

MIT Open Access Articles

A scalable architecture for ordered parallelism

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Jeffrey, Mark C. et al. "A Scalable Architecture for Ordered Parallelism." Proceedings of the 48th International Symposium on Microarchitecture (MICRO-48), December 5-9, 2015, Waikiki, Hawaii, USA, Association for Computing Machinery (ACM), December 2015 © 2015 Association for Computing Machinery (ACM)

As Published: <http://dx.doi.org/10.1145/2830772.2830777>

Publisher: Association for Computing Machinery (ACM)

Persistent URL: <http://hdl.handle.net/1721.1/111962>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of use: Creative Commons Attribution-Noncommercial-Share Alike



A Scalable Architecture for Ordered Parallelism

Mark C. Jeffrey[†] Suvinay Subramanian[†] Cong Yan[†] Joel Emer* Daniel Sanchez[†]

[†]MIT CSAIL *NVIDIA / MIT CSAIL
{mcj, suvinay, congy, emer, sanchez}@csail.mit.edu

ABSTRACT

We present Swarm, a novel architecture that exploits *ordered irregular parallelism*, which is abundant but hard to mine with current software and hardware techniques. In this architecture, programs consist of short tasks with programmer-specified timestamps. Swarm executes tasks speculatively and out of order, and efficiently speculates thousands of tasks ahead of the earliest active task to uncover ordered parallelism. Swarm builds on prior TLS and HTM schemes, and contributes several new techniques that allow it to scale to large core counts and speculation windows, including a new execution model, speculation-aware hardware task management, selective aborts, and scalable ordered commits.

We evaluate Swarm on graph analytics, simulation, and database benchmarks. At 64 cores, Swarm achieves 51–122× speedups over a single-core system, and outperforms software-only parallel algorithms by 3–18×.

Categories and Subject Descriptors

C.1.4 [Processor architectures]: Parallel architectures

Keywords

Multicore, ordered parallelism, irregular parallelism, fine-grain parallelism, synchronization, speculative execution

1. INTRODUCTION

Parallel architectures are now pervasive, but thread-level parallelism in applications is often scarce [19, 35]. Thus, it is crucial that we explore new architectural mechanisms to efficiently exploit as many types of parallelism as possible. Doing so makes parallel systems more versatile, easier to program, and, for many applications, it is the only way to improve performance.

We focus on *ordered irregular parallelism* [55], which is often abundant but hard to exploit. Programs with ordered irregular parallelism have three key features.

First, they consist of tasks that must follow a total or partial order. Second, tasks may have data dependences that are not known a priori. Third, tasks are not known in advance. Instead, tasks dynamically create children tasks and schedule them to run at a future time.

Ordered irregular parallelism is abundant in many domains, such as simulation, graph analytics, and databases. For example, consider a timing simulator for a parallel computer. Each task is an event (e.g., executing an instruction in a simulated core). Each task must run at a specific simulated time (introducing order constraints among tasks), and modifies a specific component (possibly introducing data dependences among tasks). Tasks dynamically create other tasks (e.g., a simulated memory access), possibly for other components (e.g., a simulated cache), and schedule them for a future simulated time.

Prior work has tried to exploit ordered parallelism in software [33, 34], but has found that, in current multicores, runtime overheads negate the benefits of parallelism. This motivates the need for architectural support.

To guide our design, we first characterize several applications with ordered irregular parallelism (Sec. 2). We find that tasks in these applications are as small as a few tens of instructions. Moreover, many of these algorithms rarely have true data dependences among tasks, and their maximum achievable parallelism exceeds 100×. It may seem that thread-level speculation (TLS) [28, 60, 66, 68], which speculatively parallelizes sequential programs, could exploit this parallelism. However, this is not the case due to two reasons (Sec. 3):

- Ordered irregular algorithms have little parallelism when written as sequential programs. To enforce order constraints, sequential implementations introduce *false data dependences* among otherwise independent tasks. For example, sequential implementations of timing simulators use a priority queue to store future tasks. Priority queue accesses introduce false data dependences that limit the effectiveness of TLS.
- To scale, ordered irregular algorithms need very large speculation windows, of thousands of tasks (hundreds of thousands of instructions). Prior TLS schemes use techniques that scale poorly beyond few cores and cannot support large speculation windows.

We present Swarm, an architecture that tackles these challenges. Swarm consists of (i) a task-based execution model where order constraints do not introduce false data dependences, and (ii) a microarchitecture that

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

MICRO-48, December 05–09, 2015, Waikiki, HI, USA

Copyright is held by the owner/author(s). Publication rights licensed to ACM.

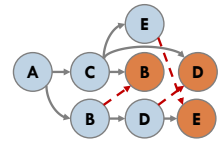
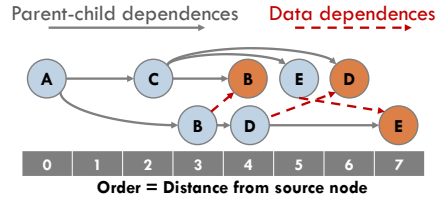
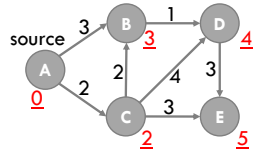
ACM 978-1-4503-4034-2/15/12 ...\$15.00

DOI: <http://dx.doi.org/10.1145/2830772.2830777>

```

prioQueue.enqueue(source, 0)
while prioQueue not empty:
  (node, dist) = prioQueue.dequeueMin()
  if node.distance not set:
    node.distance = dist
    for child in node.children:
      d = dist + distance(node, child)
      prioQueue.enqueue(child, d)
  else: // node already visited, skip

```



(a) Dijkstra’s sssp code, highlighting the non-visited and visited paths that each task may follow

(b) Example graph and resulting shortest-path distances (underlined)

(c) Tasks executed by sssp. Each task shows the node it visits. Tasks that visit the same node have a data dependence

(d) A correct speculative schedule that achieves 2x parallelism

Figure 1: Dijkstra’s single-source shortest paths algorithm (sssp) has plentiful ordered irregular parallelism.

leverages this execution model to scale efficiently (Sec. 4).

Swarm is a tiled multicore with distributed task queues, speculative out-of-order task execution, and ordered task commits. Swarm adapts prior eager version management and conflict detection schemes [48, 79], and features several new techniques that allow it to scale. Specifically, we make the following novel contributions:

- An execution model based on tasks with programmer-specified timestamps that conveys order constraints to hardware without undue false data dependences.
- A hardware task management scheme that features speculative task creation and dispatch, drastically reducing task management overheads, and implements a very large speculation window.
- A scalable conflict detection scheme that leverages eager versioning to, upon mispeculation, selectively abort the mispeculated task and its dependents (unlike prior TLS schemes that forward speculative data, which abort all later tasks).
- A distributed commit protocol that allows ordered commits without serialization, supporting multiple commits per cycle with modest communication (unlike prior schemes that rely on successor lists, token-passing, and serialized commits).

We evaluate Swarm in simulation (Sec. 5 and Sec. 6) using six challenging workloads: four graph analytics algorithms, a discrete-event simulator, and an in-memory database. At 64 cores, Swarm achieves speedups of 51–122× over a single-core Swarm system, and outperforms state-of-the-art parallel implementations of these algorithms by 2.7–18.2×. In summary, by making ordered execution scalable, Swarm speeds up challenging algorithms that are currently limited by stagnant single-core performance. Moreover, Swarm simplifies parallel programming, as it frees developers from using error-prone explicit synchronization.

2. MOTIVATION

2.1 Ordered Irregular Parallelism

Applications with ordered irregular parallelism have three main characteristics [33, 55]. First, they consist of tasks that must follow a total or partial order. Second, tasks are not known in advance. Instead, tasks dynamically create children tasks, and schedule them to run at a future time. Task execution order is different from task creation order. Third, tasks may have data dependences that are not known a priori.

Ordered irregular algorithms are common in many domains. First, they are common in graph analytics, especially in search problems [33, 55]. Second, they are important in simulating systems whose state evolves over time, such as circuits [47], computers [12, 59], networks [37, 72], healthcare systems [39], and systems of partial differential equations [32, 44]. Third, they are needed in systems that must maintain externally-imposed order constraints, such as geo-replicated databases where transactions must appear to execute in timestamp order [14], or deterministic architectures [17, 45] and record-and-replay systems [36, 77] that constrain the schedule of parallel programs to ensure deterministic execution.

To illustrate the challenges in parallelizing these applications, consider Dijkstra’s single-source shortest paths (sssp) algorithm [15, 22]. sssp finds the shortest distance between some source node and all other nodes in a graph with weighted edges. Fig. 1(a) shows the sequential code for sssp, which uses a priority queue to store tasks. Each task operates on a single node, and is ordered by its tentative distance to the source node. sssp relies on task order to guarantee that the first task to visit each node comes from a shortest path. This task sets the node’s distance and enqueues all its children. Fig. 1(b) shows an example graph, and Fig. 1(c) shows the tasks that sssp executes to process this graph. Fig. 1(c) shows the order of each task (its distance to the source node) in the x -axis, and outlines both parent-child relationships and data dependences. For example, task A at distance 0, denoted $(A, 0)$, creates children tasks $(C, 2)$ and $(B, 3)$; and tasks $(B, 3)$ and $(B, 4)$ both access node B , so they have a data dependence.

A distinctive feature of irregular parallel programs is that task creation and execution order are different: children tasks are not immediately runnable, but are subject to a global order influenced by all other tasks in the program. For example, in Fig. 1(c), $(C, 2)$ creates $(B, 4)$, but running $(B, 4)$ immediately would produce the wrong result, as $(B, 3)$, created by a different parent, must run first. Sequential implementations of these programs use scheduling data structures, such as priority or FIFO queues, to process tasks in the right order. These scheduling structures introduce false data dependences that restrict parallelism and hinder TLS (Sec. 3).

Order constraints limit non-speculative parallelism. For example, in Fig. 1(c), only $(B, 4)$ and $(D, 4)$ can run in parallel without violating correctness. A more attractive option is to use speculation to elide order

constraints. For example, Fig. 1(d) shows a speculative schedule for **sssp** tasks. Tasks in the same x -axis position are executed simultaneously. This schedule achieves $2\times$ parallelism in this small graph; larger graphs allow more parallelism (Sec. 2.2). This schedule produces the correct result because, although it elides order constraints, it happens to respect data dependences. Unfortunately, data dependences are not known in advance, so speculative execution must detect dependence violations and abort offending tasks to preserve correctness.

