if it violates a dependency, and only fix the mistake afterwards.

## 3.2 Worker Size

Previous task-parallelism work [1, 7, 8, 22] uses the whole GPU as a single worker. They maintain a task dependency graph on the CPU side and orchestrate the execution by launching a CUDA kernel for each task when dependencies are satisfied. We use *GPU-wide-worker* to describe such an organization. Tasks in those GPU-wide-worker task-parallel frameworks are usually very large to better utilize the entire GPU's resources. This organization is easy to program and has low scheduling overhead.

However, the GPU-wide-worker scheme is a poor match when task dependencies require finer management. Consider the following extreme example: Task A and Task B both contain 10,000 data items, and only a single data item in B depends on a single item in A. A GPU-wide-worker implementation must wait for A to complete before beginning work on B, even though most of the data items can be processed independently and concurrently. We use the term false dependency to describe the situation when a data item has all its dependencies satisfied, but cannot be processed because another data item in the same task has unresolved dependencies. One can reduce such false dependencies by decomposing a large task into many smaller tasks to expose more parallelism. As a result, workers should be smaller, matching the size of tasks and allowing full utilization of the GPU's resources. This approach has motivated a number of recent task-parallelism frameworks [4, 23, 26, 29], which use workers sized as either warps or Cooperative Thread Arrays (CTAs). The resulting additional complexity in scheduling many smaller workers motivates also moving scheduling decisions from the CPU to the GPU.

Most task-parallelism frameworks only provide one worker size. Our framework provides thread-, warp-, and CTA-sized workers, to support tasks of different size and different synchronization requirements. The only prior work that uses multiple granularities is Whippletree [23]. Whippletree's thread and warp worker sizes are primarily a programming model concept and suffer from synchronization penalties at the implementation level. In Whippletree's implementation, threads are still synchronized within entire CTAs, suffering from false dependencies if tasks require finer synchronization than at CTA granularity.

For graph analytics frameworks in particular, data-parallel bulk-synchronous execution models are by far the most common on GPUs because of their high GPU utilization and effective use of data-parallel load-balancing techniques (e.g. Gunrock [28], cuGraph [12], Medusa [30], SIMD-X [20], GraphBLAST [10]). Current multi-GPU task-based asynchronous graph libraries—Groute [5], Lux [18], and Galois [17]—use a data-based bulk-synchronous model for the computation kernels launched on each GPU.

## 3.3 Balance Between Data and Task Parallelism

Many parallel applications have work items that require different amounts of processing. Traditional BSP applications address this challenge with explicitly coded data-parallel load balancing techniques. We describe two different approaches in the context of Listing 1: One widely used technique is load balancing search [9], which dynamically computes the prefix-sum of workload(i).size() for  $i \in in_f$ rontier, then flattens the two for-loops into one big array and redistributes the work in the array to each CUDA thread (see Baxter [3] for details). Another popular data-parallel load balancing technique separates the work in in\_frontier into different buckets based on workload(i).size() and launches a separate kernel with the best processing strategy for that size for each bucket [21].

Task parallelism is a natural fit for these irregular applications. Workers in our framework do not directly synchronize with each other. One of the primary advantages of this lack of coordination is that it allows workers to attend to any available work items as soon as they become available ("implicit task-parallel load balancing").

Task parallelism and data parallelism are not exclusiveindividual tasks of sufficient size may also exploit data parallelism in their execution. Thus we consider a continuous spectrum with pure task-parallel and pure data-parallel load balancing at the extremes, and expect that the optimal trade-off will be application-dependent. Our framework supports two worker sizes larger than a thread (warp and CTA) and offers the programmer the ability to exploit data parallelism within each warp-sized or CTA-sized task. In the framework, workers operate on tasks asynchronously, but an individual worker itself is executed synchronously. Therefore, we can use a worker's capacity as a parameter to control the tradeoff between data and task parallelism. Given a fixed number of threads available, increasing a worker's capacity reduces the total number of workers available for a given application. At the same time, an increase in worker capacity results in more opportunities to perform data-parallel load balancing within each worker. In the extreme, setting a worker's capacity to the entire GPU leaves no room for task parallelism and is equivalent to the BSP model. We found that data-parallel load balancing inside the worker combined with task-parallel load balancing provided by Atos results in better overall load balancing (Section 6). We are not aware of any previous work that combines these two types of load balancing.

## 3.4 Kernel Strategy

Traditional GPU kernels divide a variable amount of input work into fixed-size CTAs and launch a kernel over a CTA count proportional to the amount of input work. Persistent kernels [14] decouple the relationship between data size and launched CTAs. A persistent kernel launches only enough CTAs to fill the GPU. These CTAs remain resident for the entire kernel and run a loop that maps naturally to the task-parallelism model in Listing 2.

Advantages of persistent kernels Persistent kernels reduce kernel launch overhead and CPU/GPU communication. This is particularly significant when many small kernels are required. The persistent kernel approach reduces CPU involvement in favor of programmer-written GPU logic within the persistent kernel.

**Disadvantages of persistent kernels** GPU workers in the persistent kernel concurrently pop from a shared queue; this requires atomic operations to ensure exclusive pops. Persistent kernels have higher register usage than discrete kernels (requiring extra registers to maintain the queue loop).

