

Grappa: A Latency-Tolerant Runtime for Large-Scale Irregular Applications

Jacob Nelson, Brandon Holt, Brandon Myers, Preston Briggs, Luis Ceze, Simon Kahan, Mark Oskin

University of Washington, Department of Computer Science and Engineering
{nelson, bholt, bdmyers, preston, luisceze, skahan, oskin}@cs.washington.edu

Abstract

Grappa is a runtime system for commodity clusters of multicore computers that presents a massively parallel, single address space abstraction to applications. Grappa’s purpose is to enable scalable performance of irregular parallel applications, such as branch and bound optimization, SPICE circuit simulation, and graph processing. Poor data locality, imbalanced parallel work and complex communication patterns make scaling these applications difficult.

Grappa serves both as a C++ user library and as a foundation for higher level languages. Grappa tolerates delays to remote memory by multiplexing thousands of lightweight workers to each processor core, balances load via fine-grained distributed work-stealing, increases communication throughput by aggregating smaller data requests into large ones, and provides efficient synchronization and remote operations. We present a description of the Grappa system and an evaluation of its performance on microbenchmarks.

1. Introduction

Irregular applications exhibit workloads, dependences, and memory accesses that are highly sensitive to input. Classic examples of such applications include branch and bound optimization, SPICE circuit simulation, and car crash analysis. Important contemporary examples include processing large graphs in the business, national security, and social network computing domains. For these emerging applications, reasonable response time – given the sheer amount of data – requires large multinode systems. The most broadly available multinode systems are those built from x86 compute nodes interconnected via Ethernet or InfiniBand. However, scalable performance of irregular applications on these systems is elusive for two reasons:

Poor data locality and frequent communication Data reference patterns of irregular applications are unpredictable and tend to be spread across the entire system. This results in frequent requests for small pieces of remote data. Caches are of little assistance because of low temporal and spatial locality. Prefetching is of limited value because request locations are not known early enough. Data-parallel frameworks such as MapReduce [18] are ineffective because they rely on data partitioning and regular communication patterns. Consequently, commodity networks, which are designed for large packets, achieve just a fraction of their peak bandwidth on small messages, starving application performance.

High network communication latency The performance challenges of frequent communication are exacerbated by high network latency relative to processor performance. Latency of commodity networks runs anywhere from a few to hundreds of microseconds – tens of thousands of processor clock cycles. Since irregular application tasks encounter remote references dynamically during execution and must resolve them before making further progress, stalls are frequent and lead to severely underutilized compute resources.

While some irregular applications can be manually restructured to better exploit locality, aggregate requests to increase network message size, and manage the additional challenges of load balance

and synchronization, the effort required to do so is formidable and involves knowledge and skills pertaining to distributed systems far beyond those of most application programmers. Luckily, many of the important irregular applications naturally offer large amounts of concurrency. This immediately suggests taking advantage of concurrency to tolerate the latency of data movement by overlapping computation with communication.

The fully custom Tera MTA-2 [4, 5] system is a classic example of supporting irregular applications by using concurrency to hide latencies. It had a large distributed shared memory with no caches. On every clock cycle, each processor would execute a ready instruction chosen from one of its 128 hardware thread contexts, a sufficient number to fully tolerate memory access latency. The network was designed with a single-word injection rate that matched the processor clock frequency and sufficient bandwidth to sustain a reference from every processor on every clock cycle. Unfortunately, the MTA-2’s relatively low single-threaded performance meant that it was not general enough nor cost-effective. The Cray XMT approximates the Tera MTA-2, reducing its cost but not overcoming its narrow range of applicability.

We believe we can support irregular applications with good performance and cost-effectiveness with commodity hardware for two main reasons. First, commodity multicore processors have become extremely fast with high clock rates, large caches and robust DRAM bandwidth. Second, commodity networks offer high bandwidth as long as messages are large enough. We build on these two observations and develop Grappa, a software runtime system that allows a commodity cluster of x86-based nodes connected via an InfiniBand network to be programmed as if it were a single, large, shared-memory NUMA (non-uniform memory access) machine with scalable performance for irregular applications. Grappa exploits fast processors and the memory hierarchy to provide a lightweight user-level tasking layer that supports a context switch in as little as 38ns and can sustain a large number of active workers. It bridges the commodity network bandwidth gap with a communication layer that combines short messages originating from many concurrent workers into larger packets.

As a general design philosophy, Grappa trades latency for throughput. By *increasing* latency in key components of the system we are able to: increase effective random access memory bandwidth by delaying and aggregating messages; increase synchronization rate by delegating atomic operations to gatekeeper cores, even when referencing node-local global data; and improve load balance via work-stealing. Grappa then exploits parallelism to tolerate the increased latency.

Our evaluation of Grappa shows that the core components, scheduling and communication, achieve their design goals. Thousands of workers can be efficiently context switched on a multicore processor. Grappa supports high random-access bandwidth and is able to exploit and load-balance irregular parallelism.