Recent work has tried to exploit ordered parallelism using speculative software runtimes [33, 34], but has found that the overheads of ordered, speculative execution negate the benefits of parallelism. This motivates the need for hardware support.

2.2 Analysis of Ordered Irregular Algorithms

To quantify the potential for hardware support and guide our design, we first analyze the parallelism and task structure of several ordered irregular algorithms.

Benchmarks: We analyze six benchmarks from the domains of graph analytics, simulation, and databases:

- **bfs** finds the breadth-first tree of an arbitrary graph.
- **sssp** is Dijkstra’s algorithm (Sec. 2.1).
- **astar** uses the A* pathfinding algorithm [31] to find the shortest route between two points in a road map.
- **msf** is Kruskal’s minimum spanning forest algorithm [15].
- **des** is a discrete-event simulator for digital circuits.

Each task represents a signal toggle at a gate input.

- **sil** is an in-memory OLTP database [71].

Sec. 5 describes their input sets and methodology details.

Analysis tool: We developed a pintool [46] to analyze these programs in x86-64. We focus on the instruction length, data read and written, and intrinsic data dependences of tasks, excluding the overheads and serialization introduced by the specific runtime used.

The tool uses a simple runtime that executes tasks sequentially. The tool profiles the number of instructions executed and addresses read and written (i.e., the read and write sets) of each task. It filters out reads and writes to the stack, the priority queue used to schedule tasks, and other runtime data structures such as the memory allocator. With this information, the tool finds the *critical path length* of the algorithm: the sequence of data-dependent tasks with the largest number of instructions. The tool then finds the *maximum achievable speedup* by dividing the sum of instructions of all tasks by the critical path length [78] (assuming unbounded cores and constant cycles per instruction). Note that this analysis *constrains parallelism only by true data dependences*: task order dictates the direction of data flow in a dependence, but is otherwise superfluous given perfect knowledge of data dependences.

Table 1 summarizes the results of this analysis. We derive three key insights that guide the design of Swarm: **Insight 1: Parallelism is plentiful.** These applications have at least $158\times$ maximum parallelism (**msf**), and up to $3440\times$ (**bfs**). Thus, most order constraints are superfluous, making *speculative execution* attractive. **Insight 2: Tasks are small.** Tasks are very short,

| Application | bfs | sssp | astar | msf | des | sil |
|------------------------------|--------------|--------------|--------------|-------------|--------------|-------------|
| Maximum parallelism | $3440\times$ | $793\times$ | $419\times$ | $158\times$ | $1440\times$ | $318\times$ |
| Parallelism window=1K | $827\times$ | $178\times$ | $62\times$ | $147\times$ | $198\times$ | $125\times$ |
| Parallelism window=64 | $58\times$ | $26\times$ | $16\times$ | $49\times$ | $32\times$ | $17\times$ |
| Instrs | mean | 22 | 32 | 195 | 40 | 296 |
| | 90th | 47 | 70 | 508 | 40 | 338 |
| Reads | mean | 4.0 | 5.8 | 22 | 7.1 | 50 |
| | 90th | 8 | 11 | 51 | 7 | 57 |
| Writes | mean | 0.33 | 0.41 | 0.26 | 0.03 | 10.5 |
| | 90th | 1 | 1 | 1 | 0 | 11 |
| Max TLS parallelism | $1.03\times$ | $1.10\times$ | $1.04\times$ | $158\times$ | $1.15\times$ | $45\times$ |

Table 1: Maximum achievable parallelism and task characteristics (instructions and 64-bit words read and written) of representative ordered irregular applications.

ranging from a few tens of instructions (**bfs**, **sssp**, **msf**), to a few thousand (**sil**). Tasks are also relatively uniform: 90th-percentile instructions per task are close to the mean. Tasks have small read- and write-sets. For example, **sssp** tasks read 5.8 64-bit words on average, and write 0.4 words. Small tasks incur large overheads in software runtimes. Moreover, order constraints prevent runtimes from grouping tasks into coarser-grain units to amortize overheads. *Hardware support for task management* can drastically reduce these overheads.

Insight 3: Need a large speculation window. Table 1 also shows the achievable parallelism within a limited task window. With a T -task window, the tool does not schedule an independent task until all work more than T tasks behind has finished. Small windows severely limit parallelism. For example, parallelism in **sssp** drops from $793\times$ with an infinite window, to $178\times$ with a 1024-task window, to $26\times$ with a 64-task window. Thus, for speculation to be effective, the architecture must support *many more speculative tasks than cores*.

These insights guide the design of Swarm. Our goal is to approach the maximum achievable parallelism while incurring only moderate overheads.

3. BACKGROUND ON HW SUPPORT FOR SPECULATIVE PARALLELISM

Much prior work has investigated thread-level speculation (TLS) schemes to parallelize sequential programs [25, 28, 60, 61, 66, 69]. TLS schemes ship tasks from function calls or loop iterations to different cores, run them speculatively, and commit them in program order. Although TLS schemes support ordered speculative execution, we find that two key problems prevent them from exploiting ordered irregular parallelism:

1. False data dependences limit parallelism: To run under TLS, ordered algorithms must be expressed as sequential programs, but their sequential implementations have limited parallelism. Consider the code in Fig. 1(a), where each iteration dequeues a task from the priority queue and runs it, potentially enqueueing

more tasks. Frequent data dependences *in the priority queue*, not among tasks themselves, cause frequent conflicts and aborts. For example, iterations that enqueue high-priority tasks often abort all future iterations.

Table 1 shows the maximum speedups that an ideal TLS scheme achieves on sequential implementations of these algorithms. These results use perfect speculation, an infinite task window, word-level conflict detection, immediate forwarding of speculative data, and no communication delays. Yet parallelism is meager in most cases. For example, `sssp` has $1.1\times$ parallelism. Only `msf` and `silob` show notable speedups, because they need no queues: their task orders match loop iteration order.

The root problem is that loops and method calls, the control-flow constructs supported by TLS schemes, are insufficient to express the order constraints among these tasks. By contrast, Swarm implements a more general execution model with timestamp-ordered tasks to avoid software queues, and implements hardware priority queues integrated with speculation mechanisms, avoiding spurious aborts due to queue-related references. **2. Scalability bottlenecks:** Although prior TLS schemes have developed scalable versioning and conflict detection schemes, two challenges limit their performance with large speculation windows and small tasks: *Forwarding vs selective aborts*: Most TLS schemes find it is desirable to forward data written by an earlier, still-speculative task to later reader tasks. This prevents later tasks from reading stale data, reducing mispeculations on tight data dependences. However, it creates complex chains of dependences among speculative tasks. Thus, upon detecting mispeculation, most TLS schemes abort the task that caused the violation *and all later speculative tasks* [25, 28, 61, 66, 68]. TCC [29] and Bulk [11] are the exception: they do not forward data and only abort later readers when the earlier writer commits.

We find that forwarding speculative data is crucial for Swarm. However, while aborting all later tasks is reasonable with small speculative windows (2–16 tasks are typical in prior work), Swarm has a 1024-task window, and unselective aborts are impractical. To address this, we contribute a novel conflict detection scheme based on eager version management that allows both forwarding speculative data and selective aborts of dependent tasks. *Commit serialization*: Prior TLS schemes enforce in-order commits by passing a token among ready-to-commit tasks [28, 61, 66, 68]. Each task can only commit when it has the token, and passes the token to its immediate successor when it finishes committing. This approach cannot scale to the commit throughput that Swarm needs. For example, with 100-cycle tasks, a 64-core system should commit 0.64 tasks/cycle on average. Even if commits were instantaneous, the latency incurred by passing the token makes this throughput unachievable.

To tackle this problem, we show that, by adapting techniques from distributed systems, we can achieve in-order commits without serialization, token-passing, or building successor lists.

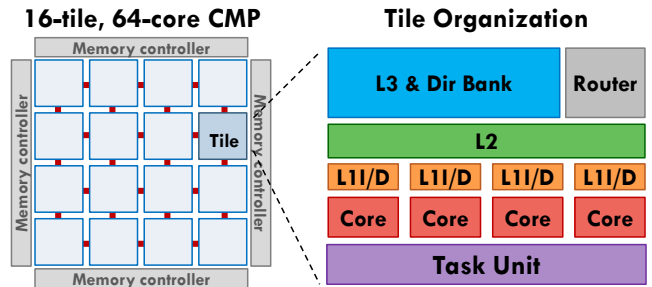


Figure 2: Swarm CMP and tile configuration.

4. SWARM: AN ARCHITECTURE FOR ORDERED PARALLELISM

Fig. 2 shows Swarm’s high-level organization. Swarm is a tiled, cache-coherent chip multiprocessor (CMP). Each tile has a group of simple cores. Each core has small, private, write-through L1 caches. All cores in a tile share a per-tile L2 cache, and each tile has a slice of a shared NUCA L3 cache. Each tile features a *task unit* that queues, dispatches, and commits tasks. Tiles communicate through a mesh NoC.

Key features: Swarm is optimized to execute short tasks with programmer-specified order constraints. Programmers define the execution order by assigning *timestamps* to tasks. Tasks can create children tasks with equal or later timestamps than their own. Tasks appear to execute in global timestamp order, but Swarm uses speculation to elide order constraints.

Swarm is coherently designed to support a large speculative task window efficiently. Swarm has no centralized structures: each tile’s task unit queues runnable tasks and maintains the speculative state of finished tasks that cannot yet commit. Task units only communicate when they send new tasks to each other to maintain load balance, and, infrequently, to determine which finished tasks can be committed.

Swarm speculates far ahead of the earliest active task, and runs tasks even if their parent is still speculative. Fig. 3(a) shows this process: a task with timestamp 0 is still running, but tasks with later timestamps and several speculative ancestors are running or have finished execution. For example, the task with timestamp 51, currently running, has three still-speculative ancestors, two of which have finished and are waiting to commit (8 and 20) and one that is still running (40).

Allowing tasks with speculative ancestors to execute uncovers significant parallelism, but may induce aborts that span multiple tasks. For example, in Fig. 3(b) a new task with timestamp 35 conflicts with task 40, so 40 is aborted and child task 51 is both aborted and discarded. These aborts are *selective*, and only affect tasks whose speculative ancestors are aborted, or tasks that have read data written by an aborted task.

We describe Swarm in a layered fashion. First, we present Swarm’s ISA extensions. Second, we describe Swarm hardware assuming that *all queues are unbounded*. Third, we discuss how Swarm handles bounded queue sizes. Fourth, we present Swarm’s hardware costs.