Intuitively, if a discrete-kernel application suffers from large leaded by the large large leaded by the large large

```
Listing 3 Atos framework APIs
4 FRAMEWORK API
template<typename T, typename COUNTER_T>
struct Queues {
Listing 3 Avids framework APIsity, int num_queues, int iteration);

template<typename T, typename COUNTER_T>
struct Queues {
Listing 3 Avids framework APIsity, int num_queues, int iteration);

template<typename Fibename Fibename Counter T, typename ... Args>
structhostusevoid launchThread (bool ifPersist, int numBlock,
intoowmThread, int typename Fibename Fibename ... Args>
__nost__ void launchThreadopedoifTPersiststintneumGhedock,
int numThread, int shareMem_size, Ffiff1, F2ff2, AAgs...aPgs);

template<typename I, typename Fibename Fibename Rfgstypename... Args>
__host__ void launchWaFable (bool ifPersist, intrumBlock,
int numThread, int shareMem_size, Ffiff1, F2ff2, AAgs...aBg);
}

template<int FETCH_SIZE, typename Fi, typename F2, typename... Args>
__host__ void launchCTA (bool ifPersist, int numBlock,
With the book of the color of the content of the color of th
```

With the discussion from Section 3 in mind, we introduce the Atos API shown in Listing 3. Launch API functions are used to launch workers who repeatedly pop tasks from the work queue; each it bather discussioned rancos of influencial propriet, which the watare Amilia kayan in Liatis maruhan 12 ki ah Muhatin mareku watar the her interest at the property of the proper the hworkers han applied timetical that cahe task appears. When the the contributed on the unique precision \$2 defined and prime and model the permitted him where the new torone maximal included where it the antique number of the sas that eact on section by which we epecure pour que, aux en presentantes de la company de la HISTORIAN CONTROL STATEMENT OF THE CONTROL OF THE C defined how varied than a verse constitute and it is steep as well now a MITCHER THE PROPERTY OF THE PR (decrese) on eminimizate de la managementante de la conservata de la conse the succession worker his principal thread, the FETCH\_SIZE should be not incondicate lact on the manisasing and the size in elevantes discressive Bry with wars is blocked for other yarkers but in-ton To Allesticated the super of the determination of the control inf speculation of the property of the propert Appending and opportunition of the description of t THE WASTE PASSITES WE GOAL SET WHICH I SETT THE ASSETT AND ASSISTED TO SET OF and hand and are depth values. If the depth of a neighbor is improved period of interest in the queue. Lastly, we pass BFSWarp() and its arguments onto launchwarp and invoke it. In this piperwe Asset Shreeldlasse nested loop problems in the ndomni professi besilen generation desperations destructed by Weishaus at haw because of perapire lementations are well-retudied, in the BSE model of the party o so we can be confident that our results are meaningful. Also, when run on particular datasets, their BSP implementations exhibit one or more of the challenges described in Section 2.1.

run on particular datasets, their BSP implementations exhibit one dristing of the chance get aged RES (worker size warp) ordeaux, France struct BFS (

int num\_nodes; Lighting Paros-based relaxed BFS (worker size: warp)

```
QHeunam<u>w</u>noblèsists;
        int num_edges;
@ER(Gegr;csr, int capacity, int num_queues) {
          in€<sup>sç</sup>depêhçsr;
        Ouwoeklinstalimits(capacity, num_queues);
                 cudaMalloc(&depth, sizeof(int) * num_nodes);
         BFS(Csr _csr, int capacity, int num_queues) {
                                           & csr
         voidrRī§$ŧartWat@6iptcquŋŖlaGk_qietegyṃThread) {
                \label{eq:workhastac} \textbf{workhastac} \\ \textbf{aueshWarefizeoffunRigok,nownThomesd}, \ \textit{0}, \ \text{BFSWarp(),} \ \\ \textbf{*this);}
}:
void BFSStartWarp(int numBlock, int numThread) {
templaFRIรี$ሂደዊ፻፭ሢጸርሃጸዊ፻፭ሢቪ ቫ교ሐBኒዕርጲሃጾቤዓሕፕዋና፩፭ቒቸው; ፅፑቴ</mark>ਔαrp(), *this);
class BFSWarp {
                                               _ void operator()(VertexId node, BFS bfs) {
          __device_
tempYqttexidpdqahd vebtexdqpth[nqdelypename SizeT = int>
classiprtwapdetoffset = bfs.csr.get_neighbor_list_start(node);
publiceT neighborlen = bfs.csr.get_neighbor_list_length(node)
                 19146 | Heighbor Ten — ประการและ Lengthon III Length(Nobe),
เพื่อใช้เ<mark>ล้า ใช้ยนี อุดะฟิโตัก-โร้จุนตร์เปลาสู่ปลู่ปลูกตร์เปลารูว (tem + 32) (
VeYteโซล์ไปลูกซ์เลียง</mark>
                 Vertextd-depthsours: depth, hoost; lies in the miles it is a significant of the miles in the miles it is a significant of the miles in 
                         VertexId old_depth = atomicMin(bfs.depth + neighbor, depth + 1);
if (old_depth > depth + 1) bfs.worklists.push_warp(neighbor);
                         syncwarp():
```

#### 5.1 Breadth-First Search

The ps Brienarth of 1881 (Search elaxed-barrier BFS ("Speculative BFS") are chostning by the state of barrier BFS ("Speculative BFS") are chostning by the state of barrier BFS ("Speculative BFS") are chostning by the state of barrier BFS ("Speculative BFS") are chostning by the state of barrier BFS ("Speculative BFS") are chostning by the state of barrier by the chostning by th

Page Rage Ranks the importance (rank) of nodes in a graph. Page Rank bear must be amportance (rank) of nodes in a graph. Page Rank bear must be importance from lating the BSF and page Rank bear age Rank has mellioned by the BSF and has completely a personal and another than 1900 to the major and has completely a personal and a property of the page Rank has a pull of the page Rank has a pull