Grappa comprises three main components: a tasking layer, distributed shared memory, and a communication layer. We will explore each of these in turn.

2. Tasking System

The basic unit of execution in Grappa is a *task*. When tasks are ready to execute, they are mapped to a *worker*, which is akin to a user-level thread. Each hardware core has a single operating system thread pinned to it.

Tasks Tasks are specified by a closure (or “function object” in C++ parlance) that holds both code to execute and initial state. The functor can be specified with a function pointer and explicit arguments, a C++ struct that overloads the parentheses operator, or a C++11 lambda construct. These objects, typically very small (on the order of 64 bytes), hold read-only values such as an iteration index and pointers to common data or synchronization objects. Task functors can be serialized and transported around the system, and eventually executed by a worker, as described next.

Workers Workers execute application and system (e.g., communication) tasks. A worker is simply a collection of status bits and a stack, allocated at a particular core. When a task is ready to execute it is assigned to a worker, which executes the task functor on its own stack. Once a task is mapped to a worker it stays with that worker until it finishes.

Scheduling During execution, a worker yields control of its core whenever performing a long-latency operation, allowing the processor to remain busy while waiting for the operation to complete. In addition, a programmer can direct scheduling explicitly. To minimize context-switch overhead, the Grappa scheduler operates entirely in user-space and does little more than store state of one worker and load that of another. When a task encounters a long-latency operation, its worker is suspended and subsequently woken when the operation completes.

Each core in a Grappa system has its own independent scheduler. The scheduler has a collection of active workers ready to execute called the *ready worker queue*. Each scheduler also has three queues of tasks waiting to be assigned a worker:

deadline task queue a priority queue of tasks that are executed according to task-specific deadline constraints;

private task queue a queue of tasks that must run on this core and is therefore not subject to stealing;

public task queue a queue of tasks that are waiting to be matched with workers. It is a local partition of a shared task pool.

Whenever a task yields, the scheduler makes a decision about what to do next. First, any task in the deadline task queue whose deadline is imminent is chosen for execution. This queue manages high priority system tasks, such as periodically servicing communication requests. Second, the scheduler determines if any workers with running tasks are ready to execute; if so, one is scheduled. Finally, if no workers are ready to run, but tasks are waiting to be matched with workers, an idle worker is woken (or a new worker is spawned), matched with a task, and scheduled.

Context switching Grappa context switches between workers non-preemptively. As with other cooperative multithreading systems, we treat context switches as function calls, saving and restoring only the callee-saved state as specified in the x86-64 ABI [6]. This involves saving six general-purpose 64-bit registers and the stack pointer, as well as the 16-bit x87 floating point control word and the SSE context/status register. Thus, the minimum amount of state a cooperative context switch routine must save, according to the ABI, is 62 bytes.

Since Grappa keeps a very large number of active workers, their context data will not fit in cache. By oversubscribing on the number of workers beyond what is required for local DRAM latency tolerance, the scheduler can ensure there is always some number of context

pointers in the ready queue. This allows the scheduler to prefetch contexts into cache using software prefetch instructions; the size of the L1 cache is sufficient to hold enough contexts to tolerate the latency to main memory. Empirically we find that prefetching the fourth worker in the scheduling order is sufficient. This prefetching constrains the types of task scheduling decisions that can be made but makes context switching effectively free of cache misses, even to hundreds of thousands of workers. We provide an analysis of our context switch performance in Section 5.1.

Work stealing When the scheduler finds no work to assign to its workers, it commences to steal tasks from other cores using an asynchronous `call_on` active message. It chooses a victim at random until it finds one with a non-zero amount of work in its public task queue. The scheduler steals half of the tasks it finds at the victim. Work stealing is particularly interesting in Grappa since performance depends on having many active worker threads on each core. Even if there are many active threads, if they are all suspended on long-latency operations, then the core is underutilized. The stealing policy must predict whether local tasks will likely generate enough new work soon; a similar problem is addressed in [41].

2.1 Expressing Parallelism

Grappa programmers focus on expressing as much parallelism as possible without concern for where it will execute. Grappa then chooses where and when to exploit this parallelism, scheduling as much work as is necessary on each core to keep it busy in the presence of system latencies and task dependences.

Grappa provides three methods for expressing parallelism. First, a single task can be created to execute in parallel with the current task by calling `spawn` with a functor. This adds it to the queue of tasks which will be executed the next time a worker is available. Second, the programmer can invoke a parallel for loop with `parallel_for`, provided that the trip count is known at loop entry. The programmer specifies a functor which takes the loop index as a parameter, and an optional threshold to control parallel overhead. Grappa does *recursive decomposition* of iterations, similar to Cilk’s `cilk_for` construct [9], and TBB’s `parallel_for` [36]. It generates a logarithmically-deep tree of tasks, stopping to execute the loop body when the number of iterations is below the required threshold. Third, parallelism can be expressed via asynchronous delegate operations, which are explained next, in Section 3.