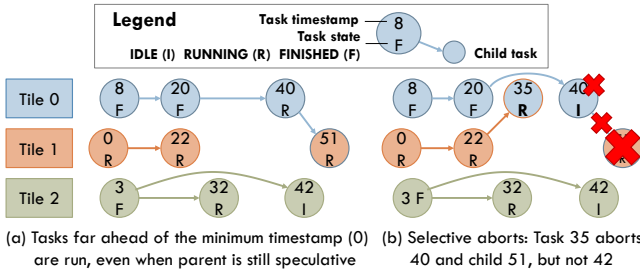


Figure 3: Example execution of sssp. By executing tasks even if their parents are speculative, Swarm uncovers ordered parallelism, but may trigger selective aborts.

4.1 ISA Extensions and Programming Model

Swarm manages and dispatches tasks using hardware task queues. A task is represented by a descriptor with the following architectural state: the task’s function pointer, a 64-bit timestamp, and the task’s arguments.

Tasks appear to run in timestamp order. Tasks with the same timestamp may execute in any order, but run *atomically*—the system lazily selects an order for them.

A task can create one or more *children tasks* with an equal or later timestamp than its own. A child is ordered after its *parent*, but children with the same timestamp may execute in any order. Because hardware must track parent-child relations, tasks may create a limited number of children (8 in our implementation). Tasks that need more children enqueue a single task that creates them.

Swarm adds instructions to enqueue and dequeue tasks. The `enqueue_task` instruction accepts a task descriptor (held in registers) as its input and queues the task for execution. A thread uses the `dequeue_task` instruction to start executing a previously-enqueued task. `dequeue_task` initiates speculative execution at the task’s function pointer and makes the task’s timestamp and arguments available (in registers). Task execution ends with a `finish_task` instruction.

`dequeue_task` stalls the core if an executable task is not immediately available, avoiding busy-waiting. When no tasks are left in any task unit and all threads are stalled on `dequeue_task`, the algorithm has terminated, and `dequeue_task` jumps to a configurable pointer to handle termination.

API: We design a low-level C++ API that uses these mechanisms. Tasks are simply functions with signature:

```
void taskFn(timestamp, args...)
```

Code can enqueue other tasks by calling:

```
enqueueTask(taskFn, timestamp, args...)
```

If a task needs more than the maximum number of task descriptor arguments, three 64-bit words in our implementation, the runtime allocates them in memory.

4.2 Task Queuing and Prioritization

The task unit has two main structures:

1. The *task queue* holds task descriptors (function pointer, timestamp, and arguments).
2. The *commit queue* holds the speculative state of tasks that have finished execution but cannot yet commit.

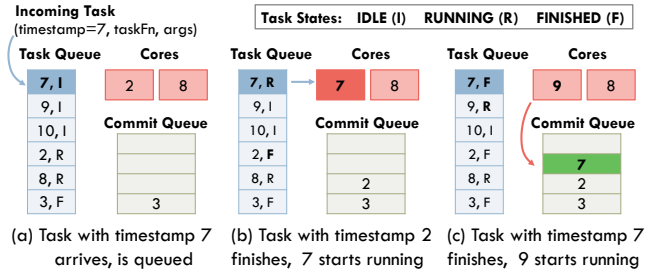


Figure 4: Task queue and commit queue utilization through a task’s lifetime.

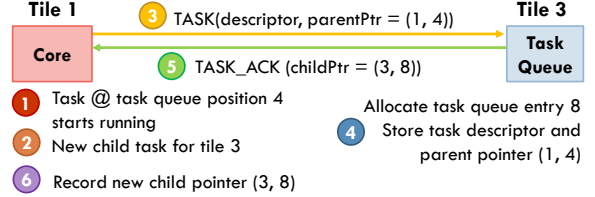


Figure 5: Task creation protocol. Cores send new tasks to other tiles for execution. To track parent-child relations, parent and child keep a pointer to each other.

Fig. 4 shows how these queues are used throughout the task’s lifetime. Each new task allocates a task queue entry, and holds it until commit time. Each task allocates a commit queue entry when it finishes execution, and also deallocates it at commit time. For now, assume these queues always have free entries. Sec. 4.7 discusses what happens when they fill up.

Together, the task queue and commit queue are similar to a reorder buffer, but at task-level rather than instruction-level. They are separate structures because commit queue entries are larger than task queue entries, and typically fewer tasks are waiting to commit than to execute. However, unlike in a reorder buffer, tasks do not arrive in priority order. Both structures manage their free space with a freelist and allocate entries independently of task priority order, as shown in Fig. 4.

Task enqueues: When a core creates a new task (through `enqueue_task`), it sends the task to a randomly-chosen target tile following the protocol in Fig. 5. Parent and child track each other using *task pointers*. A task pointer is simply the tuple (*tile, task queue position*). This tuple uniquely identifies a task because it stays in the same task queue position throughout its lifetime.

Task prioritization: Tasks are prioritized for execution in timestamp order. When a core calls `dequeue_task`, the highest-priority idle task is selected for execution. Since task queues do not hold tasks in priority order, an auxiliary *order queue* is used to find this task.

The order queue can be cheaply implemented with two small ternary content-addressable memories (TCAMs) with as many entries as the task queue (e.g., 256), each of which stores a 64-bit timestamp. With Panigrahy and Sharma’s PIDR_OPT method [54], finding the next task to dispatch requires a single lookup in both TCAMs, and each insertion (task creation) and deletion (task commit or squash) requires two lookups in both TCAMs. SRAM-

based implementations are also possible, but we find the small TCAMs to have a moderate cost (Sec. 4.8).

4.3 Speculative Execution and Versioning

The key requirements for speculative execution in Swarm are allowing fast commits and a large speculative window. To this end, we adopt *eager versioning*, storing speculative data in place and logging old values. Eager versioning makes commits fast, but aborts are slow. However, Swarm’s execution model makes conflicts rare, so eager versioning is the right tradeoff.

Eager versioning is common in hardware transactional memories [30, 48, 79], which do not perform ordered execution or speculative data forwarding. By contrast, most TLS systems use lazy versioning (buffering speculative data in caches) or more expensive multiversioning [11, 25, 28, 29, 56, 60, 61, 66, 68, 69] to limit the cost of aborts. Some early TLS schemes are eager [25, 80], and they still suffer from the limitations discussed in Sec. 3.

Swarm’s speculative execution borrows from LogTM and LogTM-SE [48, 63, 79]. Our key contributions over these and other speculation schemes are (i) conflict detection (Sec. 4.4) and selective abort techniques (Sec. 4.5) that leverage Swarm’s hierarchical memory system and Bloom filter signatures to scale to large speculative windows, and (ii) a technique that exploits Swarm’s large commit queues to achieve high-throughput commits (Sec. 4.6).

Fig. 6 shows the per-task state needed to support speculation: read- and write-set signatures, an undo log pointer, and child pointers. Each core and commit queue entry holds this state.

A successful `dequeue_task` instruction jumps to the task’s code pointer and initiates speculation. Since speculation happens at the task level, there are no register checkpoints, unlike in HTM and TLS. Like in LogTM-SE, as the task executes, hardware automatically performs conflict detection on every read and write (Sec. 4.4). Then, it inserts the read and written addresses into the Bloom filters, and, for every write, it saves the old memory value in a memory-resident undo log. Stack addresses are not conflict-checked or logged.

When a task finishes execution, it allocates a commit queue entry; stores the read and write set signatures, undo log pointer, and children pointers there; and frees the core for another task.

4.4 Virtual Time-Based Conflict Detection

Conflict detection is based on a priority order that respects both programmer-assigned timestamps and parent-child relationships. Conflicts are detected at cache line granularity.

Unique virtual time: Tasks may have the same programmer-assigned timestamp. However, conflict detection has much simpler rules if tasks follow a total order. Therefore, tasks are assigned a *unique virtual time* when they are dequeued for execution. Unique virtual time is the 128-bit tuple (*programmer timestamp*, *dequeue cycle*, *tile id*). The (*dequeue cycle*, *tile id*) pair is unique since at most one dequeue per cycle is permitted at a

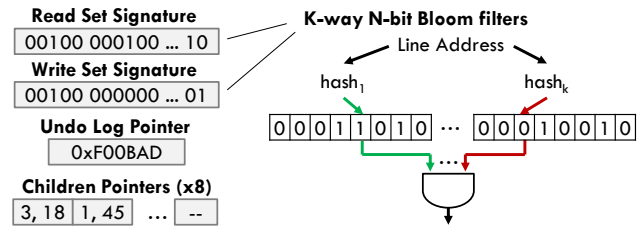


Figure 6: Speculative state for each task. Each core and commit queue entry maintains this state. Read and write sets are implemented with space-efficient Bloom filters.

tile. Conflicts are resolved using this unique virtual time, which tasks preserve until they commit.

Unique virtual times incorporate the ordering needs of programmer-assigned timestamps and parent-child relations: children always start execution after their parents, so a parent always has a smaller dequeue cycle than its child, and thus a smaller unique virtual time, even when parent and child have the same timestamp.

Conflicts and forwarding: Conflicts arise when a task accesses a line that was previously accessed by a later-virtual time task. Suppose two tasks, t_1 and t_2 , are running or finished, and t_2 has a later virtual time. A read of t_1 to a line written by t_2 or a write to a line read or written by t_2 causes t_2 to abort. However, t_2 can access data written by t_1 even if t_1 is still speculative. Thanks to eager versioning, t_2 automatically uses the latest copy of the data—there is no need for speculative data forwarding logic [25].

Hierarchical conflict detection: Swarm exploits the cache hierarchy to reduce conflict checks. Fig. 7 shows the different types of checks performed in an access:

1. The L1 is managed as described below to ensure L1 hits are conflict-free.
2. L1 misses are checked against other tasks in the tile (both in other cores and in the commit queue).
3. L2 misses, or L2 hits where a virtual time check (described below) fails, are checked against tasks in other tiles. As in LogTM [48], the L3 directory uses memory-backed *sticky bits* to only check tiles whose tasks may have accessed the line. Sticky bits are managed exactly as in LogTM.

Any of these conflicts trigger task aborts.

Using caches to filter checks: The key invariant that allows caches to filter checks is that, when a task with virtual time T installs a line in the (L1 or L2) cache, that line has no conflicts with tasks of virtual time $> T$. As long as the line stays cached with the right coherence permissions, it stays conflict-free. Because conflicts happen when tasks access lines out of virtual time order, if another task with virtual time $U > T$ accesses the line, it is also guaranteed to have no conflicts.