CEP 22 August 22 September 1, 2022 Bordeaux France

Front i.e. 10 its neighbors. The second kernel augustates all vertical and the property of the property of

```
A BOTTON OF THE STATE OF THE ST
```

```
Algorithm 1 Bylke Synchronous BFS
                                                                                                                                 The state of the s
                                                                                                                                         To the second se
```

The above three graph algorithms show different approaches to the above three graph algorithms show different approaches to the above three graph algorithms show different approaches to the above three graph algorithms show different approaches to the above three graphs are the above three graphs are the above the

The above three graph algorithms show different approaches to managing asynchrony that can all be addressed within the Atos

```
Algorithm 2 30 derilative BFS
                                                    ### Le not from the community of the com
                                                                                                                                                       frontier.append(neighbor)
       Algorithm 3 Bulk Synchronous PageRank
Algorithm 3 Bulk Synchronous PageRank
Algorithm 3 Bulk Synchronous PageRank
           Algerithm 3 Bulks Synchimmus PageRank
                                                                                                         PWHTANGHESIANL WELLVEST 1 ambda / lambda vertex neighborLen)

Find Island (September 1) ambda / lambda vertex neighborLen)

Find Island (September 1) ambda / lambda vertex neighborLen)

Find Island (September 1) ambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lambda vertex neighborLen)

Find (September 1) ambda / lambda / lam
                                                                               Algorithm 4 Asynchronous PageRank
Algorithm 4 Asynchronous PageRank
Algorithm 4 Asynchronous PageRank
Alegitan Language and Alegitan Alegitan
```

residue[check\_id%G.total\_vertices] The above three graph algorithms who we different approaches to managing asynchrony that can all be addressed within the Atos managing asynchrony that can all be addressed within the Atos framework: (1) PageBank: a robust asynchronous algorithm: (2) RESERVENCE OF A STATE OF A STA

THE STATE OF THE S (version 11.1.168) with the -O3 flag and gcc 9.3.0 with the -O3 flag. All results ignore transfer time and are averaged over 20 runs.

Atos: A Task-Parallel GPU Scheduler for Graph Analytics Atos: A Task-Parallel GPU Scheduler for Graph Analytics

#### Algorithm 5 Bulk Synchronous speculative Graph Coloring Algorithm 5 Bulk Synchronous speculative Graph Coloring

# Algorithm 6 Asynchronous speculative Graph Coloring Algorithm 6 Asynchronous speculative Graph Coloring

```
#$$\text{if First is the Color of the Color
```

6.1 Experimental Overview

**6.1** Experimental Overview
We evaluate the design principles discussed in Section 3 using three Was ratuate the design principles discussed in Section Ausing threeimplementation variants based on a rombination of Atos configud rations: (1) "persist-32" artilizen persistentikarnels with warp-sized. workers. It has no data parallel load balancing within (2) workerinstead only using our birit task-parallel load balancing (2) "persistworker siker FETCH SIZE" willizes persistent kernels with STAsized-workers. (3)1"disarsterworker-sizebFETCHASIZEantilizes Idisa creterkernels and CTA-sized workers Both CTA variage to use load) balancing search [21] (and ata-parallel load to balancing technique)inside workers in conjunction with implicit task-parallel load balancing. BFS and PageRank, we compare the performance of our impl Fore RFS and Page Ranks we sompare the performance of purpopplaneantation from Gwarack 128h a state of the sart single a GPIU RSP abased graph framework owbich extensively users datasparallel loadbalancing techniques | For Graph Goloring Guttock's independenty set grant reloring algorithm is not comparable son we faithfully implemented a BSP graph or oldring using the same speculatives greedy graph coloring algorithm a Duri BSP implementation duses, Gurrank's bucket hased data-parallel load balancing method [28], described in Section 3.3.

Table 1: Runtime With 25 peeting Swimparing 260 BSP uniferthe property of the saids beeting Swimparing to BSP Graphe property at the property of the property

a	<u> </u>	ceran			
	Dataset	BSP.	persist warp	persist CTA	discrete CTA
	soc-LiveJournal1 <sup>S</sup>	Applica	tion; BFS (BSP=0	344(x) 23)	10.7 (x1.42)
	hpllyawood_2009s	A <b>pply</b> eat	ions aris (BST)	umparakit.(81)A	4di60000000TA
	indochina_2004 <sup>S</sup>	13.2	15.6 (x0.84)	8.03 (x1,65)	7.42 (x1.79)
	roll Tive Fournal 18	BSP3	PETSIST WASP	P629434 ( 13)	174 (38482)
	500 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	<b>58:3</b> 6	2 <del>2/</del> 6 <sup>2</sup> (\$X(\$3)	49:43 (x12:28)	18.50((3258)
	HOROCHURA-20098Ar	plication	: Pate Produktis	=63203(toicluss))	4.56 <sup>2</sup> (X2.02)
-	Dataserna 2004s	BS92	profesion (No.	persist (cops)	alse te to
-	roadNetnea soe Live ournal1s	202.9	938 (x1.58) <sup>1)</sup>	4637((21378 <sup>3</sup> )	116 (x2.358)
		pp <b>sta</b> tio	n: <b>80.6e(Rands)</b> BS		15.5 (x3.56)
	indochina_2004 <sup>8</sup> Ap	pll <b>iggj</b> or	: Partish WBSF	=Germer CTA	4discrete OTA
	road_usa'''	221	169 (x1.30)	121 (X1.81)	112 (X1.95)
	Data set etallinal1s	PSP P	P658 (8th 266)P	P618 (\$22050)	sizerete CTA
	shellyweged-2009s	Aprilio	ations (Qualphs)Co	lo <b>nia 5</b> (x 1,27)	178.4.(x 1,20)
-	Denowura_2004§	BKSP	perfect (tylend)	BERGER CAA	disapetre3v20pp
-	shaddinisa uzna4s	9835	2009 (34.20)	3628 (*3.62)	4312 (x3.28)
	history Netura2009s	720p5	31692(421 <del>38</del> 6)	52B (47.8P)	494 (x4.36)7)
	roddobinaca2004s	2 <b>%</b> ppli	ca <b>tilos:(631449)</b> Co	ol <b>b84:1(§</b> 32:993)	8928 (x2.42)
-	road_usa	38.2	51.4 (x0.74)	19.3 (x1.97)	81.9 (x0.46)
-	rDataset_cam	or Pipe	cational signification	loggesist CTA	1discrete warp
-	Dataset Soc-LiveJournal1s	BSP <sub>.5</sub>	persist warp	persist C.T.A	discrete warp
-	hollywagd 2009s	077 <sub>5</sub> 9	31 <sub>4</sub> 9 <sub>4</sub> (x <sub>1</sub> 2 <sub>4</sub> )	359 <sub>1</sub> 3 <sub>1</sub> (x <sub>1</sub> 1 <sub>3</sub> 1 <sub>1</sub> )	(274 (x0.28)
		0/4/87	24√1/2/14/9/29 A	24212(1/912/91)	∠9/7 \\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\