Figure 1 shows sample code using Grappa for a parallel tree search. Note how the code looks very similar to a recursive search procedure for a shared-memory system, without regard for communication, and Grappa’s parallel loop construct allows easy parallelization of the search.

3. Distributed Shared Memory

Applications written for Grappa utilize two forms of memory: local and global. Local memory is local to a single core within a node in the system. Accesses occur through conventional pointers. Applications use local accesses for a number of things in Grappa: the stack associated with a task, accesses to localized global memory in caches (see below), and accesses to debugging infrastructure local to each system node. Local pointers cannot access memory on other cores, and are valid only on their home core.

Large data that is expected to be shared and accessed with low locality is stored in Grappa’s global memory. All global data must be accessed through calls into Grappa’s API.

Global memory addressing Grappa provides two methods for storing data in the global memory. The first is a distributed heap striped across all the machines in the system in a block-cyclic fashion. Addresses to memory in the global heap use *linear addresses*. Choosing

```

class Vertex {
    Key key;
    int64_t numChildren;
    GlobalAddress<Vertex> children;
};

void search(GlobalAddress<Vertex> vtx_addr,
           Key key, GlobalAddress<Vertex> result) {
    // blocking remote read to get vertex info
    Vertex vtx = delegate_read(vtx_addr);
    if (vtx.key == key) {
        // key found
        delegate_write(result, vtx);
    } else {
        // spawn stealable tasks for iterations
        parallel_for(0, vtx.numChildren, [=](int i) {
            // recursive search
            search(vtx.children+i, key, result);
        });
    }
}

```

Figure 1: Sample Grappa code illustrating a parallel tree search similar to the unbalanced tree search benchmark we describe later. Children are spread over the system, so each parallel recursive search performs a delegate read to get vertex data.

the block size involves trading off sequential bandwidth against aggregate random access bandwidth. The block size, which is configurable, is typically set to 64 bytes, or the size of a single hardware cache line, in order to exploit spatial locality when available.

Grappa also allows any local data on a core’s stacks or heap to be exported to the global address space to be made accessible to other cores across the system. Addresses to global memory allocated in this way use *2D global addresses*. This uses a traditional PGAS (partitioned global address space [19]) addressing model, where each address is a tuple of a rank in the job (or global process ID) and an address in that process. Any node-local data can be made accessible to other cores in the system by wrapping the address and node ID into a 2D global address. This address can then be accessed with a delegate operation and even be buffered by other cores. 2D addresses may refer to memory allocated from a single processes’ heap or from a task’s stack.

Global memory access Access to Grappa’s distributed shared memory is provided through *delegate* operations, which are short memory accesses performed at the memory location’s home node. When the data access pattern has low-locality, it is more efficient to modify the data on its home core rather than bringing a copy to the requesting core and returning it after modification. Delegate operations [30, 33] provide this capability. Applications can dispatch computation to be performed on individual machine-word sized chunks of global memory to the memory system itself. Delegates can execute arbitrary code, provided they do not block, to ensure communicator workers make progress. Provided they touch only memory owned by a single core, we can use them to perform simple *read/write* operations to global memory, as well as more complex *read-modify-write* operations (e.g., *fetch-and-add*). We use these primitive operations to implement higher-level synchronization mechanisms such as mutexes, condition variables, and full-empty bits.

Delegate operations are *always* executed at the home core of their address. The remote operation may not perform any operations that could cause a context switch; this ensures any modifications are atomic. We limit delegate operations to operate on objects in the 2D address space or objects that fit in a single block of the linear address space so they can be satisfied with a single network request. Given these restrictions, we can ensure that delegate operations for the

same address from multiple requesters are always serialized through a single core in the system, providing atomic semantics without using actual atomic operations (and thus avoiding their typical high cost).

Delegate operations can be either *blocking* or *asynchronous*. With blocking operations, the task issuing the delegate call blocks until the delegate operation completes, which is necessary, for example, to ensure that synchronization has finished before continuing. On the other hand, remote data accesses often can overlap, and delegates with no return value may not need to block the caller. To avoid unnecessary waiting, we support asynchronous delegate operations. For reads, we support a “futures”-like mechanism which allows tasks to issue reads in parallel and block on the “promises” returned. Delegate write operations may also be performed asynchronously, but synchronization is still needed to ensure that asynchronous operations have completed. Grappa provides a `GlobalCompletionEvent` synchronization object, which asynchronous operations (including tasks) can be enrolled. Tasks can block on these objects to be woken when all enrolled operations are complete.

When programmers want to operate on data structures spread across multiple nodes, accesses must be expressed as multiple delegate operations along with with appropriate synchronization operations. Grappa’s API also includes calls for gathering and scattering contiguous blocks in the global heap, but the user is responsible for ensuring correct synchronization.