However, accesses from a task with virtual time $U < T$ must trigger conflict checks, as another task with intermediate virtual time X , $U < X < T$, may have accessed the line. U ’s access does not conflict with T ’s, but may conflict with X ’s. For example, suppose a task with virtual time $X = 2$ writes line A . Then, task $T = 3$ in another core reads A . This is not a conflict with X ’s

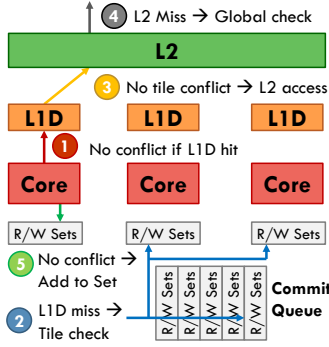


Figure 7: Local, tile, and global conflict detection for an access that misses in the L1 and L2.

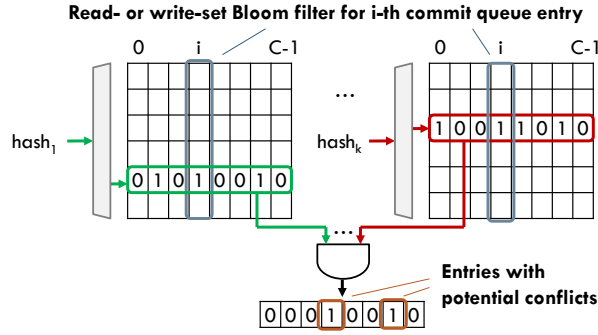


Figure 8: Commit queues store read- and write-set Bloom filters by columns, so a single access reads bit from all entries. All entries are checked in parallel.

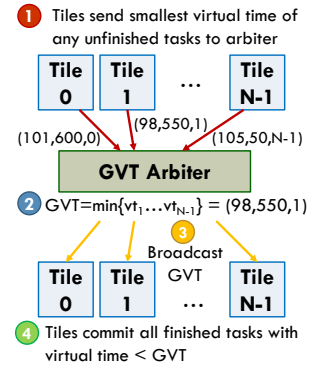


Figure 9: Global virtual time commit protocol.

write, so A is installed in T 's L1. The core then finishes T and dequeues a task $U = 1$ that reads A . Although A is in the L1, U has a conflict with X 's write.

We handle this issue with two changes. First, when a core dequeues a task with a smaller virtual time than the one it just finished, it flushes the L1. Because L1s are small and write-through, this is fast, simply requiring to flash-clear the valid bits. Second, each L2 line has an associated *canary virtual time*, which stores the lowest task virtual time that need not perform a global check. For efficiency, lines in the same L2 set share the same canary virtual time. For simplicity, this is the maximum virtual time of the tasks that installed each of the lines in the set, and is updated every time a line is installed. *Efficient commit queue checks:* Although caches reduce the frequency of conflict checks, all tasks in the tile must be checked on every L2 access and on some global checks. To allow large commit queues (e.g., 64 tasks/queue), commit queue checks must be efficient. To this end, we leverage that checking a K -way Bloom filter only requires reading one bit from each way. As shown in Fig. 8, Bloom filter ways are stored in columns, so a single 64-bit access per way reads all the necessary bits. Reading and ANDing all ways yields a word that indicates potential conflicts. For each queue entry whose position in this word is set, its virtual time is checked; those with virtual time higher than the issuing task's must be aborted.

4.5 Selective Aborts

Upon a conflict, Swarm aborts the later task and all its dependents: its children and other tasks that have accessed data written by the aborting task. Hardware aborts each task t in three steps:

1. Notify t 's children to abort and be removed from their task queues.
2. Walk t 's undo log in LIFO order, restoring old values. If one of these writes conflicts with a later-virtual time task, wait for it to abort and continue t 's rollback.
3. Clear t 's signatures and free its commit queue entry.

Applied recursively, this procedure selectively aborts all dependent tasks, as shown in Fig. 10. This scheme has two key benefits. First, it reuses the conflict-detection logic used in normal operation. Undo-log writes (e.g., A 's second `wr 0x10` in Fig. 10) are normal conflict-checked

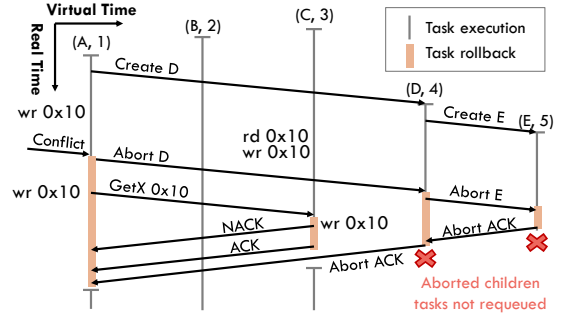


Figure 10: Selective abort protocol. Suppose $(A, 1)$ must abort after it writes `0x10`. $(A, 1)$'s abort squashes child $(D, 4)$ and grandchild $(E, 5)$. During rollback, A also aborts $(C, 3)$, which read A 's speculative write to `0x10`. $(B, 2)$ is independent and thus not aborted.

writes, issued with the task's timestamp to detect all later readers and writers. Second, this scheme does not explicitly track data dependences among tasks. Instead, it uses the conflict-detection protocol to recover them as needed. This is important, because any task may have served speculative data to many other tasks, which would make explicit tracking expensive. For example, tracking all possible dependences on a 1024-task window using bit-vectors, as proposed in prior work [13, 58], would require $1024 \times 1023 \approx 1$ Mbit of state.

4.6 Scalable Ordered Commits

To achieve high-throughput commits, Swarm adapts the virtual time algorithm [38], common in parallel discrete event simulation [21]. Fig. 9 shows this protocol. Tiles periodically send the smallest unique virtual time of any unfinished (running or idle) task to an arbiter. Idle tasks do not yet have a unique virtual time and use $(timestamp, current\ cycle, tile\ id)$ for the purposes of this algorithm. The arbiter computes the minimum virtual time of all unfinished tasks, called the *global virtual time* (GVT), and broadcasts it to all tiles. To preserve ordering, only tasks with virtual time $<$ GVT can commit.

The key insight is that, by combining the virtual time algorithm with Swarm's large commit queues, *commit costs are amortized over many tasks*. A single

GVT update often causes many finished tasks to commit. For example, if in Fig. 9 the GVT jumps from (80,100,2) to (98,550,1), all tasks with virtual time (80,100,2) < t < (98,550,1) can commit. GVT updates happen sparingly (e.g., every 200 cycles) to limit bandwidth. Less frequent updates reduce bandwidth but increase commit queue occupancy.

In addition, eager versioning makes commits fast: a task commits by freeing its task and commit queue entries, a single-cycle operation. Thus, if a long-running task holds the GVT for some time, once it finishes, commit queues quickly drain and catch up to execution.

Compared with prior TLS schemes that use successor lists and token passing to reconcile order (Sec. 3), this scheme does not even require finding the successor and predecessor of each task, and does not serialize commits.

For the system sizes we evaluate, a single GVT arbiter suffices. Larger systems may need a hierarchy of arbiters.

4.7 Handling Limited Queue Sizes

The per-tile task and commit queues may fill up, requiring a few simple actions to ensure correct operation.

Task queue virtualization: Applications may create an unbounded number of tasks and schedule them for a future time. Swarm uses an overflow/underflow mechanism to give the illusion of unbounded hardware task queues [27, 41, 64]. When the per-tile task queue is nearly full, the task unit dispatches a special, non-speculative *coalescer* task to one of the cores. This *coalescer* task removes several *non-speculative*, idle task descriptors with high programmer-assigned timestamps from the task queue, stores them in memory, and enqueues a *splitter* task that will re-enqueue the overflowed tasks.

Note that only non-speculative task queue entries can be moved to software. These (*i*) are idle, and (*ii*) have no parent (i.e., their parent has already committed). When all entries are speculative, we need another approach.

Virtual time-based allocation: The task and commit queues may also fill up with speculative tasks. The general rule to avoid deadlock due to resource exhaustion is to always prioritize earlier-virtual time tasks, aborting other tasks with later virtual times if needed. For example, if a tile speculates far ahead, fills up its commit queue, and then receives a task that precedes all other speculative tasks, the tile must let the preceding task execute to avoid deadlock. This results in three specific policies for the commit queue, cores, and task queue.

Commit queue: If task t finishes execution, the commit queue is full, and t precedes any of the tasks in the commit queue, it aborts the highest-virtual time finished task and takes its commit queue entry. Otherwise, t stalls its core, waiting for an entry.

Cores: If task t arrives at the task queue, the commit queue is full, and t precedes all tasks in cores, t aborts the highest-virtual time task and takes its core.

Task queue: Suppose task t arrives at a task unit but the task queue is full. If some tasks are non-speculative, then a *coalescer* is running, so the task waits for a free entry. If all tasks in the task queue are speculative, the enqueued request is NACK'd (instead of ACK'd as in

| | Entries | Entry size | Size | Est. area |
|-----------------------------|---------|------------|----------------|-----------------------|
| Task queue | 256 | 51 B | 12.75 KB | 0.056 mm ² |
| Commit filters queue | filters | 16×32 B | 32 KB (2-port) | 0.304 mm ² |
| | other | 36 B | 2.25 KB | 0.012 mm ² |
| Order queue | 256 | 2×8 B | 4 KB (TCAM) | 0.175 mm ² |

Table 2: Sizes and estimated areas of main task unit structures.

Fig. 5) and the *parent task* stalls, and retries the enqueue using linear backoff. To avoid deadlock, we leverage that when a task's unique virtual time matches the GVT, it is the smallest-virtual time task in the system, and cannot be aborted. This task need not keep track of its children (no child pointers), and when those children are sent to another tile, they can be overflowed to memory if the task queue is full. This ensures that the GVT task makes progress, avoiding deadlock.

4.8 Analysis of Hardware Costs

We now describe Swarm's overheads. Swarm adds task units, a GVT arbiter, and modifies cores and L2s.

