

# Design of a Multithreaded Barnes-Hut Algorithm for Multicore Clusters

## Technical Report

Junchao Zhang and Babak Behzad

Department of Computer Science, University of Illinois  
at Urbana-Champaign  
{jczhang, bbehza2}@illinois.edu

Marc Snir

Department of Computer Science, University of Illinois  
at Urbana-Champaign and MCS Division, Argonne  
National Laboratory  
snir@anl.gov

### Abstract

We describe in this paper an implementation of the Barnes-Hut algorithm on multicore clusters. Based on a partitioned global address space (PGAS) library, the design integrates intranode multithreading and internode one-sided communication, exemplifying a PGAS + X programming style. Within a node, the computation is decomposed into tasks (subtasks), and multitasking is used to hide network latency. We study the tradeoffs between locality in private caches and locality in shared caches and bring into the design the insights gained. As a result, our implementation consumes less memory per core, invokes less internode communication, and enjoys better load balancing strategies. The final code achieves up to 41% performance improvement over a non-multithreaded counterpart from our previous work. Through detailed comparison, we also show its advantage over other well-known Barnes-Hut implementations, in both programming complexity and performance.

**Keywords** Barnes-Hut, n-body, PGAS, cluster, multicore

### 1. Introduction

The current evolution of supercomputers exhibits several important trends [11, 26].

- The number of cores per node keeps increasing.
- The amount of memory per core is decreasing.
- One-sided communication (rDMA) is increasingly well supported on the interconnection networks.

The use of one-sided communication has several advantages. Such communication can have less software overhead, since code is executed on only one of the two communicating nodes; this results in lower latency, especially for short messages [25]. Also, it is easier to code irregular applications with dynamic communication patterns, using one-sided communication. In such an application, the consumer of a variable often knows the location of that variable, but the owner of the memory containing this variable may not know who is the consumer. Libraries such as SHMEM [6], Global

Arrays (GA) [20] and PGAS languages such as UPC [27] or CAF [13] use one-sided communication as their main communication mechanism.

The natural idiom for irregular applications with dynamic communication patterns is to use remote reads or get operations in order to access remote data. In order to achieve good performance, it is essential to hide the latency of the long round-trip of a remote memory access. Similar to simultaneous or concurrent multithreading in shared memory environments, such latency hiding is most conveniently achieved by descheduling tasks that are blocked on a remote access and reusing the core to run another, ready-to-execute task. This requires low-overhead task scheduling. Low-overhead task scheduling also enables efficient load balancing, ensuring that all cores are used.

Often, multiple cores on a node will use the same shared structures; memory pressure can be alleviated by keeping only one copy per node for such structures. In addition, codes such as the Barnes-Hut algorithm exhibit significant reuse of remote values; it is advantageous to keep a local copy of data brought from remote nodes, for possible reuse. Effectively, we use local memory as a software managed cache for remote memory.

We demonstrate in this paper the use of these techniques in the context of the Barnes-Hut (BH) algorithm. The main contributions of this paper are the following.

- We give the first BH design that integrates intranode multithreading and internode one-sided communication and uses multitasking to hide network latency.
- We study the interplay in BH between locality in private caches and reuse in shared caches.
- We compare BH implementations done using distinct programming models and discuss how they handle programming challenges on multicore clusters.

The rest of this paper is organized as follows. We describe the basics of the BH algorithm in Section 2 and then give details of our design in Section 3. In Section 4 we evaluate and analyze our design. In Section 5 we compare our code with two other BH implementations. Section 6 provides an overview of the library we are designing. Section 7 surveys related work. We conclude in Section 8 with a brief summary and look at future work.

### 2. Barnes-Hut Algorithm

The Barnes-Hut (BH) algorithm [5] is a fast algorithm for the n-body problem, which simulates the evolution of a system of  $n$  bodies (particles), where bodies apply forces on each other. A

simulation consists of multiple time steps. In each time step forces are computed and the bodies moved. A brute-force approach to this problem leads to  $\Theta(n^2)$  complexity. To lower the complexity, BH approximates the interaction of a body with a set of other bodies, by representing the set with a point located at its center of mass. The approximation is valid when the bodies in the set are “far enough” from the first body—with “far enough” being formalized as  $l/d < \theta$ , where  $l$  is the size of the cube containing the bodies in the set,  $d$  is the distance from the body to the center of mass, and  $\theta$  is a constant called the cell-opening criterion.

The BH algorithm partitions the 3D space hierarchically into *cells* using an octree representation. The root of the octree represents the cell that contains all bodies. Each cell is recursively divided into octants, until its number of bodies is below a fixed threshold. To compute forces on one body, the procedure begins with the root cell. If the current cell is far enough or contains only one leaf body, then we compute force with it and stop there. Otherwise, we *open* the cell and continue, recursively, with each of its children. With this hierarchical approach, the BH algorithm reduces the computation complexity to  $O(n \log n)$  [5].

*Load balancing* and *locality* are two important issues in parallel implementations of the BH algorithm. Since the input body distribution is usually nonuniform, each body interacts with a different number of cells. We cannot simply assign bodies to processors evenly. Also in BH, cells accessed during force computation for one body are likely to be accessed again for a nearby body. In systems with hierarchical memories, it is critical to allocate to one processor bodies close to each other in order to exploit this locality.

A shared memory implementation of BH from the SPLASH-2 benchmark suite [30] handles the issues as follows. It assigns each body a cost, which is the number of forces computed for this body. Since body locations change slowly per time step, one can use costs in the previous time step as estimates of costs in the current time step. SPLASH-2 BH uses an algorithm called *cost-zone* to assign bodies to threads. Octree leaves are split into  $p$  zones of consecutive leaves of roughly equal total cost, where  $p$  is the number of threads. Thread  $i$  picks bodies in the  $i$ th zone. A left-to-right traversal of the octree leaves corresponds to an ordering of the bodies along a space filling curve (SFC); we call this order *SFC order* [28]. Each thread is allocated a segment from the curve.

In a previous work [31], we implemented BH in UPC [27] on distributed memory, using one UPC thread per core. The implementation, which we call UPC BH from now on, inherited the ideas from SPLASH-2. A time step in UPC BH includes four phases, which are separated by barriers; each phase executes in parallel:

**Build Octree:** Threads build subtrees for subspaces assigned to them and compute the center of mass of each cell. Then threads hook subtrees together to form the global octree. For details, see [31].

**Partition Octree:** Each thread is assigned an array of bodies, according to the cost-zone algorithm.

**Compute Forces:** Threads compute forces for their bodies by traversing the octree from the root, then update their body costs. (See more details in Section 3.2.)

**Advance Bodies:** Threads advance their bodies by computing new velocities, positions, and so on. They also compute the boundaries of the new root cell for the next time step.