Memory consistency model discussion As mentioned earlier, all synchronization operations are done via delegate operations. Since they all execute on their home core in some serial order, they are guaranteed to be globally linearizable [24], with their updates visible to all cores across the system in the same order. In addition, only one synchronous delegate will be in flight at a time from a particular task. Therefore, synchronization operations from a particular task are not subject to reordering. Consequently, all synchronization operations execute in program order and are made visible in the same order to all cores in the system. These properties are sufficient to guarantee a memory model that offers sequential consistency for data-race-free programs [1] (all accesses to shared data are separated by synchronization). This is the memory model that underpins C/C++ [10, 26].

Note, however, that if the application code uses explicit buffers or asynchronous delegates to access shared data, all updates must be published back to the home core before the synchronization operation that protects the data is performed. This is done using release operations on cached regions and using the `GlobalCompletionEvent` object to determine that asynchronous delegates have completed.

4. Communication Support

Grappa’s communication support has two layers: user-level messaging interface based on active messages; and network-level transport that supports request aggregation for better communication bandwidth.

Active message interface At the upper (user-level) layer, Grappa implements asynchronous active messages [42]. Each message consists of a function pointer, an optional argument payload, and an optional data payload.

Message aggregation In our experiments the vast majority of application-level messages are small, between 32 and 64 bytes. Our measurements confirm manufacturers’ published data [15]; with packets of this size, the available bisection bandwidth is only a small fraction (3%) of the peak bisection bandwidth. As mentioned earlier, commodity networks including InfiniBand achieves their peak bisection bandwidth only when the packet sizes are relatively large – on the order of multiple kilobytes. The reason for this discrepancy is the combination of overheads associated with handling each packet

(in terms of bytes that form the actual packet, processing time at the card, multiple round-trips on the PCI Express bus and processing on the CPU within the driver stack). Consequently, to make the best use of the network, we must convert small messages into large ones.

Message processing mechanics Since communication is very frequent in Grappa, aggregating and sending messages efficiently is very important. To achieve that, Grappa makes careful use of caches, prefetching, and lock-free synchronization operations.

Each processing core of a system node maintains an array of outgoing message lists. The array size is the number of system cores in the Grappa system. The outgoing message lists and messages are located in a region of memory shared across all cores in a Grappa node (thus enabling cores to peek at each other’s message lists). When a task sends a message, it allocates a buffer (typically on its stack), determines the destination system node, and links the buffer into the corresponding linked list.

Each processing core in a given system node is responsible for aggregating and sending the resulting messages from all cores on that node to a set of destination nodes. Each core periodically executes a task responsible for sending messages. This task examines the private (to each core) message lists for each destination node it is responsible for managing and, if the list is long enough or a message has waited past a time-out period, all messages to a given destination system node from that source system node are sent. Aggregating and sending a message involves manipulating a set of shared data-structures (the message lists). This is done using CAS (compare-and-swap) operations to avoid high synchronization costs. Note that we use a per-core array of message lists that is only periodically modified across processor cores after experimentally determining that this approach was faster (sometimes significantly) than a global per-system node array of message lists.

Each node has a region of memory with send buffers where the final aggregated messages are built. These buffers are visible to the network card, and messages are sent with user-mode operations only. When the worker responsible for outbound messages to a given system node has received a sufficient number of message send requests or a timeout is reached, the linked list of messages is walked and messages are copied to a send buffer. This process requires careful prefetching because most of the outbound messages are *not* in the processor cache at this time (recall that a core can be aggregating messages originating from other cores in the same node). Once the send buffer has been formed, it is handed off to GASNet for transfer to the remote system node. RDMA is used if the underlying network supports it.

Once the remote system node has received the message buffer, a management task is spawned to manage the unpacking process. The management task spawns a task on each core at the receiving system to simultaneously unpack messages destined for that core. Upon completion, these unpacking tasks synchronize with the management task. Once all cores have processed the message buffer, the management task sends a reply to the sending system node indicating the successful delivery of the messages.

5. Evaluation

We implemented the Grappa in C++ for the Linux operating system. The core runtime system system is about 15K lines of code. We ran these experiments on a cluster of AMD Interlagos processors. Nodes have 32 2.1-GHz cores in two sockets, 64GB of memory, and 40Gb Mellanox ConnectX-2 InfiniBand network cards. Nodes are connected via a QLogic InfiniBand switch.

5.1 Basic Grappa Performance

User-level context switching Fast context switching is at the heart of Grappa’s latency tolerance abilities. We assess context switch

overheads using a simple microbenchmark that runs a configurable number of workers on a single core, where each worker increments values in a large array.

Figure 2a shows the average context switch time as the number of workers grow. At our standard operating point ($\approx 1K$ workers), context switch time is on the order of 50ns. As we add workers, the time increases slowly, but levels off: we also ran with 500,000 workers (10 times what is shown in the figure) and found that context switch time was around 75ns. In comparison, for the same yield test using kernel-level Pthreads on a single core, the switch time is 450ns for a few threads and 800ns for 1000–32000 threads.