Table 2.005.000 datasets used (not our continued by 13.2 (x2.5) (x3.5) (

Table 2: Summary of datasets used in out experiments. Graph Table 2: Summary verifies these Diams. Diams outdeen degree of the summary of the

	mony wood_boos	1.1271			11,107	11,107	100	
	indochina_2004 <sup>s</sup>	7.4M	191M	26	256,425	6,984	. 8	
	road_usa <sup>m</sup>	23.9M	57M	6.809	Magx.	Magx.	Avg.	
	rbatinet_cam	Vertices	Edges	Diam.	indeg.	outdeg.	degree	
-	Dataset LiveJournal1s	Vertices	Edges	Diam <sub>20</sub>	index5	outdeg2	$degree_4$	
_	hollywood 2009 <sup>s</sup>	$4^{1}8M$	68M	$2b^{1}$	13,9057	20,2927	195	
	rindochina_20048	1 <sup>7</sup> 1 <sup>4</sup> M	191M	126	256425	11,46%	105 <sup>8</sup>	
Gun	roek sbacket-	bas <b>ëd</b> Ma	ta <sup>57</sup> Mr	alle <b>r</b> to	a&b&babar	icingen	ethod $\frac{1}{2}$ [28	3]
1	road Net nea		57M1	6,809	92	0 1/2	222	-
aesc		11 31 3x4	E 3.4	0.40	10	10	2	

We run the three case studies on three scale-free and two meshlike detasets (Table 2) and studies itself the custimer expedition resultafine beautiful (Three dynamical fine for the confine (speedup) results like three stip (Table 12) and surprise the runtime (speedup) results 602 fou Portformanions Ghadlenges in Three Study

# 6.2 Cestormance Challenges in Three Study

6.2 Performance Challenges in Three Study Each of oil Study cases embodies a subset of the BSP performance challengers Sausseds in Santianes 1. subjet in the aspectorian afe Tachne februisitudy in a the Ambordina presudent Western Barbaires albert Barb Rhellenensedishukadin badida howhich waforms auzehnie Bop destinational antique that the state of the ptirformer new tatal recommendate in the seation of his they were ignerdarkoool and admed interblane: This trace algorithm in involvenitors' acide for extate rutter in meighbore direction reasoning the billion exercises of destedile atablead in the lance in this its sur is need be men drawn area. vandenfundstatafeld naderlowertex dagraede grien en is his bellowerte et e mastatike kacan tich Peploke makri Figur de gree and Brenegel buth do ugte par sancing rein Table Post pen BSP Pinaple mentation of the 1th ren alconrithmall Grenting Problems by Firmstraktaning in orindistressible preseggiost tilending Motive atte Bollang termentation to the then the aced goe rither Lovi the the physical strain in the control of the control parases canfil be readed and their problems of the frontier peroblems a Telefrice detectors of a birted continued based the branch for the problem of this time. belongseen on last part characteristic hard a language of the street of the last control of the last contr number of BSRt iterations and chigh a yet ago dagge cleading to a large amount of work per iteration). In contrast, Gunrock on

Table 3: Summary of performance challenges for each case study.

	BFS	PageRank	Graph Coloring
Scale-Free	Load Imbalance	Load Imbalance	Load Imbalance + Small Frontier
Mesh-Like	Small Frontier	None	None

Table 4: Upper: Workload ratio of three Atos implementations relative to Gunrock's implementations for BFS and PageRank. A workload ratio of *n* means that our implementation does *n* times as much work as Gunrock. Lower: the workload ratio of four implementations relative to the input graph's total vertex count for graph coloring.

	Application: BFS		Application: PageRank				
Dataset	persist	persist	discrete	persist	persist	discrete	
Dataset	warp	CTA	CTA	warp	CTA	CTA	
soc-LiveJournal1 <sup>8</sup>	1.43	1.06	1.01	0.73	0.72	0.72	
hollywood_2009 <sup>s</sup>	2.26	1.19	1.07	1.08	1.18	0.9	
indochina_2004 <sup>8</sup>	1.28	1.00	1.00	0.76	0.73	0.75	
road_usa <sup>m</sup>	3.56	1.05	1.04	0.79	0.79	0.92	
roadNet_ca <sup>m</sup>	2.05	1.02	1.04	1.18	1.11	0.97	
Application: Graph Coloring							
Dataset	BSP	F	ersist	persist	disc	crete	
Dataset	DSP	warp		CTA	W	arp	
soc-LiveJournal1 <sup>s</sup>	1.17	3.31 1.15		1.74	2.78		
hollywood_2009 <sup>s</sup>	3.31			5.24	37.3	37.34 16.97	
indochina_2004 <sup>s</sup>	1.96			4.45	16.9		
road_usa <sup>m</sup>	1.22	1	.00	1.46	1.41	1	
JNT_4 m	255	1	00	1.71	0.4	4	

mesh-like datasets *does exhibit* the small frontier problem, because these datasets have high diameters and small average degree; consequently, there is a large number of iterations, with little work per iteration, leading to low throughput over many iterations.