Table 2 shows the per-entry sizes, total queue sizes, and area estimates for the main task unit structures: task queue, commit queue, and order queue. All numbers are for one per-tile task unit. We assume a 16-tile, 64-core system as in Fig. 2, with 256 task queue entries (64 per core) and 64 commit queue entries (16 per core). We use CACTI [70] for the task and commit queue SRAM areas (using 32 nm ITRS-HP logic) and scaled numbers from a commercial 28 nm TCAM [3] for the order queue area. Task queues use single-port SRAMs. Commit queues use several dual-port SRAMs for the Bloom filters (Fig. 8), which are 2048-bit, 8-way in our implementation, and a single-port SRAM for all other state (unique virtual time, undo log pointer, and child pointers).

Overall, these structures consume 0.55 mm² per 4-core tile, or 8.8 mm² per chip, a minor cost. Enqueues and dequeues access the order queue TCAM, which consumes ~70pJ per access [50]. Moreover, queue operations happen sparingly (e.g. with 100-cycle tasks, one enqueue and dequeue every 25 cycles), so energy costs are small.

The GVT arbiter is simple. It buffers a virtual time per tile, and periodically broadcasts the minimum one.

Cores are augmented with **enqueue/dequeue/finish-task** instructions (Sec. 4.1), the speculative state in Fig. 6 (530 bytes), a 128-bit unique virtual time, and logic to insert addresses into Bloom filters and to, on each store, write the old value to an undo log. Finally, the L2 uses a 128-bit canary virtual time per set. For an 8-way cache with 64 B lines, this adds 2.6% extra state.

In summary, Swarm's costs are moderate, and, in return, confer significant speedups.

5. EXPERIMENTAL METHODOLOGY

Modeled system: We use an in-house microarchitectural, event-driven, sequential simulator based on Pin [46] to model a 64-core CMP with a 3-level cache hierarchy. We use simple IPC-1 cores with detailed timing models for caches, on-chip network, and main memory (adapted from zsim [62]), and also model Swarm

| | |
|---------------------|---|
| Cores | 64 cores in 16 tiles (4 cores/tile), 2 GHz, x86-64 ISA, IPC-1 except misses and Swarm instructions |
| L1 caches | 16 KB, per-core, split D/I, 8-way, 2-cycle latency |
| L2 caches | 256 KB, per-tile, 8-way, inclusive, 7-cycle latency |
| L3 cache | 16 MB, shared, static NUCA [40] (1 MB bank/tile), 16-way, inclusive, 9-cycle bank latency |
| Coherence | MESI, 64 B lines, in-cache directories, no silent drops |
| NoC | 4×4 mesh, 256-bit links, X-Y routing, 3 cycles/hop |
| Main mem | 4 controllers at chip edges, 120-cycle latency |
| Queues | 64 task queue entries/core (4096 total), 16 commit queue entries/core (1024 total) |
| Swarm instrs | 5 cycles per enqueue/dequeue/finish_task |
| Conflicts | 2048-bit 8-way Bloom filters, H_3 hash functions [10] Tile checks take 5 cycles (Bloom filters) + 1 cycle for every timestamp compared in the commit queue |
| Commits | Tiles and GVT arbiter send updates every 200 cycles |
| Spills | Coalescers fire when a task queue is 75% full Coalescers spill up to 15 tasks each |

Table 3: Configuration of the 64-core CMP.

| | Software baselines | Input | Seq run-time |
|--------------|-----------------------|--------------------------|--------------|
| bfs | PBFS [43] | hugetric-00020 [5, 16] | 3.68 Bcycles |
| sssp | Bellman-Ford [33, 55] | East USA roads [1] | 4.42 Bcycles |
| astar | Own | Germany roads [53] | 2.08 Bcycles |
| msf | PBBS [7] | kronecker_logn16 [5, 16] | 2.16 Bcycles |
| des | Chandy-Misra [33, 55] | csaArray32 [55] | 3.05 Bcycles |
| silos | Silo [71] | TPC-C, 4 whs, 32 Ktxns | 2.93 Bcycles |

Table 4: Benchmark information: source of baseline implementations, inputs, and run-time of the serial version.

features (e.g., conflict checks, aborts, etc.) in detail. Table 3 details the modeled configuration.

Benchmarks: We use the six benchmarks mentioned in Sec. 2.2: **bfs**, **sssp**, **astar**, **msf**, **des**, and **silos**. Table 4 details their provenance and input sets.

For most benchmarks, we use tuned serial and state-of-the-art parallel versions from existing suites (Table 4). We then port each serial implementation to Swarm. Swarm versions use fine-grain tasks, but use the same data structures and perform the same work as the serial version, so differences between serial and Swarm versions stem from parallelism, not other optimizations.

We wrote our own tuned serial and Swarm **astar** implementations. **astar** is notoriously difficult to parallelize—to scale, prior work in parallel pathfinding sacrifices solution quality for speed [8]. Thus, we do not have a software-only parallel implementation.

We port **silos** to show that Swarm can extract ordered parallelism from applications that are typically considered unordered. Database transactions are unordered in **silos**. We decompose each transaction into many small ordered tasks to exploit *intra-transaction parallelism*. Tasks from different transactions use disjoint timestamp ranges to preserve atomicity. This exposes significant fine-grain parallelism within and across transactions.

Input sets: We use a varied set of inputs, often from standard collections such as DIMACS (Table 4). **bfs** operates on an unstructured mesh; **sssp** and **astar** use large road maps; **msf** uses a Kronecker graph; **des** simulates an array of carry-select adders; and **silos** runs the TPC-C benchmark on 4 warehouses.

All benchmarks have serial run-times of over two billion cycles (Table 4). We have evaluated other inputs (e.g., random and scale-free graphs), and qualitative dif-

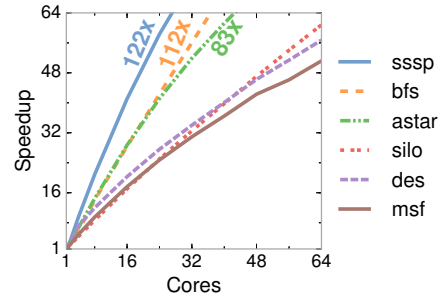


Figure 11: Swarm self-relative speedups on 1-64 cores. Larger systems have larger queues and caches, which affect speedups and sometimes cause superlinear scaling.

ferences are not affected. Note that some inputs can offer plentiful trivial parallelism to a software algorithm. For example, on large, shallow graphs (e.g., 10 M nodes and 10 levels), a simple bulk-synchronous **bfs** that operates on one level at a time scales well [43]. But we use a graph with 7.1 M nodes and 2799 levels, so **bfs** must speculate across levels to uncover enough parallelism.

For each benchmark, we fast-forward to the start of the parallel region (skipping initialization), and report results for the full parallel region.

Idealized memory allocation: Dynamic memory allocation is not simulated in detail, and a scalable solution is left to future work. Only **des** and **silos** tasks allocate memory frequently, and data dependences in the system’s memory allocator serialize them. In principle, we could build a task-aware allocator with per-core memory pools to avoid serialization. However, building high-performance allocators is complex [26, 65]. Instead, the simulator allocates and frees memory in a task-aware way. Freed memory is not reused until the freeing task commits to avoid spurious dependences. Each allocator operation incurs a 30-cycle cost. For fairness, *serial and software-parallel implementations also use this allocator*. We believe this simplification will not significantly affect **des** and **silos** results when simulated in detail.

6. EVALUATION

We first compare Swarm with alternative implementations, then analyze its behavior in depth.

6.1 Swarm Scalability

Fig. 11 shows Swarm’s performance on 1- to 64-core systems. In this experiment, *per-core* queue and L2/L3 capacities are kept constant as the system grows, so *systems with more cores have higher queue and cache capacities*. This captures performance per unit area.

Each line in Fig. 11 shows the speedup of a single application over a 1-core system (i.e., its self-relative speedup). At 64 cores, speedups range from 51× (**msf**) to 122× (**sssp**), demonstrating high scalability. In addition to parallelism, the larger queues and L3 of larger systems also affect performance, causing super-linear speedups in some benchmarks (**sssp**, **bfs**, and **astar**). We tease apart the contribution of these factors in Sec. 6.3.

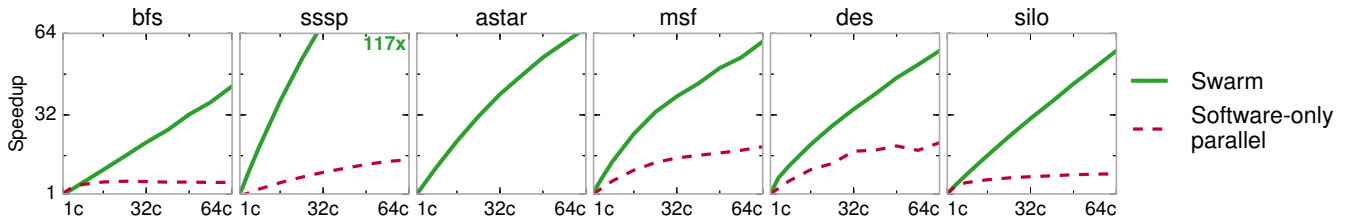


Figure 12: Speedup of Swarm and state-of-the-art software-parallel implementations from 1 to 64 cores, relative to a tuned serial implementation running on a system of the same size.

6.2 Swarm vs Software Implementations

Fig. 12 compares the performance of the Swarm and software-only versions of each benchmark. Each graph shows the speedup of the Swarm and software-parallel versions over the tuned serial version running on a system of the same size, from 1 to 64 cores. As in Fig. 11, queue and L2/L3 capacities scale with the number of cores.