The force computation phase, which performs  $O(n \log n)$  operations, usually dominates the performance; other phases perform  $O(n)$  operations. This paper extends UPC BH, mainly in the force computation phase.

### 3. Multithreaded BH Design

UPC BH is a porting of SPLASH-2 BH to UPC [27]. A naive porting resulted in abysmal performance, but a sequence of optimiza-

tions resulted in dramatic improvement. In UPC BH, each core runs a persistent UPC thread, so that the code took only limited advantage of shared memory within nodes. This paper extends our previous work by using multithreading and multitasking within each process, for better latency hiding and load balancing and less off-node communication and memory use. Ideally, we would have implemented the new code in UPC again. Since UPC is based on C, however, it lacks support for generic programming, which is needed for us to abstract common services that are also useful for other applications. We therefore designed PPL, a C++ template library atop the Berkeley UPC runtime [1], and reimplemented BH in PPL. We refer to this implementation as PPL BH. For the time being, it is enough to know that PPL has the same memory model as UPC: Each process in PPL has its heap divided into two parts: a private heap and a local part of a global heap. While the private heap can be accessed only by threads local to the process, the global heap part can be accessed by any thread. Accesses to the local part of the global heap is much faster than accesses to remote parts. PPL has a generic global pointer structure that can point to any remote memory locations. Dereferencing global pointers may require remote reads. We use global pointers for the links in the BH octree. (We talk more about PPL in Section 6.)

In this section, we first introduce our test platform and test methodology, which will be used in experiments in this section and thereafter. Then we give an overview of the force computation in UPC BH for comparison. After that we describe PPL BH force computation in detail. At the end of this section, we study how to achieve both load balance and cache efficiency within a node.

#### 3.1 Test Platform and Test Methodology

We did all experiments on an x86 Linux InfiniBand cluster. Each compute node has two hex-core Intel Xeon 5650 CPUs running at 2.67 GHz. The six cores on a CPU have a private 32 KB L1 data cache and a private 256 KB L2 cache but share a 12 MB L3 cache. We used gcc4.6.3 as our C/C++ compiler and used the runtime of Berkeley UPC 2.14.2 as the communication library below PPL. The input bodies were generated by the Plummer model [2] as in SPLASH-2 BH. The octree was built with each leaf having at most 10 bodies. We ran 22 time steps and timed only the last 20 steps, with a time step = 0.025 seconds. All computations were done in double precision. The number of bodies  $n$  and the cell opening criterion  $\theta$  were varied. Since the focus of this paper is the force computation phase, which consumes most of the execution time, we report only average time per step of this phase.

#### 3.2 UPC BH Force Computation

In UPC BH, there is one UPC thread per core in a multicore node. The octree is distributed among threads. For load balancing, bodies are partitioned across all UPC threads using the cost-zone algorithm; there is no distinction between threads that belong to the same process and threads that belong to distinct processes. Every UPC thread is assigned an array of bodies sorted in SFC order with equal total cost. Two important optimizations in UPC BH are caching and computation/communication overlapping. Each thread caches cells visited, as these are likely to be reused to compute forces for subsequent bodies. The cached cells form a local partial octree, which is a snapshot of the global octree. We use pointer swizzling to have child pointers in a cached cell point to either its children in the original octree or cached copies of these children, depending on whether the children have been cached or not. The cached data is discarded at the end of the force computation phase. UPC threads do not share the cached data even when they are on the same node. If multiple threads on a node need the same off-node cells, they fetch them separately and communicate multiple times.

To overlap computation and communication, UPC BH takes advantage of the two-level parallelism in BH: Force computations for different bodies are independent and can be done in parallel; interactions between a body and different cells are also independent and can be done in parallel except that all forces acting on one body need to be summed together. At the beginning, each thread caches the octree root locally. Each thread maintains a work list of bodies. To compute forces for a body, a thread traverses the octree from the cached root. If a cell needs to be opened and its children are not cached, the thread will invoke a nonblocking communication to fetch the children; meanwhile the thread traverses other paths in the octree or just picks up another body from the work list. Threads periodically check pending nonblocking communications to complete them. For completed ones, threads resume interactions with the cells just fetched back.

### 3.3 PPL BH Force Computation

To make our description easier, we define some terms first. A cell is *localized* if its children have all been cached. Each cell has a `localized` flag to indicate whether it is localized or not; this flag is initially cleared. To localize a cell, we fetch its children, swizzle its global child pointers to local pointers pointing to the cached children, and then set the `localized` flag. The localization is a split-phase operation that includes making a nonblocking communication request and completing the request. Hence, we also add a `requested` flag in each cell to indicate that the cell has been requested but is not yet localized; this prevents making multiple requests for the same cell. The `requested` flags of all cells are initially cleared, too.

A *task* consists of the tree traversal and force computation for one body. During the traversal, if a cell needs to be opened but is not localized, we generate a *subtask* to handle the interactions between the body and the subtree rooted at that cell. Thanks to parallelisms in force computation, all tasks and subtasks can be executed in parallel, and synchronization is needed only to properly add together the forces acting on one body. Our objective is to orchestrate tasks and subtasks efficiently.

In PPL BH, usually there is one process per node. Upon entering the force computation phase, each process gets an array of bodies through the cost-zone algorithm and spawns one thread per core. This gives us the opportunity to allocate tasks and subtasks dynamically to threads and to share cached copies of cells across all threads. However, dynamic allocation and sharing can lead to increased synchronization overheads and reduced locality. (We study these tradeoffs in Section 3.4.) The general rule is that threads take bodies from the array, generate tasks, and execute them. If a thread is blocked in a task's execution by an *unlocalized* cell during tree traversal, it will generate a subtask to encapsulate the context and continue the traversal along other paths if possible; otherwise it will generate new tasks and execute them. Although subtasks from the same task can be executed simultaneously by distinct threads, we do not do so; instead, we execute a task and all its descendant subtasks on the same thread. The reason is that the large number of bodies is sufficient to ensure that threads are always busy; moreover, this choice avoids the need to synchronize reductions across multiple threads.

PPL BH has the same optimizations as does UPC BH: caching and computation/communication overlapping. We considered three ways of splitting work across threads, as shown in Figure 3.1.

**(1) Equal & Centralized** (Figure 3.1(a)): In the first approach, all threads are equally involved in computation and communication. As shown, all threads share a map (`c2s_map`), which stores, for each cell, the list of subtasks blocked on an open request for that cell. The map acts as a hub for subtask registering and releasing. When a thread wants to open an unlocalized cell, it gen-

erates a subtask and registers it in the map under that cell. If the cell is not requested, the thread also makes a nonblocking request to fetch children of the cell and puts a handle to the request in a list (`pending_requests`). Threads periodically check the list to see whether any request is completed. For completed requests, they will push registered subtasks to runnable subtask queues (`runnable_subtaskq`) on threads. Each subtask carries a thread id (`tid`) that indicates the thread responsible for executing the subtask—hence the queue the subtask joins when it becomes runnable. Threads query their own `runnable_subtaskq` whenever the task they currently execute becomes blocked.