PageRank: In Figure 2, for Gunrock PageRank, both scale-free and mesh-like datasets do not exhibit the small frontier problem, as they have high throughput over most of the execution time and converge in fewer than 35 iterations (though Indochina-2004 exhibits a long flat tail in the latter half of execution).

Graph Coloring: In Figure 3, for BSP graph coloring, scale free datasets have low throughput for more than 70% of execution time, and thus have the small frontier problem. Mesh-like datasets terminate in fewer than 40 iterations, and have short tails, and thus do not have the small frontier problem. This is because on scale-free datasets, the high-degree vertices will need to be recolored many times, leading to a large number of iterations, during which the frontier contains a few high-degree vertices with color conflicts. In contrast, mesh-like datasets have low average degree, and are less likely to have color conflicts.

## 6.3 Relaxing Barriers

As discussed in Sections 2–3, relaxing barriers exposes more concurrency, giving higher throughput and shorter execution time. However, relaxing barriers may result in extra work. If the performance improvement from increased concurrency outweighs the cost of extra work, we obtain a net performance gain.

There are two key factors influencing this tradeoff. First, we find that in the presence of a small frontier problem, the increase in concurrency from relaxing barriers is always more significant than the cost of extra work. Second, on naturally unordered algorithms such

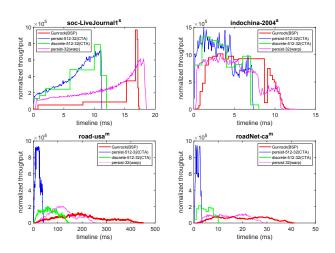


Figure 1: Normalized throughput vs. time on BFS. The top charts are scale-free; bottom charts are mesh-like.

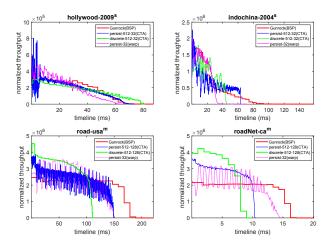


Figure 2: Normalized throughput vs. time on PageRank.

as PageRank, one can always relax the barrier: although the barrier gives the BSP implementation a more predictable convergence rate, the barrier generally does not make the convergence faster.

Table 4 summarizes the extra work for three study cases. Figures 1, 2 and 3 plot the throughput of four implementations (BSP + three Atos variants) of three study cases against timeline for four datasets. Notably, these plots show the *normalized throughput*, which is the measured throughput divided by the overwork factors in Table 4. This gives a fair measure of overall performance, as it incorporates both the benefits of improved concurrency (higher absolute throughput) and the cost of extra work. Essentially, normalized throughput measures "useful" throughput rather than raw absolute throughput. We provide detailed analysis below for each application.

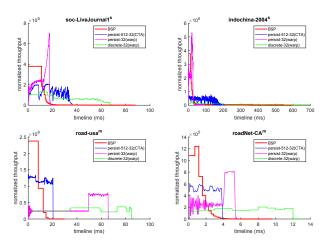


Figure 3: Normalized throughput vs. time on graph coloring.

BFS: Figure 1 shows that for the two mesh-like datasets, all 3 Atos implementations achieve considerably higher normalized throughput than Gunrock. Why? Table 3 shows that Gunrock on mesh-like datasets has a severe small frontier problem. Therefore, the increase in concurrency in the 3 Atos implementations offers a significant performance advantage. Table 4 indicates that the persistent-warp implementation generates 3.5x extra work vs. Gunrock, but despite this extra work, Atos's normalized throughput is still significantly higher than Gunrock. Scale-free graphs, on the other hand, exhibit more parallelism and do not suffer from the small-frontier problem. Atos's fastest implementations are still faster than Gunrock's, but not nearly as much as for the mesh networks.

On all BFS experiments, Atos's CTA implementations are faster than its warp ones. Atos's CTA implementations use a combination of task-parallel and data-parallel load balancing techniques (see Section 6.4 for details), and thus have better load balancing than its warp implementations, which only use task-parallel load balancing. This leads to higher GPU utilization and hence higher absolute throughput. Second, CTA implementations produce less extra work than warp (see Table 4). Due to better load balancing in CTA, the workload of each worker has lower variance. If a worker receives too much work, there will be a long delay before the vertices' updated depths are visible to other workers; this increases the likelihood that downstream vertices are first reached via other sub-optimal paths, which leads to extra work.

**PageRank:** Unlike BFS, PageRank is naturally unordered, as it satisfies Dijkstra's don't care non-determinism [11]. Therefore, relaxing the barrier in the outer loop does not generate any misspeculations and hence results in no wasted work. In fact, Table 4 shows that the Atos implementations perform *less* work than Gunrock in general. This is because the BSP barrier forces each vertex to be processed at most once per iteration. By relaxing this barrier, the Atos implementations can update certain important vertices (e.g., vertices with high centrality) more frequently than other vertices, thus leading to more efficient propagation of rank.