If we calculate aggregate context switch *rate* of all cores in a node, we find that with prefetching, Grappa context switching is limited *not* by memory latency, as normally assumed, but rather memory bandwidth. Specifically, we empirically found that 4 cache lines (1 for worker struct and 3 for stack data) was sufficient to avoid cache misses in the microbenchmark. Every context switch then requires 8 cache-line transfers. The off-chip bandwidth of a single socket in our system is 270M cache lines per second [32, 33]. This implies that, in the limit, we can sustain at most 34M context switches per second per socket (a context-switch time of 29ns).

In summary, our tasking layer is able to efficiently sustain very high concurrency and, as we will show later, the amount of concurrency sustained is sufficient for the latencies Grappa needs to hide.

Global memory and communication We measure the performance of Grappa’s global memory and communication layers using a faithful implementation of the giga updates per second (GUPS) benchmark, which measures cluster-wide random access bandwidth. Read-modify-write updates are dispatched at random to a global large array. This benchmark stresses the communication layer of Grappa separately from the scheduler, because only a single worker is used per system node. Figure 2b shows that Grappa is able to sustain well over a billion updates per second with 64 nodes. Note also that when aggregation is turned off, the update rate is nearly flat. Clearly aggregation is instrumental for good communication performance.

This compares very favorably to published results [25] for other high-end HPC systems. Though the actual computation done by GUPS is not useful, irregular, data-intensive applications typically must be able to sustain a high rate of random accesses in order to, for example, visit and mark vertices during a graph traversal. High random access rate in a distributed setting has been a long-standing challenge in HPC.

Putting it all together with Unbalanced Tree Search (UTS) Unbalanced Tree Search (UTS) is a benchmark for evaluating the programmability and performance of systems for parallel applications that require dynamic load balancing [35]. It involves traversing an unbalanced implicit tree: at each vertex, its number of children is sampled from some probability distribution, and this number of new nodes are added to a work queue to be visited. While this benchmark captures irregular, dynamic *computation*, we actually want to evaluate performance of algorithms with irregular *memory* access patterns. Thus we augment UTS by using the existing traversal code to create a large tree in memory, and then we traverse the in-memory tree. In our modified UTS, the timed portion is this traversal of the in-memory tree. This in-memory traversal has no knowledge of the tree structure beforehand.

Figure 2c shows the overall performance of Grappa running UTS. This experiment demonstrates that Grappa’s context switching and communication layers can be used together, while balancing workload, to run an irregular application efficiently.

Visiting vertices in the distributed tree requires mostly remote accesses, and because each vertex in the tree must be visited before

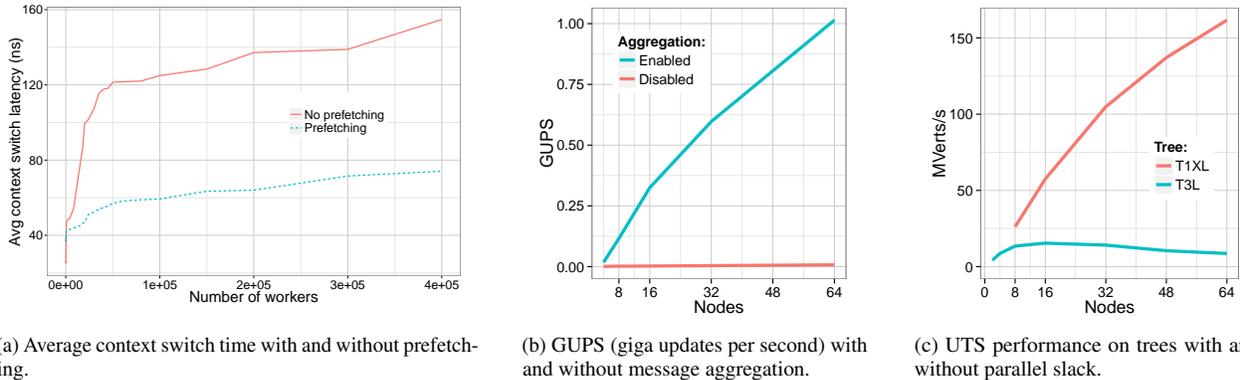


Figure 2: Grappa performance.

it can be expanded, blocking remote reads are required. In this case, we are forced to context switch to tolerate the remote access and continue aggregating. Figure 1 shows a closely analogous tree search in Grappa.

We look at two classes of trees, T1 and T3, from the original benchmark. T1 trees are very shallow and wide (i.e., significant parallel slack [40]), while T3 trees are very deep (i.e., little parallel slack). Given that access to each vertex is a random access, the critical path to search T3 trees is very long, hence the low performance and scalability. On such trees, we do not expect there to be sufficient concurrency for any system, including Grappa, to achieve high throughput – at the 16-node data point, the average active tasks per core over the search is an order of magnitude larger for T1 than for T3. Given the lack of parallelism, scaling up only serves to reduce throughput by distributing the tree to more machines. On the other hand, Grappa performs and scales very well for T1 trees.