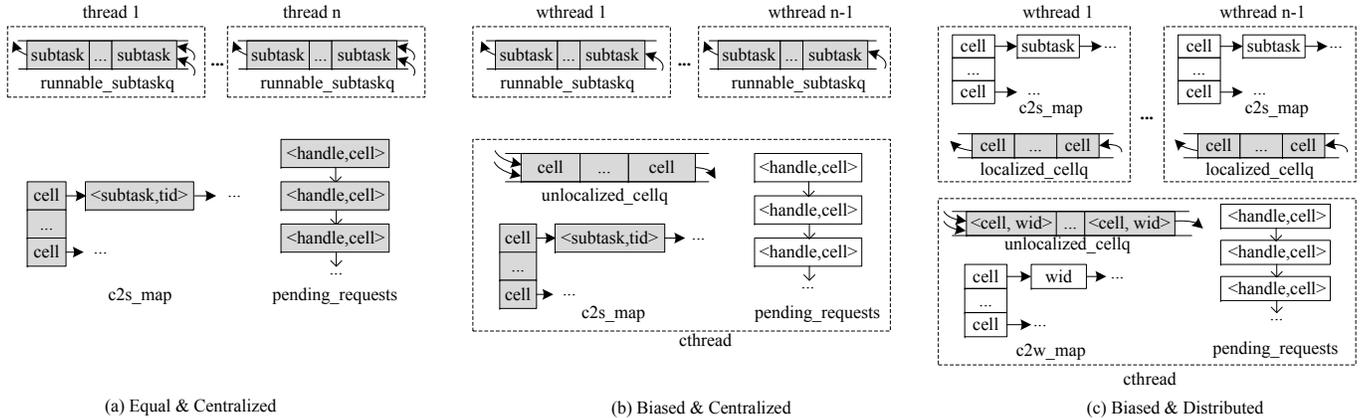
The biggest problem of this approach is synchronization. First, accesses to the `localized` flags and the `requested` flags must be atomic so that all threads have a consistent view of the cached octree. Second, all data structures, such as `c2s_map` and the `pending_requests` queues, must be concurrent. The overall overhead could be high even with an efficient implementation of these data structures.

**(2) Biased & Centralized** (Figure 3.1(b)): In the second approach, one thread is designated as the communication thread (`cthread`), with the remaining as worker threads (`wthreads`). All threads share the concurrent map (`c2s_map`) as in the first approach. Only the `cthread` is in charge of communication. Thus, communication management is easier, and data structures such as `pending_requests` need not be concurrent anymore. Every `wthread` has a runnable subtask queue. When a `wthread` wants to localize a cell and finds it was already requested, the `wthread` registers a subtask for the cell in `c2s_map`. Otherwise, if the cell was not requested previously, the `wthread` will also mark the cell as requested and push it in a concurrent queue (`unlocalized_cellq`). The operation of marking cells as requested needs not be atomic. The same cell can be pushed into `unlocalized_cellq` multiple times by different threads.

The `cthread` pops cells from `unlocalized_cellq`, checks and updates their `requested` flags again to remove the redundancy, and issues only one request per cell. The `cthread` also checks pending requests. If a request is completed, it marks the corresponding cell as localized, looks up `c2s_map`, pushes subtasks registered under the cell back to the queues (`runnable_subtaskq`) of their owner `wthreads`, then deletes the entry in `c2s_map`. Each subtask has a worker id (`wid`) so that the `cthread` knows which queue to choose. The `cthread` and `wthreads` must be properly synchronized if they are operating on the same entry in `c2s_map`. For example, when a cell in `c2s_map` was deleted by the `cthread`, no `wthread` should have chance to insert it again. In this approach, every `runnable_subtaskq` is now a single-producer single-consumer (SPSC) queue while `unlocalized_cellq` is a multiproducer single-consumer (MPSC) queue.

This approach is similar to thread scheduling in an operating system: Multiple threads can wait for a same signal. Once the signal arrives, all threads registered under this signal are woken up.

**(3) Biased & Distributed** (Figure 3.1(c)): This approach has the same `cthread` and `wthreads` as in the second approach. But this time the cell-to-subtask map (`c2s_map`) is distributed among threads and is not concurrent anymore. If a `wthread` wants to localize a cell, it looks up its private `c2s_map` to see whether the cell has already been requested *by itself*. If so, it just registers the subtask in its map; otherwise, it also pushes the cell along with a worker id (`wid`) into a concurrent queue (`unlocalized_cellq`). An unlocalized cell may be pushed into the queue multiple times by different threads. But the `cthread` makes only one communication request for each cell. The `cthread` uses a private map (`c2w_map`) to map cells to `wthreads` that have requested them. The `cthread` manages communication. When a request is completed, it marks the cell as localized and pushes it back to the `wthreads`' queue (`localized_cellq`).



**Figure 3.1.** Data structures used in the three approaches. Shaded structures are accessed concurrently.

On the other side, wthreads pop cells from their queue, look up their map, and execute subtasks registered under the cells. Note that accesses to the `localized` flags of cells need not to be synchronized between the `cthread` and wthreads. If the `cthread` set the `localized` flag of a cell and the new value is not immediately observed by a wthread, the wthread may superfluously push the cell into `unlocalized_cellq`. When the `cthread` pops up the cell, it will check the cell's flag. Of course, the `cthread` will find the flag is true because it was set by itself before. If the flag is set, the `cthread` just rebounds the cell back to its owner wthread. Sooner or later, the new value will be seen by wthreads, thanks to the hardware cache coherence protocol. It turns out that all synchronizations in this approach can be done through either SPSC queues or MPSC queues, which can be implemented efficiently by using lock-free data structures [12].

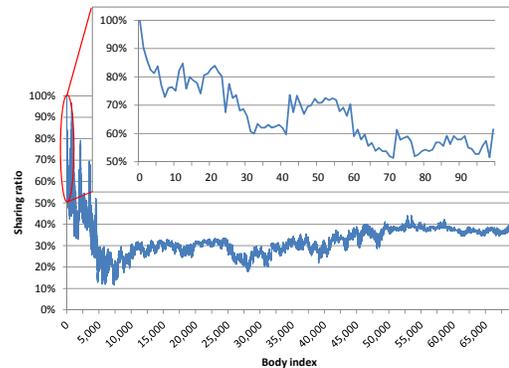
In our code, all maps are implemented as hash maps. We use pointers to cells as keys for the hash maps, which we found to be efficient in practice.

### 3.4 Intranode Task Scheduling

We now study how to schedule tasks within a node, in other words, how to distribute bodies to threads, with the aim of achieving both load balance and cache efficiency. A distinguishing feature of multicore CPUs is that they usually have their last level cache shared among cores on the chip. Hence, besides private cache efficiency, it is interesting to know whether we can improve shared cache efficiency by synergistic task scheduling.