Swarm outperforms the serial versions by 43–117 \times , and the software-parallel versions by 2.7–18.2 \times . We analyze the reasons for these speedups for each application. **bfs**: Serial **bfs** does not need a priority queue. It uses an efficient FIFO queue to store the set of nodes to visit. At 1 core, Swarm is 33% slower than serial **bfs**; however, Swarm scales to 43 \times at 64 cores. By contrast, the software-parallel version, PBFS [43], scales to 6.0 \times , then slows down beyond 24 cores. PBFS only works on a single level of the graph at a time, while Swarm speculates across multiple levels.

sssp: Serial **sssp** uses a priority queue. Swarm is 32% faster at one core, and 117 \times faster at 64 cores. The software-parallel version uses the Bellman-Ford algorithm [15]. Bellman-Ford visits nodes out of order to increase parallelism, but wastes work in doing so. Threads in Bellman-Ford communicate infrequently to limit overheads [33], wasting much more work than Swarm’s speculative execution. As a result, Bellman-Ford **sssp** scales to 14 \times at 64 cores, 8.1 \times slower than Swarm.

astar: Our tuned serial **astar** uses a priority queue to store tasks [15]. Swarm outperforms it by 2% at one core, and by 66 \times at 64 cores.

msf: The serial and software-parallel **msf** versions sort edges by weight to process them in order. Our Swarm implementation instead does this sort implicitly through the task queues, enqueueing one task per edge and using its weight as the timestamp. This allows Swarm to overlap the sort and edge-processing phases. Swarm outperforms the serial version by 70% at one core and 61 \times at 64 cores. The software-parallel **msf** uses software speculation via deterministic reservations [7], and scales to 19 \times at 64 cores, 3.1 \times slower than Swarm.

des: Serial **des** uses a priority queue to simulate events in time order. Swarm outperforms the serial version by 23% at one core, and by 57 \times at 64 cores. The software-parallel version uses the Chandy-Misra-Bryant (CMB) algorithm [47, 67]. CMB exploits the simulated communication latencies among components to safely execute some events out of order (e.g., if two nodes have a 10-cycle simulated latency, they can be simulated up to 9 cycles away). CMB scales to 21 \times at 64 cores, 2.7 \times

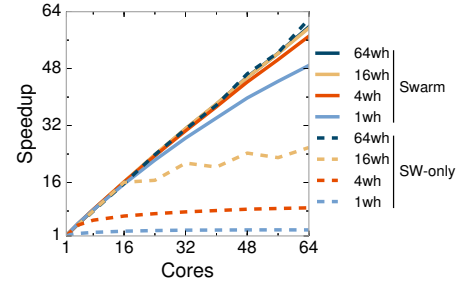


Figure 13: Speedup of Swarm and software silo with 64, 16, 4, and 1 TPC-C warehouses.

slower than Swarm. Half of Swarm’s speedup comes from exploiting speculative parallelism, and the other half from reducing overheads.

silo: Serial **silo** runs database transactions sequentially without synchronization. Swarm outperforms serial **silo** by 10% at one core, and by 57 \times at 64 cores. The software-parallel version uses a carefully optimized protocol to achieve high transaction rates [71]. Software-parallel **silo** scales to 8.8 \times at 64 threads, 6.4 \times slower than Swarm. The reason is fine-grain parallelism: in Swarm, each task reads or writes at most one tuple. This exposes parallelism within and across database transactions, and reduces the penalty of conflicts, as only small, dependent tasks are aborted instead of full transactions.

Swarm’s benefits on **silo** heavily depend on the amount of coarse-grain parallelism, which is mainly determined by the number of TPC-C warehouses. To quantify this effect, Fig. 13 shows the speedups of Swarm and software-parallel **silo** with 64, 16, 4, and 1 warehouses. With 64 warehouses, software-parallel **silo** scales linearly up to 64 cores and is 4% faster than Swarm. With fewer warehouses, database transactions abort frequently, limiting scalability. With a single warehouse, software-parallel **silo** scales to only 2.7 \times . By contrast, Swarm exploits fine-grain parallelism within each transaction, and scales well even with a single warehouse, by 49 \times at 64 cores, 18.2 \times faster than software-parallel **silo**.

Overall, these results show that Swarm outperforms a wide range of parallel algorithms, even when they use application-specific optimizations. Moreover, Swarm implementations use no explicit synchronization and are simpler, which is itself valuable.

6.3 Swarm Analysis

We now analyze the behavior of different benchmarks in more detail to gain insights about Swarm.

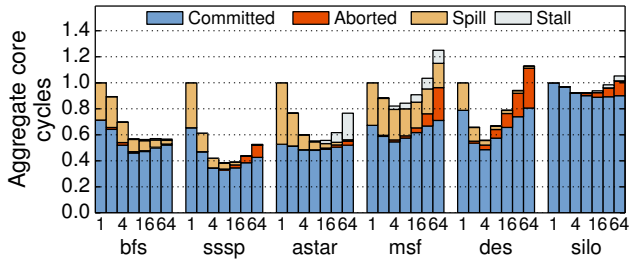


Figure 14: Breakdown of total core cycles for Swarm systems with 1 to 64 cores. Most time is spent executing tasks that are ultimately committed.

| Speedups | 1c vs 1c-base | 64c vs 1c-base | 64c vs 1c |
|----------------------|---------------|----------------|-----------|
| Swarm baseline | 1× | 77× | 77× |
| + unbounded queues | 1.4× | 87× | 61× |
| + 0-cycle mem system | 5× | 274× | 54× |

Table 5: gmean speedups with progressive idealizations: unbounded queues and a zero-cycle memory system (1c-base = 1-core Swarm baseline without idealizations).

Cycle breakdowns: Fig. 14 shows the breakdown of aggregate core cycles. Each set of bars shows results for a single application as the system scales from 1 to 64 cores. The height of each bar is the sum of cycles spent by all cores, normalized by the cycles of the 1-core system (lower is better). With linear scaling, all bars would have a height of 1.0; higher and lower bars indicate sub- and super-linear scaling, respectively. Each bar shows the breakdown of cycles spent executing tasks that are ultimately committed, tasks that are later aborted, spilling tasks from the hardware task queue (using coalescer and splitter tasks, Sec. 4.7), and stalled.

Swarm spends most of the cycles executing tasks that later commit. At 64 cores, aborted work ranges from 1% (**bfs**) to 27% (**des**) of cycles. All graph benchmarks spend significant time spilling tasks to memory, especially with few cores (e.g., 47% of cycles for single-core **astar**). In all benchmarks but **msf**, spill overheads shrink as the system grows and task queue capacity increases; **msf** enqueues millions of edges consecutively, so larger task queues do not reduce spills. Finally, cores rarely stall due to full or empty queues. Only **astar** and **msf** spend more than 5% of cycles stalled at 64 cores: 27% and 8%, respectively.

Fig. 14 also shows the factors that contribute to super-linear scaling in Fig. 11. First, larger task queues can capture a higher fraction of runnable tasks, reducing spills. Second, larger caches can better fit the working set, reducing the cycles spent executing committed tasks (e.g., **sil**). However, beyond 4–8 cores, the longer hit latency of the larger NUCA L3 counters its higher hit rate in most cases, increasing execution cycles.

Speedups with idealizations: To factor out the impact of queues and memory system on scalability, we consider systems with two idealizations: unbounded queues, which factor out task spills, and an ideal memory system with 0-cycle delays for all accesses and messages. Table 5 shows the gmean speedups when these idealizations

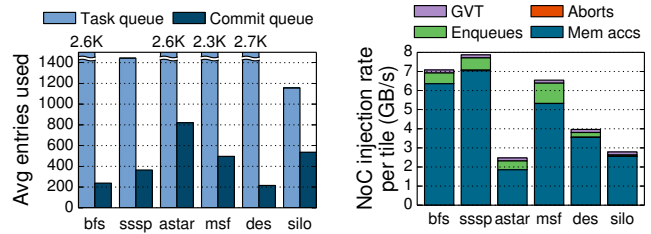


Figure 15: Average task and commit queue occupancies for 64-core Swarm.

Figure 16: Breakdown of NoC traffic per tile for 64-core, 16-tile Swarm.

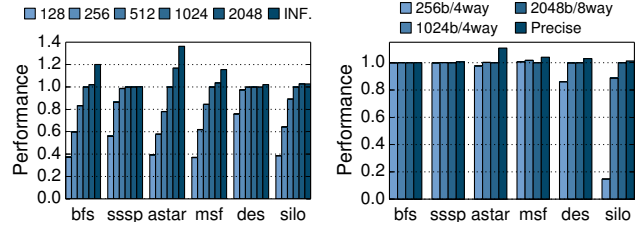


Figure 17: Sensitivity of 64-core Swarm to commit queue and Bloom filter sizes.

are progressively applied. The left and middle columns show 1- and 64-core speedups, respectively, over the 1-core baseline (without idealizations). While idealizations help both cases, they have a larger impact on the 1-core system. Therefore, the 64-core speedups relative to the 1-core system *with the same idealizations* (right column) are lower. With all idealizations, this speedup is purely due to exploiting parallelism; 64-core Swarm is able to mine 54× parallelism on average (46×–63×).

Queue occupancies: Fig. 15 shows the average number of task queue and commit queue entries used across the 64-core system. Both queues are often highly utilized. Commit queues can hold up to 1024 finished tasks (64 per tile). On average, they hold from 216 in **des** to 821 in **astar**. This shows that cores often execute tasks out of order, and these tasks wait a significant time until they commit—a large speculative window is crucial, as the analysis in Sec. 2.2 showed. The 4096-entry task queues are also well utilized, with average occupancies between 1157 (**sil**) and 2712 (**msf**) entries.

Network traffic breakdown: Fig. 16 shows the NoC traffic breakdown at 64 cores (16 tiles). The cumulative injection rate per tile remains well below the saturation injection rate (32 GB/s). Each bar shows the contributions of memory accesses (between the L2s and L3) issued during normal execution, tasks enqueues to other tiles, abort traffic (including child abort messages and rollback memory accesses), and GVT updates. Task enqueues, aborts, and GVT updates increase network traffic by 15% on average. Thus, Swarm imposes small overheads on traffic and communication energy.

Conflict detection energy: Conflict detection requires Bloom filter checks—performed in parallel over commit queue entries (Fig. 7)—and for those entries where the Bloom filter reports a match, a virtual time check to see whether the task needs to be aborted. Both

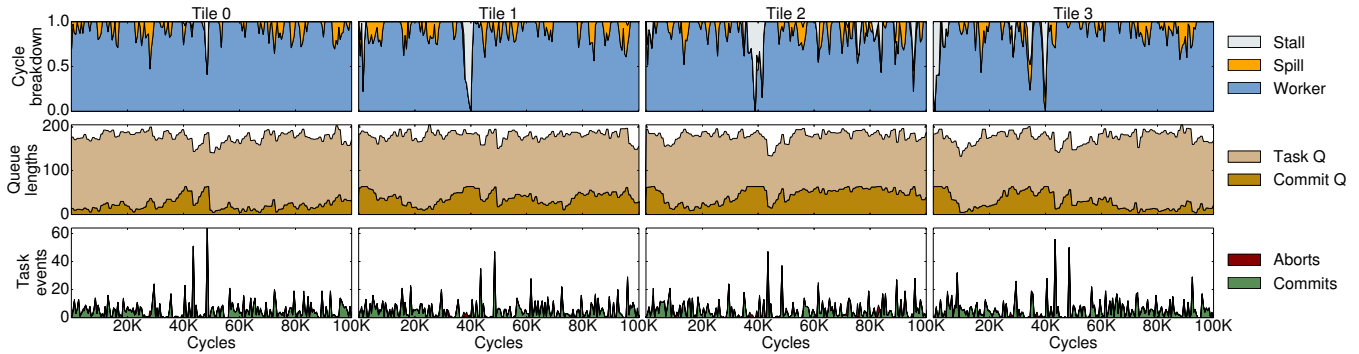


Figure 18: Execution trace of astar on 16-core (4-tile) Swarm over a 100 Kcycle interval: breakdown of core cycles (top), queue lengths (middle), and task commits and aborts (bottom) for each tile.

events happen relatively rarely. Each tile performs one Bloom filter check every 8.0 cycles on average (from 2.5 cycles in *msf* to 13 cycles in *bfs*). Each tile performs one timestamp check every 49 cycles on average (from 6 cycles in *msf* to 143 cycles in *astar*). Hence, Swarm’s conflict detection imposes acceptable energy overheads. **Canary virtual times:** To lower overheads, all lines in the same L2 set share a common canary virtual time. This causes some unnecessary global conflict checks, but we find the falsely unfiltered checks are infrequent. At 64 cores, using precise per-line canary virtual times reduces global conflict checks by 10.3% on average, and improves application performance by less than 1%.