6. Related Work

Multithreading Grappa uses multithreading to tolerate memory latency. This is a well known technique. Hardware implementations include the Denelcor HEP [38], Tera MTA [5], Cray XMT [21], Simultaneous multithreading [39], MIT Alewife [2], Cyclops [3], and even GPUs [20].

Grappa’s closest ancestor is the Threaded Abstract Machine [16], a runtime system for prototyping dataflow execution models on distributed memory supercomputers. The Active Messages [42] work that grew out of this project inspired our communication layer. One of the conclusions of this work [17] was that context switch costs can be low only when contexts are in cache, and that latency tolerance was not sufficient to guarantee performance on commodity processors. Grappa demonstrates that times have changed: modern commodity processors have sufficient bandwidth and prefetch capacity to stream contexts from DRAM and sustain a very large number of active contexts.

Grappa implements its own software-based multithreading with a lightweight user-mode task scheduler to multiplex *thousands* of tasks on a single processing core. The main difference between Grappa’s support for lightweight threads and prior work such as QThreads [43] and Capriccio [7] is context prefetching, which is needed for good performance when multiplexing such a large number of tasks.

Software distributed shared memory. Many approaches to SDSM have been explored, including IVY [28], Munin [8, 12], TreadMarks [27], and Blizzard [37]. Much of the innovation has been focused around reducing the synchronization cost of doing updates. Grappa’s delegate-based approach to updates avoids synchronization overhead entirely, providing sequential consistency for data-race-

free programs without the cost of a full coherence protocol. Grappa accepts the random access nature of irregular applications and optimizes for throughput rather than low latency.

Partitioned Global Address Space languages. The high-performance computing community has largely discarded the coherent distributed shared memory approach in favor of the Partitioned Global Address Space (PGAS) model. Example include Split-C [15], Chapel [13], X10 [14], Co-array Fortran [34] and UPC [19]. Grappa shares many parts of its design philosophy with these languages, but differs in that most PGAS languages expect programmers to modify algorithms to take advantage of locality by processing node-local data as much as possible; access to data stored on other nodes is possible but is seen as something best avoided. Grappa optimizes for random access to data anywhere in the cluster, while still allowing the exploitation of available locality.

Distributed graph processing systems. While Grappa is a general runtime system for any large-scale concurrent application, its features are well-suited for graph analysis. Other distributed graph processing frameworks include Pregel [31] and GraphLab/PowerGraph [22, 29]. Both adopt vertex-oriented programming models that work well for some application domains like machine learning but are too restrictive for general computation.

While the bulk-synchronous MapReduce [18] model and related systems such as Hadoop [23] are not a good fit for irregular or graph problems, the ideas have been extended by systems such as HaLoop [11] and Spark [44] to better fit iterative graph-based machine learning problems. Again, the programming model is restricted compared to Grappa, in exchange for good streaming IO support and the ability to tolerate node failures.

7. Conclusion

Irregular computations are both important and challenging to execute quickly. Scaling these applications easily on commodity hardware has been a historical challenge. Grappa simplifies this task for software developers and compiler writers. Grappa’s key aspect is extreme latency tolerance, which not only hides network latency but also enables the system to spend time on sophisticated work stealing and network optimizations, trading latency for even more throughput.

Our evaluation of Grappa shows that the core components, scheduling and communication, achieve their design goals. Thousands of workers can be efficiently context switched on a multicore processor, limited by DRAM bandwidth. Aggregating messages enables Grappa to achieve over 1.0 GUPS on 64 nodes. Grappa is able to exploit and load-balance irregular parallelism in UTS. Work is ongoing to compare Grappa’s performance with that of other systems.