Figure 2 shows that all three Atos implementations compact the workload and process it with higher normalized throughput (persist-CTA has a higher profiling cost). Though PageRank does not suffer from the small frontier problem, the three Atos implementations nonetheless have superior performance over Gunrock, because relaxing the barrier increases concurrency. In addition, relaxing the barrier lowered the overall workload in practice, even though in theory it may lead to a more unpredictable convergence rate.

**Graph Coloring:** Unlike BFS and PageRank, all graph coloring implementations (including BSP) use a speculative approach (greedy graph coloring) and thus all have extra work. Table 4 summarizes the multiplicative factor of extra work, which is defined as a ratio vs. the number of vertices in the graph (the lowest possible workload). Atos's persist warp has the least extra work; on some datasets, the extra work is less than 1%, which means after the first color assignment, only 1% of vertices have a color conflict and must be recolored. Atos's discrete warp has the most extra work (on hollywood-2009, 37.34x). The extra work is due to the combination of two factors:

1. Conflicts tend to arise when neighboring vertices are colored concurrently: From Section 5.3, given a vertex, the algorithm first checks its neighbors' colors, then assigns a color to the vertex that does not conflict with its neighbors. The color assignment is speculative because it is done using possibly outdated color information from the vertex's neighbors. When neighboring vertices are colored simultaneously, they read outdated colors from each other, leading to conflicts and recoloring.

2. Consecutive vertices on the work queue are likely to be neighbors: On many if not most graphs, the vertex ID is semantically meaningful: vertices whose vertex ID are numerically close are more likely to be neighbors. At the beginning of graph coloring, all vertices are initially inserted onto the work queue in order of vertex ID.

Since consecutive vertices on the work queue tend to be assigned colors concurrently, the above implies a high likelihood of color conflicts. We verify that the large amount of extra work is indeed due to semantically meaningful vertex IDs: running the exact same experiment with randomly permuted vertex IDs, the amount of extra work drops to less than 1.5x for all four implementations on all datasets. ID permutation leads to the following runtime improvements (in ms) on scale-free datasets:

Impl.	soc-LiveJournal1	hollywood	indochina	
discrete-warp persist-CTA BSP	$63 \rightarrow 31$ $36 \rightarrow 21$ $96 \rightarrow 89$	$274 \rightarrow 26$ $59 \rightarrow 28$ $77 \rightarrow 61$	$2073 \rightarrow 222$ $184 \rightarrow 50$ $673 \rightarrow 485$	

The BSP implementation has a more modest improvement because BSP's thread-warp-CTA load balancing scheme [21] already divides each bucket into three individually-load-balanced subbuckets, reducing inter-bucket conflicts. Persist-warp has little change as there is almost no extra work even before permutation. Notably, after permutation, all three Atos variants are faster than BSP implementation on scale-free datasets.

Comparing persist-warp and persist-CTA: persist-CTA has better load balancing, allowing for more (potentially adjacent) vertices to be colored simultaneously, resulting in more extra work than persist-warp. We verify this from Table 4. Roughly speaking, the amount of extra work for persist-CTA is more significant on

scale-free graphs, as a vertex can have a large number of neigl bors, leading to more potential conflicts. Therefore, persist-CT. outperforms persist-warp on mesh-like graphs, where the increase concurrency and better load-balancing outweigh the cost of waste work; conversely, on scale-free datasets, persist-CTA is slower tha persist-warp, because the cost of extra work is too high (see Tables and 4).

Comparing persist-warp and discrete-warp: Discrete-war has more extra work than persist-warp, hurting its performanc for two reasons.

- (1) The scheduling policies of discrete and persistent kernel are different. When kernels are launched from the CPU (discret kernel strategy), the kernel launched earlier always has a highe scheduling (hardware) priority than the kernel launched later. This effectively causes vertices to always be colored in roughly the same order as their initial ordering (by vertex ID, which causes many conflicts). In contrast, the persistent kernel only incurs one kernel launch and warps within it are scheduled by the hardware scheduler, whose decisions are much less ordered by vertex ID. Thus persistwarp has fewer coloring conflicts caused by adjacencies and hence less overwork.
- (2) Discrete-warp has lower register usage than persistent-warp (72 vs. 42), so persist-warp only achieves 43% occupancy per SM and discrete-warp achieves 62%. Therefore the discrete-wrap assigns colors to more vertices simultaneously, leading to a greater likelihood of conflicts than persist-warp. Unlike our other applications, in graph coloring, the cost of extra work largely reduces the benefit of increased concurrency. On scale-free datasets (without random permutation), our highest performance is achieved with lower concurrency and less overwork (persist-warp variant), which achieves a lower absolute throughput but a higher normalized throughput (and hence higher performance overall).

## 6.4 Worker Size and Trade-off between Taskand Data-Parallelism Load Balancing

As discussed in Section 3, Atos enables the user to trade off between task-parallelism and data-parallelism load balancing by adjusting the worker size and FETCH\_SIZE. Atos's persist-CTA, like persistwarp, uses a persistent kernel to exploit task parallelism, but now the task-parallel work units are fewer and larger (the size of a CTA) and we can leverage more data parallelism within a CTA. In most cases, persist-CTA outperforms persist-warp with both higher normalized/absolute throughput, except for the graph coloring on scale-free datasets, where it achieves only higher absolute throughput. Figure 4 illustrates this tradeoff for BFS and PageRank on soc-LiveJournal (scale-free) and road\_usa (mesh-like). We exclude graph coloring because it can only be run with one CTA size, due to high register usage (72) and high shared memory usage (46 KB).