There are two common distribution strategies to start with: cyclic distribution and block distribution. In the former, body  $i$  will be assigned to core  $i \bmod p$ , where  $i$  is the body index in the body array and  $p$  is the number of cores. In the latter, each core is assigned a block of consecutive bodies from the body array. The locality in BH says that cells visited by one body are very likely to be visited again by nearby bodies. We did the following experiments to measure the locality in these two distributions.

We ran UPC BH with  $n = 1M, \theta = 0.5$  and 16 processes (threads in UPC terms), and examined process 8 (P8) in the fourth time step. The process was assigned about 69K bodies. We tagged cells visited during force computation for the first body. Then we found out how many of them were visited again by the second body, the third body, and so on and calculated a sharing ratio for each body with respect to the first body. The result is shown in Figure 3.2, which also includes a zoom-in picture for the first 100 bodies on P8. We can see that BH has very good locality when SFC order is used. For example, the second body shares about 90% of the cells



**Figure 3.2.** Sharing ratio for the 69K bodies on P8.

the first body visited. Even the 100th body still shares more than 50% with the first body. The sharing ratio decreases with distance, but the trend is not strict. The reason is that distance in a 1D SFC does not strictly correspond to distance in a 3D space. The curve does not hit zero because top cells of the octree are visited by all bodies. So there is still a good sharing (20–40%) even at distant parts of the curve. Curves on other processes have similar shapes, so we do not show them here.

Using the same concept, we studied sharing between blocks. The configuration was the same. But this time we cut the 69K bodies into 12 blocks of equal total costs (note that we have 12 cores per node). We tagged cells visited by the first block of bodies and calculated the sharing ratio of the second block, the third block, and so on with respect to the first block. Figure 3.3 shows the block-sharing ratio curves measured on processes 7 and 8. Compared with Figure 3.2, the ratio drops greatly. Hence, blocks of bodies will touch different bottom parts of the octree so that they have less sharing than before. We can also notice that the ratio varies a lot between different processes, which indicates different regions of the space have distinct sharing property. But generally the ratio is low.

The above experiments suggest that block distribution should benefit L1 cache, since computations on successive bodies will reuse the same cells. By the same token, cyclic distribution should

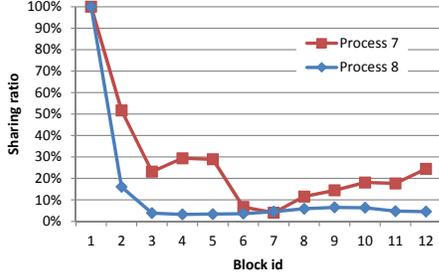


Figure 3.3. Block sharing ratio on P7 and P8.

Table 1. Average traversal length and octree size in number of cells

n		1M	2M	4M	8M	16M
Traversal length	$\theta=0.5$	1973	2074	2178	2280	2375
	$\theta=0.7$	853	892	928	962	996
Octree size		$1.4 \times 10^6$	$2.7 \times 10^6$	$5.4 \times 10^6$	$1.1 \times 10^7$	$2.2 \times 10^7$

benefit L3 cache, since cached cells are likely to be accessed by all threads in a small time interval. To test this hypothesis, we modified SPLASH-2 BH (itself uses block distribution), added cyclic distribution, and ran it with 6 threads on one CPU of our platform. We bound threads to cores and did cache profiling with PAPI [19]. Table 3 shows the miss rates for both cases. Note that the table uses *local miss rates* at each level, defined as the miss count of this cache divided by the reference count of this cache. We took the average miss rate on six cores. Surprisingly, no significant differences between cyclic distribution and block distribution can be seen in either performance or miss rate. L1 miss rates of both are very high, while L3 miss rates of both are fairly low. (In the table, L2 miss rates are high because most of the locality in L2 is absorbed by L1, which makes L2 miss rates less interesting here.)

To understand this, in the fourth time step, we measured the octree traversal length per body and the total octree size with various  $n$  and  $\theta$ , in number of cells. The traversal length of a body is the number of cells visited by the body during its traversal. Table 1 shows the results.

From Table 1 we can see that, as expected, the tree size is proportional to the number of bodies, while the traversal length is proportional to the log of the number of bodies. In the code, the average size of a cell or body is about two cache lines (i.e., 128 bytes, including fields for 3D positions, mass, child pointers, etc.). So, with  $n = 1M, \theta = 0.5$ , the average traversal size is about 246 KB. The number explains why in the previous block distribution test the L1 miss rate was high: since the size of a traversal exceeds L1 capacity (32 KB), accesses to cells in the second traversal generate L1 misses. It also explains why the two experiments have nearly identical L3 miss rates: since the size of a traversal is much smaller than that of L3, it has enough capacity to cache simultaneously many successive traversals of each thread, even if there is less sharing between threads. So we can conjecture that the benefit of cache sharing will show up only when the size of a single traversal approaches the L3 quota per core (i.e., 2 MB in our case). We could increase  $n$  to meet the condition, but  $n$  would be an incredibly large number. Instead, we decreased  $\theta$ . Table 2 shows the result for  $n = 1M, \theta=0.2$  on one CPU. (Note that  $\theta=0.2$  is not a realistic choice for BH; we use it only to test the conjecture.) In this test, the average traversal length is about 21,260 cells, or 2.6 MB. Now we can observe a big difference between cyclic distribution and block distribution: the former’s L3 miss rate is 1.25%, while the latter’s is 15.93%. That results in a big performance difference, as also shown in Table 2.

Table 2. Performance when traversal size > L3,  $n = 1M, \theta = 0.2$

	Time(s)	L1 Miss	L2 Miss	L3 Miss
Cyclic	108.88	7.33%	53.32%	1.25%
Block	142.36	6.91%	56.66%	15.93%

Table 3. Performance when traversal size < L3,  $n = 1M, \theta = 0.5$

	Time(s)	L1 Miss	L2 Miss	L3 Miss
Cyclic w/o SFC	12.66	8.19%	76.77%	7.65%
Cyclic	10.59	8.24%	69.73%	0.14%
Block	10.84	8.23%	69.46%	0.14%
Block w/ Tiling	8.28	0.20%	74.14%	6.36%