6.4 Sensitivity Studies

We explore Swarm’s sensitivity to several design parameters at 64 cores:

Commit queue size: Fig. 17(a) shows the speedups of different applications as we sweep aggregate commit queue entries from 128 (8 tasks per tile) to unbounded; the default is 1024 entries. Commit queues are fundamental to performance: fewer than 512 entries degrade performance considerably. More than 1024 entries confer moderate performance boosts to some applications. We conclude that 1024 entries strikes a good balance between performance and implementation cost for the benchmarks we study.

Bloom filter configuration: Fig. 17(b) shows the relative performance of different Bloom filter configurations. The default 2048-bit 8-way Bloom filters achieve performance within 10% of perfect conflict detection. Smaller Bloom filters cause frequent false positives and aborts in *silos* and *des*, which have the tasks with the largest footprint. However, *bfs*, *sssp*, and *msf* tasks access little data, so they are insensitive to Bloom filter size.

Frequency of GVT updates: Swarm is barely sensitive to the frequency of GVT updates. As we vary the period between GVT updates from 50 cycles to 800 cycles (the default is 200 cycles), performance at 64 cores drops from 0.1% in *sssp* to 3.0% in *msf*.

6.5 Swarm Case Study: astar

Finally, we present a case study of *astar* running on a 16-core, 4-tile system to analyze Swarm’s time-varying behavior. Fig. 18 depicts several per-tile metrics,

sampled every 500 cycles, over a 100 Kcycle interval: the breakdown of core cycles (top row), commit and task queue lengths (middle row), and tasks commit and abort events (bottom row). Each column shows these metrics for a single tile.

Fig. 18 shows that task queues are highly utilized throughout the interval. As task queues approach their capacity, coalescer tasks kick in, spilling tasks to memory. Commit queues, however, show varied occupancy. As tasks are executed out of order, they use a commit queue entry until they are safe to commit (or are aborted). Most of the time, commit queues are large enough to decouple execution and commit orders, and tiles spend the vast majority of time executing worker tasks.

Occasionally, however, commit queues fill up and cause the cores to stall. For example, tiles stall around the 40 Kcycle mark as they wait for a few straggler tasks to finish. The last of those stragglers finishes at 43 Kcycles, and the subsequent GVT update commits a large number of erstwhile speculative tasks, freeing up substantial commit queue space. These events explain *astar*’s sensitivity to commit queue size as seen in Fig. 17(a).

Finally, note that although queues fill up rarely, commits tend to happen in bursts throughout the run. This shows that fast commits are important, as they enable Swarm to quickly turn around commit queue entries.

7. ADDITIONAL RELATED WORK

Prior work has studied the limits of instruction-level parallelism under several idealizations, including a large or infinite instruction window, perfect branch prediction and memory disambiguation, and simple program transformations to remove unnecessary data dependences [4, 9, 18, 20, 24, 42, 49, 57, 74]. Similar to our limit study, these analyses find that parallelism is often plentiful ($>1000\times$), but very large instruction windows are needed to exploit it ($>100K$ instructions [42, 57]). Our oracle tool focuses on task-level parallelism, so it misses intra-task parallelism, which is necessarily limited with short tasks. Instead, we focus on removing superfluous dependences in scheduling data structures, uncovering large amounts of parallelism for irregular applications.

Several TLS schemes expose timestamps to software for different purposes, such as letting the compiler schedule loop iterations in Stampede [68], speculating across

barriers in TCC [29], and supporting out-of-order spawn of speculative function calls in Renau et al. [61]. These schemes work well for their intended purposes, but cannot queue or buffer tasks with arbitrary timestamps—they can only spawn new work if there is a free hardware context. Software scheduling would be required to sidestep this limitation, which, as we have seen, would introduce false data dependences and limit parallelism.

Prior work in fine-grain parallelism has developed a range of techniques to reduce task management overheads. Active messages lower the cost of sending tasks among cores [52, 73]. Hardware task schedulers such as Carbon [41] lower overheads further for specific problem domains. GPUs [76] and Anton 2 [27] feature custom schedulers for non-speculative tasks. By contrast, Swarm implements speculative hardware task management for a different problem domain, ordered parallelism.

Prior work has developed shared-memory priority queues that scale with the number of cores [2, 75], but they do so by relaxing priority order. This restricts them to benchmarks that admit order violations, and loss of order means threads often execute useless work far from the critical path [33, 34]. Nikas et al. [51] use hardware transactional memory to partially parallelize priority queue operations, accelerating `sssp` by $1.8\times$ on 14 cores. Instead, we dispense with shared-memory priority queues: Swarm uses distributed priority queues, load-balanced through random enqueues, and uses speculation to maintain order.

Our execution model has similarities to parallel discrete-event simulation (PDES) [21]. PDES events run at a specific virtual time and can create other events, but cannot access arbitrary data, making them less general than Swarm tasks. Moreover, state-of-the-art PDES engines have overheads of tens of thousands of cycles per event [6], making them impractical for fine-grain tasks. Fujimoto proposed the Virtual Time Machine (VTM), tailored to the needs of PDES [23], which could reduce these overheads. However, VTM relied on an impractical memory system that could be indexed by address and time.

8. CONCLUSIONS

We have presented Swarm, a novel architecture that unlocks abundant but hard-to-exploit irregular ordered parallelism. Swarm relies on a novel execution model based on timestamped tasks that decouples task creation and execution order, and a microarchitecture that performs speculative, out-of-order task execution and implements a large speculation window efficiently. Programs leverage Swarm’s execution model to convey new work to hardware as soon as it is discovered rather than in the order it needs to run, exposing a large amount of parallelism. As a result, Swarm achieves order-of-magnitude speedups on ordered irregular programs, which are key in emerging domains such as graph analytics, data mining, and in-memory databases [34, 55, 71]. Swarm hardware could also support thread-level speculation and transactional execution with minimal changes.

Swarm also opens several research avenues. First,

Swarm’s techniques may benefit a broader class of applications. For instance, Swarm could be applied to automatically parallelize general-purpose programs more effectively than prior TLS systems. Second, we have shown that co-designing the execution model and microarchitecture is a promising approach to uncover parallelism. Investigating new or more general execution models may expose additional parallelism in other domains. Third, Swarm shows that globally sequenced execution can be scalable even with fine-grained tasks. With additional work, Swarm’s techniques could be scaled to multi-chip and multi-machine systems. These topics are the subject of our ongoing and future work.

9. ACKNOWLEDGMENTS

We sincerely thank Nathan Beckmann, Harshad Kasture, Anurag Mukkara, Li-Shiuan Peh, Po-An Tsai, Guowei Zhang, and the anonymous reviewers for their helpful feedback. M. Amber Hassaan and Donald Nguyen graciously assisted with Galois benchmarks. This work was supported in part by C-FAR, one of six SRC STAR-net centers by MARCO and DARPA, and by NSF grant CAREER-1452994. Mark Jeffrey was partially supported by a MIT EECS Jacobs Presidential Fellowship and an NSERC Postgraduate Scholarship.

10. REFERENCES

- [1] “9th DIMACS Implementation Challenge: Shortest Paths,” 2006.
- [2] D. Alistarh, J. Kopinsky, J. Li *et al.*, “The SprayList: A scalable relaxed priority queue,” in *PPoPP*, 2015.
- [3] *4096x128 ternary CAM datasheet (28nm)*, Analog Bits, 2011.
- [4] T. Austin and G. Sohi, “Dynamic dependency analysis of ordinary programs,” in *ISCA-19*, 1992.
- [5] D. Bader, H. Meyerhenke, P. Sanders *et al.*, Eds., *10th DIMACS Implementation Challenge Workshop*, 2012.
- [6] P. Barnes Jr, C. Carothers, D. Jefferson *et al.*, “Warp speed: executing time warp on 1,966,080 cores,” in *PADS*, 2013.
- [7] G. Blleloch, J. Fineman, P. Gibbons *et al.*, “Internally deterministic parallel algorithms can be fast,” in *PPoPP*, 2012.
- [8] S. Brand and R. Bidarra, “Multi-core scalable and efficient pathfinding with Parallel Ripple Search,” *Computer Animation and Virtual Worlds*, 23(2), 2012.
- [9] M. Butler, T.-Y. Yeh, Y. Patt *et al.*, “Single instruction stream parallelism is greater than two,” in *ISCA-18*, 1991.
- [10] J. Carter and M. Wegman, “Universal classes of hash functions (extended abstract),” in *STOC-9*, 1977.
- [11] L. Ceze, J. Tuck, J. Torrellas *et al.*, “Bulk disambiguation of speculative threads in multiprocessors,” in *ISCA-33*, 2006.
- [12] S. Chandrasekaran and M. Hill, “Optimistic simulation of parallel architectures using program executables,” in *PADS*, 1996.
- [13] M. Cintra and D. Llanos, “Toward efficient and robust software speculative parallelization on multiprocessors,” in *PPoPP*, 2003.
- [14] J. Corbett, J. Dean, M. Epstein *et al.*, “Spanner: Google’s globally distributed database,” *ACM TOCS*, 31(3), 2013.
- [15] T. Cormen, C. Leiserson, R. Rivest *et al.*, *Introduction to Algorithms*, 3rd ed. MIT Press, 2009.
- [16] T. Davis and Y. Hu, “The University of Florida sparse matrix collection,” *ACM TOMS*, 38(1), 2011.
- [17] J. Devietti, B. Lucia, L. Ceze *et al.*, “DMP: deterministic shared memory multiprocessing,” in *ASPLOS-XIV*, 2009.
- [18] K. Ebcioglu, E. Altman, M. Gschwind *et al.*, “Optimizations and oracle parallelism with dynamic translation,” in *MICRO-32*, 1999.