References

- [1] Sarita V. Adve and Mark D. Hill. Weak ordering – A new definition. In *ISCA-17*, 1990.
- [2] Anant Agarwal, Ricardo Bianchini, David Chaiken, Kirk L. Johnson, David Kranz, John Kubiawicz, Beng-Hong Lim, Ken Mackenzie, and Donald Yeung. The MIT Alewife machine: Architecture and performance. In *22nd Annual International Symposium on Computer Architecture*, June 1995.
- [3] George Almási, Călin Caşcaval, José G. Castaños, Monty Denneau, Derek Lieber, José E. Moreira, and Henry S. Warren, Jr. Dissecting Cyclops: A detailed analysis of a multithreaded architecture. *SIGARCH Computer Architecture News*, 31:26–38, March 2003.
- [4] Gail Alverson, Robert Alverson, David Callahan, Brian Koblenz, Allan Porterfield, and Burton Smith. Exploiting heterogeneous parallelism on a multithreaded multiprocessor. In *Proceedings of the 6th international conference on Supercomputing*, ICS '92, pages 188–197, New York, NY, USA, 1992. ACM.
- [5] Robert Alverson, David Callahan, Daniel Cummings, Brian Koblenz, Allan Porterfield, and Burton Smith. The Tera computer system. In *Proceedings of the 4th International Conference on Supercomputing*, ICS '90, pages 1–6, New York, NY, USA, 1990. ACM.
- [6] AMD64 ABI. <http://www.x86-64.org/documentation/abi-0.99.pdf>, July 2012.
- [7] Rob Von Behren, Jeremy Condit, Feng Zhou, George C. Necula, and Eric Brewer. Capriccio: Scalable threads for internet services. In *In Proceedings of the 19th ACM Symposium on Operating Systems Principles*, pages 268–281. ACM Press, 2003.
- [8] John K. Bennett, John B. Carter, and Willy Zwaenepoel. Munin: Distributed shared memory based on type-specific memory coherence. In *Proceedings of the Second ACM SIGPLAN symposium on Principles and Practice of Parallel Programming*, PPOPP '90, pages 168–176, New York, NY, USA, 1990. ACM.
- [9] Robert D. Blumofe, Christopher F. Joerg, Bradley C. Kuszmaul, Charles E. Leiserson, Keith H. Randall, and Yuli Zhou. Cilk: An efficient multithreaded runtime system. In *Proceedings of the Fifth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPOPP '95, pages 207–216, New York, NY, USA, 1995. ACM.
- [10] Hans-J. Boehm. A Less Formal Explanation of the Proposed C++ Concurrency Memory Model. C++ standards committee paper WG21/N2480 = J16/07-350, <http://www.open-std.org/JTC1/SC22/WG21/docs/papers/2007/n2480.html>, December 2007.
- [11] Yingyi Bu, Bill Howe, Magdalena Balazinska, and Michael D. Ernst. Haloop: Efficient iterative data processing on large clusters. *Proc. VLDB Endow.*, 3(1-2):285–296, September 2010.
- [12] John B. Carter, John K. Bennett, and Willy Zwaenepoel. Implementation and performance of Munin. In *Proceedings of the Thirteenth ACM Symposium on Operating Systems Principles*, SOSPP '91, pages 152–164, New York, NY, USA, 1991. ACM.
- [13] Bradford L. Chamberlain, David Callahan, and Hans P. Zima. Parallel programmability and the Chapel Language. *International Journal of High Performance Computing Application*, 21(3):291–312, August 2007.
- [14] Philippe Charles, Christian Grothoff, Vijay Saraswat, Christopher Donawa, Allan Kielstra, Kemal Ebcioglu, Christoph von Praun, and Vivek Sarkar. X10: An object-oriented approach to non-uniform cluster computing. In *Proceedings of the 20th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications*, OOPSLA '05, pages 519–538, New York, NY, USA, 2005. ACM.
- [15] David E. Culler, Andrea Dusseau, Seth Copen Goldstein, Arvind Krishnamurthy, Steven Lumetta, Thorsten von Eicken, and Katherine Yelick. Parallel programming in Split-C. In *Proceedings of the 1993 ACM/IEEE conference on Supercomputing*, Supercomputing '93, pages 262–273, New York, NY, USA, 1993. ACM.
- [16] David E. Culler, Seth Copen Goldstein, Klaus Erik Schauer, and Thorsten von Eicken. TAM – A compiler controlled threaded abstract machine. *Journal of Parallel and Distributed Computing*, 18:347–370, July 1993.
- [17] David E. Culler, Klaus E. Schauer, and Thorsten von Eicken. Two fundamental limits on dataflow multiprocessing. In *Proceedings of the IFIP WG10.3. Working Conference on Architectures and Compilation Techniques for Fine and Medium Grain Parallelism*, PACT '93, pages 153–164, Amsterdam, The Netherlands, The Netherlands, 1993. North-Holland Publishing Co.
- [18] Jeffrey Dean and Sanjay Ghemawat. MapReduce: Simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, January 2008.
- [19] Tarek El-Ghazawi, William Carlson, Thomas Sterling, and Katherine Yelick. *UPC: Distributed Shared Memory Programming*. John Wiley and Sons, Inc., Hoboken, NJ, USA, 2005.
- [20] Kayvon Fatahalian and Mike Houston. A closer look at GPUs. *Communications of the ACM*, 51:50–57, October 2008.
- [21] John Feo, David Harper, Simon Kahan, and Petr Konecny. Eldorado. In *Proceedings of the 2nd Conference on Computing Frontiers*, CF '05, pages 28–34, New York, NY, USA, 2005. ACM.
- [22] Joseph E. Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. PowerGraph: Distributed graph-parallel computation on natural graphs. In *Proceedings of the 10th USENIX conference on Operating Systems Design and Implementation*, OSDI'12, pages 17–30, Berkeley, CA, USA, 2012. USENIX Association.
- [23] Hadoop. Hadoop website <http://hadoop.apache.org>.
- [24] Maurice P. Herlihy and Jeannette M. Wing. Linearizability: A correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems*, 12(3):463–492, 1990.
- [25] HPCC. HPCC random-access benchmark http://icl.cs.utk.edu/hpcc/hpcc_results.cgi.
- [26] ISO/IEC JTC1/SC22/WG21. ISO/IEC 14882, Programming Language, C++ (Committee Draft). <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2008/n2800.pdf>, 2008.
- [27] Pete Keleher, Alan L. Cox, Sandhya Dwarkadas, and Willy Zwaenepoel. TreadMarks: Distributed shared memory on standard workstations and operating systems. In *Proceedings of the USENIX Winter 1994 Technical Conference*, WTEC'94, pages 115–131, Berkeley, CA, USA, 1994. USENIX Association.
- [28] Kai Li and Paul Hudak. Memory coherence in shared virtual memory systems. *ACM Trans. Comput. Syst.*, 7(4):321–359, November 1989.
- [29] Yucheng Low, Joseph Gonzalez, Aapo Kyrola, Danny Bickson, Carlos Guestrin, and Joseph M. Hellerstein. Distributed GraphLab: A framework for machine learning and data mining in the cloud. *PVLDB*, 2012.
- [30] Roberto Lubliner, Jisheng Zhao, Zoran Budimlic, Swarat Chaudhuri, and Vivek Sarkar. Delegated isolation. In *OOPSLA '11*, pages 885–902, 2011.
- [31] Grzegorz Malewicz, Matthew H. Austern, Aart J.C. Bik, James C. Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski. Pregel: A system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data*, SIGMOD '10, pages 135–146, New York, NY, USA, 2010. ACM.
- [32] Anirban Mandal, Rob Fowler, and Allan Porterfield. Modeling memory concurrency for multi-socket multi-core systems. In *Proceedings of the 2010 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS2010)*, pages 66–75, White Plains, NY, March 2010. IEEE.
- [33] Jacob Nelson, Brandon Myers, A. H. Hunter, Preston Briggs, Luis Ceze, Carl Ebeling, Dan Grossman, Simon Kahan, and Mark Oskin. Crunching large graphs with commodity processors. In *Proceedings of the 3rd USENIX Conference on Hot Topics in Parallelism*, HotPar '11, pages 10–10, Berkeley, CA, USA, 2011. USENIX Association.
- [34] Robert W. Numrich and John Reid. Co-array Fortran for parallel programming. *SIGPLAN Fortran Forum*, 17(2):1–31, August 1998.
- [35] Stephen Olivier, Jun Huan, Jinze Liu, Jan Prins, James Dinan, P. Sadayappan, and Chau-Wen Tseng. UTS: An unbalanced tree search benchmark. In *Proceedings of the 19th International Conference on Languages and Compilers for Parallel Computing*, LCPC'06, pages 235–250, Berlin, Heidelberg, 2007. Springer-Verlag.