## 6.5 Kernel Strategy

From Section 3, the chief advantage of the persistent kernel is removing the overhead associated with kernel invocation, which is most significant for fine-grained tasks that involve many small kernel launches. Based on the performance results in Table 1 and

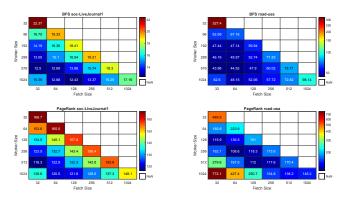


Figure 4: Runtime (ms) heatmap plotted with different worker size and fetch size for BFS and PageRank on soc-LiveJournal<sup>s</sup> and road\_usa<sup>m</sup>. Note only the lower triangle is valid.

Figure 1, the performance gap between persistent kernel and discrete kernel is particularly large for BFS on mesh-like graphs, which require many small kernel launches due to the high diameter and small workload per iteration. Graph coloring on indochina-2004 also shows a large kernel launch overhead. Using a random permutation of vertex IDs (see Section 6.3), Atos's persistent variant is 4.3x faster than its discrete variant.

## 7 CONCLUSION

In this paper, we present our task-parallel GPU dynamic scheduling framework, Atos, and analyze its performance across numerous design parameters on three case studies. Our analysis provides the following guidelines on what applications are suitable to run in a capable task-parallel framework, as well as what Atos configurations to use, given an application's characteristics:

- If the dynamic application either exhibits the small frontier problem or has load imbalance, Atos will have a performance advantage.
- (2) If the application exhibits the small frontier problem, it should be run with a persistent kernel.
- (3) If the application exhibits load imbalance, it should be run with both task- and data-parallelism load balancing in tandem to achieve better performance. For different applications, the optimal tradeoff point varies.
- (4) By relaxing the outer loop dependency in the application, Atos increases concurrency at the cost of extra work due to mis-speculation, or less predictable convergence rates. The optimal tradeoff between the increased concurrency and additional cost is application-dependent. When an application is naturally unordered (e.g., PageRank) or has the small frontier problem (e.g., BFS on mesh-like datasets and graph coloring on scale-free datasets), the increased concurrency usually outweighs the cost. Conversely, on problems such as BFS on scale-free graphs or graph coloring on mesh-like graphs, the cost of extra work can hurt performance. The best way to reduce extra work is application-dependent and

may include better load balancing (e.g., BFS) or reducing concurrency (e.g., graph coloring).

## **ACKNOWLEDGMENTS**

This work is supported by the National Science Foundation (NSF) under projects CCF-1823034, CCF-1823037, and OAC-1740333; by the Department of Defense Advanced Research Projects Agency (DARPA) under projects HR0011-18-3-0007 and FA8650-18-2-7835; by an NVIDIA gift and hardware donations; and by the Advanced Scientific Computing Research (ASCR) program within the Office of Science of the DOE under contract number DE-AC02-05CH11231 and the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration.

## **REFERENCES**

- Cédric Augonnet, Jérôme Clet-Ortega, Samuel Thibault, and Raymond Namyst. 2010. Data-Aware Task Scheduling on Multi-Accelerator based Platforms. In 2010 IEEE 16th International Conference on Parallel and Distributed Systems. IEEE, 291–298. https://doi.org/10.1109/ICPADS.2010.129
- [2] Michael Bauer, Sean Treichler, Elliott Slaughter, and Alex Aiken. 2012. Legion: Expressing Locality and Independence with Logical Regions. In SC'12: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis. IEEE, 1–11. https://doi.org/10.1109/SC.2012.71
- [3] Sean Baxter. 2013. moderngpu: Load-Balancing Search. https://moderngpu.github.io/loadbalance.html. Accessed: 2022-07-29.
- [4] Mehmet E. Belviranli, Seyong Lee, Jeffrey S. Vetter, and Laxmi N. Bhuyan. 2018. Juggler: A Dependence-Aware Task-Based Execution Framework for GPUs. In Proceedings of the 23rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '18). 54-67. https://doi.org/10.1145/3178487. 3178492
- [5] Tal Ben-Nun, Michael Sutton, Sreepathi Pai, and Keshav Pingali. 2017. Groute: An Asynchronous Multi-GPU Programming Model for Irregular Computations. In Proceedings of the 22nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '17). 235–248. https://doi.org/10.1145/3018743. 3018756
- [6] Guy E. Blelloch and Gary W. Sabot. 1990. Compiling Collection-Oriented Languages onto Massively Parallel Computers. J. Parallel and Distrib. Comput. 8, 2 (Feb. 1990), 119–134. https://doi.org/10.1016/0743-7315(90)90087-6
- [7] Daniel Cederman and Philippas Tsigas. 2008. On Dynamic Load-Balancing on Graphics Processors. In Graphics Hardware (GH '08). 57–64. https://doi.org/10. 2312/EGGH/EGGH08/057-064
- [8] Long Chen, Oreste Villa, Sriram Krishnamoorthy, and Guang R Gao. 2010. Dynamic Load Balancing on Single- and Multi-GPU Systems. In 2010 IEEE International Symposium on Parallel & Distributed Processing (IPDPS 2010). IEEE. https://doi.org/10.1109/IPDPS.2010.5470413
- [9] Andrew Davidson, Sean Baxter, Michael Garland, and John D. Owens. 2014. Work-Efficient Parallel GPU Methods for Single Source Shortest Paths. In Proceedings of the 28th IEEE International Parallel and Distributed Processing Symposium (IPDPS 2014). 349–359. https://doi.org/10.1109/IPDPS.2014.45
- [10] Timothy A. Davis. 2019. Algorithm 1000: SuiteSparse: GraphBLAS: Graph algorithms in the language of sparse linear algebra. ACM Transactions on Mathematical Software (TOMS) 45, 4 (2019), 1–25. https://doi.org/10.1145/3322125
- [11] Edsger W. Dijkstra. 1976. A Discipline of Programming. Pearson.
- [12] Alex Fender, Brad Rees, and Joe Eaton. 2022. RAPIDS cuGraph. In Massive Graph Analytics. Chapman and Hall/CRC, Chapter 17, 483–493. https://doi.org/10.1201/ 9781003033707-22
- [13] Assefaw Hadish Gebremedhin and Fredrik Manne. 2000. Scalable parallel graph coloring algorithms. Concurrency: Practice and Experience 12, 12 (Nov. 2000), 1131–1146. https://doi.org/10.1002/1096-9128(200010)12:12<1131::AID-CPE528> 3.0 CO:2-2
- [14] Kshitij Gupta, Jeff Stuart, and John D. Owens. 2012. A Study of Persistent Threads Style GPU Programming for GPGPU Workloads. In Proceedings of Innovative Parallel Computing (InPar '12). https://doi.org/10.1109/InPar.2012.6339596
- [15] Shawn Hargreaves. 2004. Generating Shaders from HLSL Fragments. In ShaderX3: Advanced Rendering with DirectX and OpenGL. Chapter 7.3, 555–568.
- [16] Muhammad Amber Hassaan, Martin Burtscher, and Keshav Pingali. 2011. Ordered vs. Unordered: a Comparison of Parallelism and Work-efficiency in Irregular Algorithms. In Proceedings of the 16th ACM Symposium on Principles and Practice of Parallel Programming (San Antonio, TX, USA) (PPoPP '11). 3-12. https://doi. org/10.1145/1941553.1941557