- [19] H. Esmaeilzadeh, E. Blem, R. St Amant *et al.*, “Dark silicon and the end of multicore scaling,” in *ISCA-38*, 2011.
- [20] E. Fatehi and P. Gratz, “ILP and TLP in shared memory applications: a limit study,” in *PACT-23*, 2014.
- [21] A. Ferscha and S. Tripathi, “Parallel and distributed simulation of discrete event systems,” U. Maryland, Tech. Rep., 1998.
- [22] M. Fredman and R. Tarjan, “Fibonacci heaps and their uses in improved network optimization algorithms,” in *FOCS*, 1984.
- [23] R. Fujimoto, “The virtual time machine,” in *SPAA*, 1989.
- [24] S. Garold, “Detection and parallel execution of independent instructions,” *IEEE Trans. Comput.*, 19(10), 1970.
- [25] M. J. Garzarán, M. Prvulovic, J. M. Llbería *et al.*, “Trade-offs in buffering speculative memory state for thread-level speculation in multiprocessors,” in *HPCA-9*, 2003.
- [26] S. Ghemawat and P. Menage, “TCMalloc: Thread-caching malloc <http://goog-perftools.sourceforge.net/doc/tcmalloc.html>.”
- [27] J. Grossman, J. Kuskin, J. Bank *et al.*, “Hardware support for fine-grained event-driven computation in Anton 2,” in *ASPLOS-XVIII*, 2013.
- [28] L. Hammond, M. Willey, and K. Olukotun, “Data speculation support for a chip multiprocessor,” in *ASPLOS-VIII*, 1998.
- [29] L. Hammond, V. Wong, M. Chen *et al.*, “Transactional memory coherence and consistency,” in *ISCA-31*, 2004.
- [30] T. Harris, J. Larus, and R. Rajwar, “Transactional memory,” *Synthesis Lectures on Computer Architecture*, 2010.
- [31] P. Hart, N. Nilsson, and B. Raphael, “A formal basis for the heuristic determination of minimum cost paths,” *IEEE Trans. on Systems Science and Cybernetics*, 4(2), 1968.
- [32] M. A. Hassaan, D. Nguyen, and K. Pingali, “Brief announcement: Parallelization of asynchronous variational integrators for shared memory architectures,” in *SPAA*, 2014.
- [33] M. A. Hassaan, M. Burtscher, and K. Pingali, “Ordered vs. unordered: a comparison of parallelism and work-efficiency in irregular algorithms,” in *PPoPP*, 2011.
- [34] M. A. Hassaan, D. Nguyen, and K. Pingali, “Kinetic Dependence Graphs,” in *ASPLOS-XX*, 2015.
- [35] M. Hill and M. Marty, “Amdahl’s Law in the Multicore Era,” *Computer*, 41(7), 2008.
- [36] D. Hower, P. Montesinos, L. Ceze *et al.*, “Two hardware-based approaches for deterministic multiprocessor replay,” *Comm. ACM*, 2009.
- [37] T. Issariyakul and E. Hossain, *Introduction to network simulator NS2*. Springer, 2011.
- [38] D. Jefferson, “Virtual time,” *ACM TOPLAS*, 7(3), 1985.
- [39] J. Jun, S. Jacobson, J. Swisher *et al.*, “Application of discrete-event simulation in health care clinics: A survey,” *Journal of the operational research society*, 50(2), 1999.
- [40] C. Kim, D. Burger, and S. Keckler, “An adaptive, non-uniform cache structure for wire-delay dominated on-chip caches,” in *ASPLOS-X*, 2002.
- [41] S. Kumar, C. Hughes, and A. Nguyen, “Carbon: architectural support for fine-grained parallelism on chip multiprocessors,” in *ISCA-34*, 2007.
- [42] M. Lam and R. Wilson, “Limits of control flow on parallelism,” in *ISCA-19*, 1992.
- [43] C. Leiserson and T. Schardl, “A work-efficient parallel breadth-first search algorithm,” in *SPAA*, 2010.
- [44] A. Lew, J. Marsden, M. Ortiz *et al.*, “Asynchronous variational integrators,” *Arch. Rational Mech. Anal.*, 167(2), 2003.
- [45] T. Liu, C. Curtsinger, and E. D. Berger, “Dthreads: efficient deterministic multithreading,” in *SOSP-23*, 2011.
- [46] C. Luk, R. Cohn, R. Muth *et al.*, “Pin: building customized program analysis tools with dynamic instrumentation,” in *PLDI*, 2005.
- [47] J. Misra, “Distributed discrete-event simulation,” *ACM Computing Surveys (CSUR)*, 18(1), 1986.
- [48] K. Moore, J. Bobba, M. Moravan *et al.*, “LogTM: Log-based transactional memory,” in *HPCA-12*, 2006.
- [49] A. Nicolau and J. Fisher, “Using an oracle to measure potential parallelism in single instruction stream programs,” in *MICRO-14*, 1981.
- [50] K. Nii, T. Amano, N. Watanabe *et al.*, “A 28nm 400MHz 4-Parallel 1.6Gsearch/s 80Mb Ternary CAM,” in *ISSCC*, 2014.
- [51] K. Nikas, N. Anastopoulos, G. Goumas *et al.*, “Employing transactional memory and helper threads to speedup Dijkstra’s algorithm,” in *ICPP*, 2009.
- [52] M. Noakes, D. Wallach, and W. Dally, “The J-Machine multicompiler: an architectural evaluation,” in *ISCA-20*, 1993.
- [53] OpenStreetMap, “<http://www.openstreetmap.org>.”
- [54] R. Panigrahy and S. Sharma, “Sorting and searching using ternary CAMs,” *IEEE Micro*, 23(1), 2003.
- [55] K. Pingali, D. Nguyen, M. Kulkarni *et al.*, “The tao of parallelism in algorithms,” in *PLDI*, 2011.
- [56] L. Porter, B. Choi, and D. Tullsen, “Mapping out a path from hardware transactional memory to speculative multithreading,” in *PACT-18*, 2009.
- [57] M. Postiff, D. Greene, G. Tyson *et al.*, “The limits of instruction level parallelism in SPEC95 applications,” *Comp. Arch. News*, 27(1), 1999.
- [58] X. Qian, B. Sahelices, and J. Torrellas, “OmniOrder: Directory-based conflict serialization of transactions,” in *ISCA-41*, 2014.
- [59] S. Reinhardt, M. Hill, J. Larus *et al.*, “The Wisconsin Wind Tunnel: virtual prototyping of parallel computers,” in *SIGMETRICS*, 1993.
- [60] J. Renau, K. Strauss, L. Ceze *et al.*, “Thread-level speculation on a CMP can be energy efficient,” in *ICS’05*, 2005.
- [61] J. Renau, J. Tuck, W. Liu *et al.*, “Tasking with out-of-order spawn in TLS chip multiprocessors: microarchitecture and compilation,” in *ICS’05*, 2005.
- [62] D. Sanchez and C. Kozyrakis, “ZSim: Fast and Accurate Microarchitectural Simulation of Thousand-Core Systems,” in *ISCA-40*, 2013.
- [63] D. Sanchez, L. Yen, M. Hill *et al.*, “Implementing signatures for transactional memory,” in *MICRO-40*, 2007.
- [64] D. Sanchez, R. Yoo, and C. Kozyrakis, “Flexible architectural support for fine-grain scheduling,” in *ASPLOS-XV*, 2010.
- [65] S. Schneider, C. Antonopoulos, and D. Nikolopoulos, “Scalable locality-conscious multithreaded memory allocation,” in *ISMM-5*, 2006.
- [66] G. Sohi, S. Breach, and T. Vijaykumar, “Multiscalar processors,” in *ISCA-22*, 1995.
- [67] L. Soule and A. Gupta, “An evaluation of the Chandy-Misra-Bryant algorithm for digital logic simulation,” *ACM TOMACS*, 1(4), 1991.
- [68] J. G. Steffan, C. Colohan, A. Zhai *et al.*, “A scalable approach to thread-level speculation,” in *ISCA-27*, 2000.
- [69] J. G. Steffan and T. Mowry, “The potential for using thread-level data speculation to facilitate automatic parallelization,” in *HPCA-4*, 1998.
- [70] S. Thoziyoor, N. Muralimanohar, J. H. Ahn *et al.*, “CACTI 5.1,” HP Labs, Tech. Rep. HPL-2008-20, 2008.
- [71] S. Tu, W. Zheng, E. Kohler *et al.*, “Speedy transactions in multicore in-memory databases,” in *SOSP-24*, 2013.
- [72] A. Varga and A. Şekercioglu, “Parallel simulation made easy with OMNeT++,” 2003.
- [73] T. Von Eicken, D. Culler, S. Goldstein *et al.*, “Active messages: a mechanism for integrated communication and computation,” in *ISCA-19*, 1992.
- [74] D. Wall, “Limits of instruction-level parallelism,” in *ASPLOS-IV*, 1991.
- [75] M. Wimmer, F. Versaci, J. Träff *et al.*, “Data structures for task-based priority scheduling,” in *PPoPP*, 2014.
- [76] C. Wittenbrink, E. Kilgariff, and A. Prabhu, “Fermi GF100 GPU architecture,” *IEEE Micro*, 31(2), 2011.
- [77] M. Xu, R. Bodik, and M. Hill, “A flight data recorder for enabling full-system multiprocessor deterministic replay,” in *ISCA-30*, 2003.
- [78] C. Yang and B. Miller, “Critical path analysis for the execution of parallel and distributed programs,” in *ICDCS*, 1988.
- [79] L. Yen, J. Bobba, M. Marty *et al.*, “LogTM-SE: Decoupling hardware transactional memory from caches,” in *HPCA-13*, 2007.
- [80] Y. Zhang, L. Rauchwerger, and J. Torrellas, “Hardware for speculative parallelization of partially-parallel loops in DSM multiprocessors,” in *HPCA-5*, 1999.