- [36] James Reinders. *Intel Threading Building Blocks: Outfitting C++ for Multi-Core Processor Parallelism*. O'Reilly Media, 2007.
- [37] Ioannis Schoinas, Babak Falsafi, Alvin R. Lebeck, Steven K. Reinhardt, James R. Larus, and David A. Wood. Fine-grain access control for distributed shared memory. In *Proceedings of the sixth international conference on Architectural support for programming languages and operating systems*, ASPLOS VI, pages 297–306, New York, NY, USA, 1994. ACM.
- [38] Burton J. Smith. Architecture and applications of the HEP multiprocessor computer system. In *Proceedings of SPIE 0298, Real-Time Signal Processing IV, 241*, volume 298, pages 241–248, July 1982.
- [39] Dean M. Tullsen, Susan J. Eggers, and Henry M. Levy. Simultaneous multithreading: Maximizing on-chip parallelism. In *Proceedings of the 22nd Annual International Symposium on Computer Architecture*, ISCA '95, pages 392–403, New York, NY, USA, 1995. ACM.
- [40] Leslie G. Valiant. A bridging model for parallel computation. *Commun. ACM*, 33(8):103–111, August 1990.
- [41] Rob V. van Nieuwpoort, Thilo Kielmann, and Henri E. Bal. Efficient load balancing for wide-area divide-and-conquer applications. In *Proceedings of the eighth ACM SIGPLAN symposium on Principles and practices of parallel programming*, PPOPP '01, pages 34–43, New York, NY, USA, 2001. ACM.
- [42] Thorsten von Eicken, David E. Culler, Seth Copen Goldstein, and Klaus Erik Schauer. Active messages: A mechanism for integrated communication and computation. In *Proceedings of the 19th Annual International Symposium on Computer Architecture*, ISCA '92, pages 256–266, New York, NY, USA, 1992. ACM.
- [43] Kyle B. Wheeler, Richard C. Murphy, and Douglas Thain. Qthreads: An API for programming with millions of lightweight threads. In *IPDPS*, pages 1–8. IEEE, 2008.
- [44] Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, and Ion Stoica. Spark: Cluster computing with working sets. In *Proceedings of the 2nd USENIX Conference on Hot Topics in Cloud Computing*, HotCloud'10, pages 10–10, Berkeley, CA, USA, 2010. USENIX Association.