- [17] Vishwesh Jatala, Roshan Dathathri, Gurbinder Gill, Loc Hoang, V. Krishna Nandivada, and Keshav Pingali. 2020. A Study of Graph Analytics for Massive Datasets on Distributed Multi-GPUs. In 2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS). 84–94. https://doi.org/10.1109/IPDPS47924.2020. 00019
- [18] Zhihao Jia, Yongkee Kwon, Galen Shipman, Pat McCormick, Mattan Erez, and Alex Aiken. 2017. A Distributed Multi-GPU System for Fast Graph Processing. Proc. VLDB Endow. 11, 3 (Nov. 2017), 297–310. https://doi.org/10.14778/3157794. 3157700
- [19] Milind Kulkarni, Keshav Pingali, Bruce Walter, Ganesh Ramanarayanan, Kavita Bala, and L. Paul Chew. 2007. Optimistic Parallelism Requires Abstractions. In Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation. 211–222. https://doi.org/10.1145/1250734.1250759
- [20] Hang Liu and H. Howie Huang. 2019. SIMD-X: Programming and Processing of Graph Algorithms on GPUs. In 2019 USENIX Annual Technical Conference (USENIX ATC 19). USENIX Association, Renton, WA, 411–428. https://www. usenix.org/conference/atc19/presentation/liu-hang
- [21] Duane Merrill, Michael Garland, and Andrew Grimshaw. 2012. Scalable GPU Graph Traversal. In Proceedings of the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '12). 117–128. https://doi.org/10. 1145/2145816.2145832
- [22] Christopher J. Rossbach, Jon Currey, Mark Silberstein, Baishakhi Ray, and Emmett Witchel. 2011. PTask: Operating System Abstractions To Manage GPUs as Compute Devices. In Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles (SOSP '11). 233–248. https://doi.org/10.1145/2043556.2043579
- [23] Markus Steinberger, Michael Kenzel, Pedro Boechat, Bernhard Kerbl, Mark Dokter, and Dieter Schmalstieg. 2014. Whippletree: Task-Based Scheduling of Dynamic Workloads on the GPU. ACM Transactions on Graphics 33, 6, Article 228 (Nov. 2014), 11 pages. https://doi.org/10.1145/2661229.2661250
- [24] Gunrock team. 2017. Throughput vs. Frontier Size. https://gunrock.github.io/docs/#/analysis/frontier\_size. Accessed: 2022-07-29.
- [25] Stanley Tzeng, Brandon Lloyd, and John D. Owens. 2012. A GPU Task-Parallel Model with Dependency Resolution. IEEE Computer 45, 8 (Aug. 2012), 34–41. https://doi.org/10.1109/MC.2012.255
- [26] Stanley Tzeng, Anjul Patney, and John D. Owens. 2010. Task Management for Irregular-Parallel Workloads on the GPU. In Proceedings of High Performance Graphics (HPG '10). 29–37. https://doi.org/10.2312/EGGH/HPG10/029-037
- [27] Leslie G. Valiant. 1990. A Bridging Model for Parallel Computation. Commun. ACM 33, 8 (Aug. 1990), 103–111. https://doi.org/10.1145/79173.79181
- [28] Yangzihao Wang, Yuechao Pan, Andrew Davidson, Yuduo Wu, Carl Yang, Leyuan Wang, Muhammad Osama, Chenshan Yuan, Weitang Liu, Andy T. Riffel, and John D. Owens. 2017. Gunrock: GPU Graph Analytics. ACM Transactions on Parallel Computing 4, 1 (Aug. 2017), 3:1–3:49. https://doi.org/10.1145/3108140
- [29] Tsung Tai Yeh, Amit Sabne, Putt Sakdhnagool, Rudolf Eigenmann, and Timothy G. Rogers. 2017. Pagoda: Fine-Grained GPU Resource Virtualization for Narrow Tasks. In Proceedings of the 22nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '17). 221–234. https://doi.org/10.1145/3018743.3018754
- [30] Jianlong Zhong and Bingsheng He. 2014. Medusa: Simplified Graph Processing on GPUs. IEEE Transactions on Parallel and Distributed Systems 25, 6 (2014), 1543–1552. https://doi.org/10.1109/TPDS.2013.111