These experiments suggest that with reasonable input parameters and hardware, it is *not important* to exploit cache sharing on multicores in BH. In contrast, the important thing is to improve L1 efficiency. To verify that conclusion, we implemented the tiling technique in [14] in SPLASH-2 BH. Block distribution is still used, but each thread now takes a tile of bodies at a time from its block. When visiting a cell, bodies in the tile will interact in turn with that cell. Some bodies may need to open the cell, while others may not. Therefore, we need to dynamically mask member bodies in a tile. As a result, the octree is traversed only once for all bodies in a tile. In this way, when a cell is brought into L1 by one body, it will be reused by other bodies in the same tile who need it. As we saw in Figure 3.2, the reuse probability is high. Row “Block w/ Tiling” in Table 3 shows the result with a tile size of 128. We can see the dramatic L1 miss reduction and performance improvement. Note that the L3 miss rate seems high after tiling. The reason is that the L3 reference count drops greatly (not shown in the table). In reality, the L3 miss count did not increase much. We tried different tile sizes and found that the performance is not sensitive to size when it is in range of 64 to 10K. If bigger than that, the tile itself will overflow the L1, lowering performance. We therefore choose a tile size of 128 hereafter. Table 3 has a row “Cyclic w/o SFC,” which is also interesting here. In this case, we used cyclic distribution, but bodies in the input array were randomly ordered. We see a high L3 miss rate and degraded performance, because the BH locality was not respected. Different parts of the octree are now randomly touched by unrelated bodies, resulting in a much bigger memory footprint in L3 than before.

We now see that we should increase the task size in PPL BH. Fortunately, doing so requires only small changes in the algorithms described in Section 3.3. Now a task computes forces for a tile of bodies. A task has a bit mask to indicate which bodies need to interact with a cell. When a body in a tile needs to open an unlocalized cell, the task generates a subtask and relay, the bit mask to it, so that when the subtask is resumed, it knows which bodies to pick up.

Note that the distributions we discussed are static in the sense that bodies are assigned to threads before force computation starts. This approach is fine for the shared-memory SPLASH-2 BH but not good for the distributed-memory PPL BH since, besides computation, there is communication. With a static block distribution, different threads need different amount of remote data, resulting in load imbalance. To smooth this variation, we adopted dynamic scheduling in PPL BH. We tried two approaches similar to OpenMP *dynamic* scheduling and *guided* scheduling for work sharing loops [21]. In the former, worker threads request a fixed-sized chunk of bodies from the body array, work on the chunk, then request another chunk until no bodies are left. In the latter, worker threads request a chunk of bodies with length proportional to the number of unassigned bodies divided by the number of worker

threads. We chose *guided* since it shows a little better performance than *dynamic*.

## 4. Evaluation and Analysis

In this section, we test and analyze PPL BH with the various multithreading approaches discussed in Section 3. Because of the obvious synchronization overhead in the Equal & Centralized approach (we can see it even in the Biased & Centralized approach), we implemented only the last two approaches. We quantitatively measure benefits of PPL BH and then present its performance.

### 4.1 Comparison of the Two Biased Approaches

We tested PPL BH with the last two biased approaches with various configurations. We note that the performance of the Biased & Centralized approach is 1.6 to 3.5 times worse than that of the Biased & Distributed approach. Both use a data-driven style: subtasks become runnable only when their requested data arrives. But the first biased approach uses a centralized cell-to-subtask map: all subtasks are registered in the map at the cthread, which notifies wthreads which subtasks become runnable. In the second biased approach, when a cell is localized, the cthread notifies the wthreads that have requested the cell. Then those wthreads look up their private maps to release subtasks depending on this cell. Although this approach can result in the same cell being registered at multiple threads, it reduces the pressure on the cthread, as the wthreads filter requests for the same cell.

We did profiling with  $n = 1M$ ,  $\theta = 0.5$  on 8 nodes. In the Biased & Centralized approach without tiling, on average dozens to hundreds of subtasks from wthreads are going to be registered under one cell in the cthread’s map; the maximal number is huge, ranging from 8,000 to 31,000. In the Biased & Distributed approach, however, maximally 6 to 11 wthreads request the same cell from the cthread simultaneously; on average, the number is 1. Clearly, through distributed subtask management we saved much of the traffic between the cthread and wthreads and thus achieve better performance. With tiling, the number of tasks/subtasks decreases so that the phenomenon is not that significant. The centralized still lags, however, likely because of the synchronization overhead in `c2s_map` (note that it is a concurrent map). Because of its superior performance, we will henceforth use PPL with the Biased & Distributed approach.

### 4.2 Less Memory Consumption

The octree in UPC BH is distributed among threads. Each thread has a local part of the octree (which we call the *local tree*). Local trees are linked together by shared pointers. In force computation, threads traverse the octree and cache visited cells. The part of the octree visited during force computation for all bodies assigned to a thread is called the thread’s *locally essential tree* (LET). Table 4 shows the average ratio of the LET size to local tree size while varying  $n$ ,  $\theta$ , and the number of threads.

Let us look at one configuration, namely, 16 nodes, each with 12 cores, for a total of 192 threads for UPC BH. With  $n = 1M$ ,  $\theta = 0.5$ , the ratio of LET size to local tree size is 3.58. In PPL BH, however, we create only one process per node. Threads spawned by a process share the same LET; this will be the same LET created in UPC BH when one runs 1 thread per node, where the ratio is 1.80. It means that PPL BH saves about half of the memory consumed by cached remote cells by sharing the LET. By the same reasoning, we see that for  $n = 1M$ ,  $\theta = 0.5$ , the percentages of memory saved through multithreading with 32 nodes and 64 nodes are  $(4.68 - 2.12)/4.68 = 55\%$  and  $(6.37 - 2.58)/6.37 = 59\%$ , respectively. We could expect bigger savings when more cores are put on a chip. From Table 4, we notice that when  $n$  increases, the

**Table 4.** Ratio of locally essential tree (LET) size to local tree size

No. of Threads	$n = 1M$		$n = 2M$		$n = 4M$	
	$\theta = 0.5$	$\theta = 0.7$	$\theta = 0.5$	$\theta = 0.7$	$\theta = 0.5$	$\theta = 0.7$
16 (1x16)	1.80	1.43	1.65	1.37	1.51	1.29
32 (1x32)	2.12	1.59	1.91	1.49	1.69	1.38
64 (1x64)	2.58	1.81	2.21	1.64	1.92	1.49
192 (12x16)	3.58	2.28	2.90	1.96	2.42	1.73
384 (12x32)	4.68	2.79	3.64	2.30	2.91	1.96
768 (12x64)	6.37	3.54	4.76	2.81	3.66	2.30

**Table 5.** Percentage of off-node communication saved in PPL BH with respect to UPC BH

No. of Nodes	$n = 1M$		$n = 2M$		$n = 4M$	
	$\theta = 0.5$	$\theta = 0.7$	$\theta = 0.5$	$\theta = 0.7$	$\theta = 0.5$	$\theta = 0.7$
16	44%	33%	36%	22%	28%	13%
32	51%	42%	43%	32%	36%	25%
64	57%	49%	50%	40%	43%	32%

ratio decreases; with bigger  $\theta$ , which translates into less accuracy and fewer cells opened, the ratio will also decrease.

### 4.3 Less Off-Node Communication

Each thread in UPC BH has its own LET, and the off-node remote cells fetched by one thread will not be reused by other threads on the same node. In PPL BH, this is not a problem anymore. We can quantitatively measure the saving in communication. For example, given 16 nodes and 12 cores per node, we run UPC BH with 192 threads. We distinguish on-node cells and off-node cells fetched during force computation and sum the number of off-node cells fetched by each thread. Then we run PPL BH on 16 nodes and collect the same number. Comparing these two numbers, we can know how big the saving is. Table 5 shows the savings on 16, 32, and 64 nodes. For example, with  $n = 1M$ ,  $\theta = 0.5$ , and 64 nodes, PPL BH saves about 57% of the off-node communication compared with UPC BH. In strong scaling, with more nodes, the saving is bigger; with other parameters fixed, increasing  $n$  or  $\theta$  reduces the saving.

### 4.4 Better Load Balancing

Since the force computation phase is synchronized between processes, its execution time is determined by the longest process. As we know, UPC BH inherited from the shared memory SPLASH-2 BH code the cost-zone load-balancing algorithm. However, this algorithm is computation-centric. On distributed memory the need to access remote cells can disturb the balance. Because of SFC ordering, boundary processes on a node usually require more remote cells than do interior processes. Considering computation/communication overlapping, the effect is hard to estimate upfront and thus is better attacked by dynamic scheduling enabled by multithreading. Let us look at an example. With  $n = 4M$ ,  $\theta = 0.5$  and 64 nodes, we measured the execution time variation of the 12 threads, using the formula  $(MaxTime - AverageTime)/AverageTime$  on each node. For UPC BH, it ranges from 0.5% to 71.0%. For PPL BH, however, it ranges only from 0.0% to 2.5%. Obviously, PPL BH achieves better intranode load balancing.

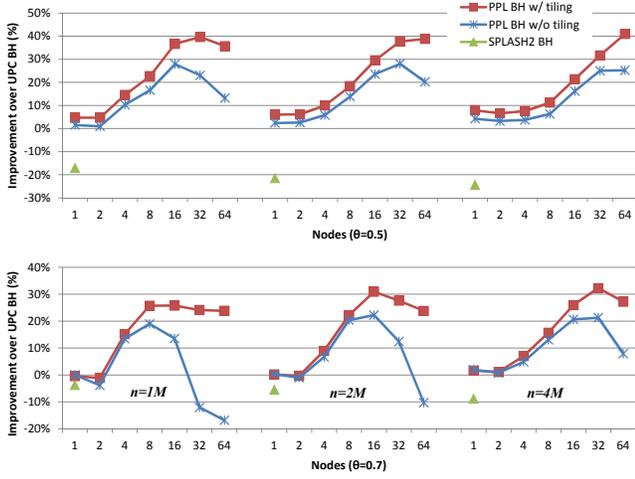


Figure 4.1. Performance improvement over UPC BH.

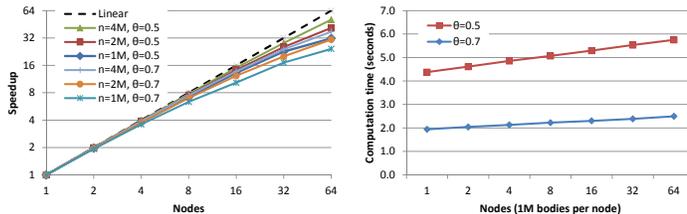


Figure 4.2. Strong scaling (left) and weak scaling (right) of PPL BH.

#### 4.5 PPL BH Performance

We used our best UPC BH implementation as a baseline and compared PPL BH with it. The performance improvement is shown in Figure 4.1, which also includes SPLASH-2 BH’s performance on one node. Surprisingly, UPC BH and PPL BH have much higher single-node performance than does SPLASH-2 BH, even though they perform extra operations. Note that we cache cells even if they are in the local part of the octree. (This step is necessary because otherwise accesses from remote processes to the local part would read invalid swizzled pointers.) While caching the octree entails additional memory copies, it reduces L1 misses, probably because of the better memory layout. Thus, both UPC BH and PPL BH perform better than SPLASH-2 BH even without tiling. Also, PPL BH significantly outperforms UPC BH. The improvement is larger for higher node counts. For example, with  $n = 4M$ ,  $\theta = 0.5$  and 64 nodes, we get the highest improvement, at 41%. We can also observe the trade-off between the benefits of multitasking and their overhead, especially when the computation density degrades. For example, with  $n = 1M$ ,  $\theta = 0.7$  and 64 nodes, without tiling, PPL BH degrades about 17%. From Table 1, we know that the average cost of a body decreases by more than half when  $\theta$  is changed from 0.5 to 0.7. With tiling, however, by having better locality and less tasks, PPL BH’s performance goes up.

Figure 4.2 shows scaling of PPL BH with tiling. For strong scaling, using performance of one node as baseline, at 64 nodes,  $n = 4M$ , it achieves a speedup of 51, 38 for  $\theta = 0.5$  and 0.7, respectively. For weak scaling, we keep 1M bodies per node, and the computation time increases logarithmically, as expected.

## 5. Comparison with Other BH Implementations

BH is a challenging application that has drawn much attention. It is worth comparing implementations done in different programming models and analyzing how they deal with critical programming issues such as overlapping computation and communication, balancing load, preserving locality, and using multicore nodes. In this section, we compare PPL BH with two other codes, PEPC and Charm++ BH, which are implemented in distinct programming models. We describe and compare each with PPL BH and then summarize.

### 5.1 PEPC

The Pretty Efficient Parallel Coulomb (PEPC) solver<sup>1</sup>, written in MPI + Fortran 2003, is a state-of-the-art parallel BH code developed and widely used at the Jülich Supercomputing Center. PEPC was reported to be able to run with up to 2 billion particles and 256K cores [29]. PEPC has multiple front-ends for various disciplines. They all share the same kernel—a tree code implementing the *hashed octree* scheme pioneered by Warren and Salmon [28]. The scheme computes unique keys for cells and particles according to their positions and then stores them in a hash table. Particles are weighted by their costs as in SPLASH-2 BH and sorted by their keys, resulting in an SFC. The curve is then partitioned among processors to achieve both locality and load balancing. MPI two-sided communication is used to fetch remote cells during tree traversal. The fetching process sends keys of parents to remote processes, which compute child keys, look up their hash tables, and send the children back. Processes do periodic synchronization to exchange data. This approach is well explained in [28].

As multicore cluster emerges, PEPC ships a hybrid MPI + Pthreads tree code [29], for the same reasons as we do in PPL BH. Its design is similar to PPL BH but without one-sided communication, multitasking and tiling. Usually there is one process per CPU, which spawns multiple *worker threads* while continuing as a *communicator thread*. Workers dynamically grab fixed-sized chunks of particles from the list of particles assigned to this process and compute forces for them. For each particle, workers maintain a *todo list*, which contains locally available cells to interact with, and a *defer list*, which contains unlocalized parent cells. Worker threads push into queues at their communicator threads requests for child cells. Communicator threads communicate with each other to satisfy these requests. Once child cells are sent back, communicators cache them locally and tags their parents in the hash table. Workers have to periodically poll tags of cells on defer lists and try to move their children to todo lists. PEPC allows overbooking cores. Communicator threads can periodically yield cores to workers, a function that is not implemented in PPL BH.

Hash keys in PEPC function as global pointers in PPL BH, except that dereferencing keys needs hash table lookups and participation of remote threads. In PPL BH, with a global name space and one-sided communication, these added complications are avoided, resulting in less programming complexity and lower runtime overhead. Also, PPL BH abstracts computation into tasks/subtasks and uses a data-driven style task scheduling in contrast to tag polling in PEPC.

### 5.2 Charm++ BH

Charm++ BH, written in Charm++ [17], is a parallel BH code developed at the University of Illinois. We took the code from the Charm++ benchmarks, which won the 2011 HPC Challenge Class 2 [16]. Charm++ is a C++-based parallel programming system that implements a programming model based on message-driven, migratable objects. It features measurement-based auto-

<sup>1</sup> <http://www.fz-juelich.de/ias/jsc/pepc>

matic load balancing and automatic computation/communication overlapping, through overdecomposition. The migratable objects are called *chares* in Charm++. Chares are activated by remote invocations and execute without preemption. Communication latency is hidden by having multiple chares for each core. The Charm++ runtime can instrument chares, measure their execution time and migrate them among processors to balance their loads.

Charm++ provides an SMP mode. When enabled, it spawns multiple threads within a node (process, actually). Each thread becomes a *processing element* (PE) that handles a set of chares. Ideally, Charm++ encourages programmers to think of chares as virtual processors, so that the code is largely independent of the number of physical processors and the number of cores within each. But, for certain optimizations (such as data sharing and message reduction in BH), this ideal model is not feasible. Charm++ provides two language constructs: *group* and *nodegroup*, which are collections of chares; there is one chare per PE in a group, and one chare per *process* in a nodegroup. Lacking a global name space, Charm++ BH adopts the hashed octree scheme again and shares its weaknesses.

To lower runtime overhead, Charm++ BH creates a chare for a group of particles instead of one. During tree-building, Charm++ BH creates an auxiliary space partitioning tree (similar to the top part of an octree). Each leaf represents a subcube of the space, enclosing a number of particles that is below a threshold set by the user. Each such leaf is handled by one chare (named *TreePiece*, or *TP*). There are usually dozens of TPs per PE. TPs build local partial octrees with their particles and compute forces for them by traversing the entire octree. During traversal, TPs may need to access other TPs on the same node. To avoid this intranode communication, Charm++ BH designs a nodegroup (*TreeMerger*). When the SMP mode is enabled, it merges all local trees in a node and forms a larger local octree, which is then shared by all TPs in the node. TPs on the same PE may also need to access the same remote cells. To save this duplicate internode communication, Charm++ BH designs a group (named *DataManager*, or *DM*). The TPs are similar to our worker threads, and the DM is similar to a communication thread, except that the TPs and DM on a PE are executed by a single thread (i.e., the PE itself). The DM maintains a software cache so that duplicate off-node requests from the same PE are screened out. But note that DMs are per-PE objects. Duplicate requests from different PEs on the same node are not filtered. Table 5 showed that this approach can result in a significant amount of superfluous communication. We believe Charm++ BH could design the DM at node level to remedy that, but only after it handles thread synchronization problems as we discussed and fixed in PPL BH.

The Charm++ load balancer is triggered every few time steps. At the end of such steps, the balancer computes the center of mass of the TPs and weights these mass points with the TP costs measured by the runtime. It then maps the mass points (hence TPs) to PEs using the well-known locality-preserving orthogonal recursive bisection (ORB) method [24]. With a new TP-to-PE map, chare migration is triggered. The approach is nice, but we found an issue with it in our experiments. Because of particle movement, the space partitioning tree (hence the TPs) can change at each step. If this happens, it means we will use an outdated TP-to-PE map until the next balancing step, with imperfect load-balancing. It is not clear how to balance the size of the TPs and the frequency of load balancing so as to optimize performance.

Thanks to overdecomposition, each PE has many TPs. Particles are sorted in an SFC order as in PEPC, and each TP owns a segment of the curve. In Charm++ BH, all TPs are active objects. A TP on a PE periodically yields the core to give its partners a chance to run. From the PE's view, however, particles are handled in a somewhat

**Table 6.** Computation time(s) of PEPC-mini, Charm++ BH and PPL BH,  $n = 1M$ ,  $\theta = 0.5$

Nodes	1	2	4	8	16	32	64
PEPC-mini	15.86	7.03	3.59	1.98	1.00	0.64	0.57
Charm++ BH	6.27	3.83	2.24	1.45	1.03	0.74	0.63
PPL BH	4.38	2.22	1.14	0.61	0.33	0.19	0.14

arbitrarily order. As we have shown in Section 3.4, this damages cache locality. It is not clear how the scheduling order could be controlled to improve locality.

### 5.3 Comparison

To squeeze performance from BH on multicore clusters, all implementations did nontrivial work. But by leveraging intranode multitasking and internode one-side communication, we encapsulated many complicated issues clearly and gave a simple but high-performance design. For simplicity, we use lines of code to indicate code complexity. They are about 3600, 8000, 25000 lines of code in PPL BH, Charm++ BH, and PEPC-mini (mentioned later), respectively. To measure performance, we made all implementations compatible to SPLASH-2 BH so that we could use the same parameters and input files.<sup>2</sup> For PEPC, we started from PEPC-mini, a skeleton molecular dynamics front-end of PEPC. With little effort, we changed its interaction from Coulomb to gravitation and its expansion from multipole to monopole and got a code comparable to ours. The critical PEPC tree code is not modified. We compiled PEPC-mini with OpenMPI-1.6.4 and the Intel Fortran Compiler 13.1. We ran it with two processes per node, six workers per process (since this configuration achieved the best performance). For Charm++ BH, we used Charm++ 6.4 and triggered the ORB load balancer every five time steps. Test results for  $n = 1M$ ,  $\theta = 0.5$  are shown in Table 6. We varied  $n$  and  $\theta$ , but the relative performance did not change much. We can see PPL BH's impressive performance advantage over the two other codes. We believe the heavy hash table lookups hurt PEPC's performance.<sup>3</sup>

## 6. PPL Library

Through a concrete application, the Barnes-Hut algorithm, we demonstrated the potential of a parallel programming style that combines intranode multitasking with internode one-sided communication and uses task preemption to hide communication latency. Although we worked on only one application, we believe that many of the abstractions (e.g., wthreads/cthread separation, queue-based synchronization, distributed task management, tagged tree nodes) in our design can be reused for other applications. We are in process of designing a library, PPL, to facilitate this kind of programming. PPL is designed as a C++ template library atop one-sided communication. We borrow ideas from the Intel Threading Building Block (TBB) [23]. At the top of the library are parallel algorithms such as `parallel_for()`, `parallel_reduce()`, which are used by programmers to express parallelism in their applications. At the middle is the task scheduler, which does load balancing and latency hiding primarily through multitasking. At the bottom are global data structures in a PGAS memory model. Currently, we have implemented these primitives in PPL:

*global variable:* `template <typename T> class gvar`. A `gvar`'s home is on process 0, but all processes have a local copy of it. The local copy is updated by a call to `gvar.cache()`. An

<sup>2</sup> With one exception, leaves in PEPC's octree can contain only one particle. We configured PPL BH accordingly and found small performance variations, so we ignored this.

<sup>3</sup> The JSC group is working on this issue.

assignment to `gvar` updates the copy on process 0 – there is no coherence protocol; the user is responsible to avoid data races. This construct captures the write-once read-many access pattern common in applications. We used a `gvar` to store the octree root.

*global vector:* `template <typename T> class gvec.` A `gvec` is like a co-array in Fortran 2008: every process has a vector of the same size. A process can access elements of remote vectors. Collective operations are defined on `gvec`. We used a `gvec` to store costs of subspaces in octree-building [31].

*global pointer:* `template <typename T> class gptr.` A global pointer can point to any memory locations in the global heap of any process. Dereferencing a global pointer may incur remote reads or writes. We used global pointers to link octree cells.

Generic high-level data structures such as trees can also be defined in a global address space. However, an important departure from data structures on shared memory is that we should also define caching properties for global data structures on distributed memory, like the one we showed in `gvar`, such that variables can be cached without changing names to reference them, a feature needed for productivity but lacking in current PGAS languages. Besides TBB Task's standard `execute()` interface, tasks in PPL also provide interfaces such as `IsDataReady()`, `ExtractRemoteAddress()` to let the runtime know whether the needed data is locally available and, if not, what the global address of the data is so worker threads can forward their requests to communication threads. We are now refining interfaces of PPL.

## 7. Related Work

The first parallel BH algorithm on distributed memory was developed by Salmon [24]. This work pioneered the locally essential tree (LET) method, which we can think of as a workaround for the problem caused by lack of a global address space and irregularity in BH. A LET for a process is the part of octree that will be traversed during force computation for bodies on this process. In this method, processes estimate which cells in their local octree are needed by other processes by using a relaxed cell opening criterion, then exchange cells and build their LET. With a LET in hand, force computation can proceed without communication. Obviously, this method loses the opportunity to overlap computation and communication. As a remedy, Warren and Salmon came up with another method, the hashed octree scheme [28], as we discussed in Section 5.1.

Truong Duy et al. [10] presented a hybrid MPI + OpenMP BH implementation. They used OpenMP to parallelize the *for* loop for force computations for bodies on a node, in order to avoid intranode communication and achieve better load balancing, as we also do in PPL BH. They also tested various OpenMP schedules for work-sharing loops and found *dynamic* was slightly better than *static* and *guided*. A similar work appeared in [22], which implemented a multipurpose n-body code with MPI + OpenMP hybrid programming. However, both works are based on the LET method, sharing its inefficiency.

Dinan et al. [9] introduced a hybrid parallel programming model that combines MPI and UPC. This model consists of UPC groups, and the intergroup communication is done through MPI. Processes in a UPC group can access each other's shared heap as normal UPC programs do, thus in effect increasing the amount of memory accessible to an MPI process. Using this model, they modified a UPC-only BH code and got a twofold speedup at the expense of a 2% increase in code size. However, this hybrid scheme's performance comes from replicating the *entire* octree in each UPC group (in other words, from reduced remote data references). They did not discuss optimizations we found crucial to BH's performance. It is not clear whether this code is scalable.

Dekate et al. [8] described a BH implementation in ParalleX [15]. They suggest four main characteristics of a scalable and high-performance n-body simulation: data-driven computation, dynamic load balancing, data locality, and variable workload. To this end, they have implemented BH in HPX, a C++ implementation of ParalleX, making use of many light-weight threads to increase parallelism; work-queues to balance the load on the processors; interaction lists to improve data locality; and futures to make use of asynchronous operations. They also make use of manager threads and communication threads for the force-calculation phase. Some of these ideas are similar to what has been shown in this paper. However, they did evaluations only on shared-memory machines.

Jo et al. [14] described a point blocking optimization for traversal code, which can be thought as a counterpart of the classic loop tiling transformation for irregular applications. They introduced a transformation framework to automatically detect such optimization opportunities and autotuning techniques to determine appropriate parameters for the transformation. Our body tiling optimization was inspired by this work. However, we also measured BH locality in different body distributions and studied the interplay between locality in private caches and reuse in shared caches, which they did not mention.

Exploiting multithreading on multicores in high-performance computing has been extensively studied in different contexts. The linear algebra library PLASMA [3], a multicore version of LAPACK [4], has dynamic scheduling as one of its crucial elements [18]. It uses pthreads as its thread library and supports two scheduling strategies: static scheduling and dynamic scheduling. The dynamic scheduling strategy of PLASMA, implemented as QUARK, makes use of queues of tasks, from which worker threads pop and execute tasks. DAGuE [7], which extends PLASMA to distributed memory, has another commonality with our work, as it uses a separate thread, called the Asynchronous Communication Engine, for doing internode communication. However, PLASMA and DAGuE have been designed for dense linear algebra, where the communication pattern is known up-front and is regular. Neither is true for BH.

## 8. Summary

We have shown how one-sided communication, message-driven task scheduling, and caching of remote data can be combined to implement the Barnes-Hut algorithm with superior performance. We believe that a library implementing this programming model will prove useful for other applications as well, in terms of both ease of programming and performance. We plan to pursue this direction in future work. We have also shown the complex interplay, in multicore systems, between task scheduling and cache hit rate, at different levels of the cache hierarchy. The proper choice of a scheduling policy is extremely dependent on machine parameters and input parameters; on the other hand, the code behavior does not change rapidly across iterations. Hence, it is likely that runtime autotuning could be used to properly select the scheduling parameters.

## Acknowledgments

This work was supported by the U.S. Department of Energy, Office of Science, Advanced Scientific Computing Research, under Contract DE-AC02-06CH11357, and by the U.S. Department of Energy Sandia National Lab grant 1205852.

## References

- [1] Berkeley UPC. <http://upc.lbl.gov>.

- [2] SJ Aarseth, M. Henon, and R. Wielen. A comparison of numerical methods for the study of star cluster dynamics. *Astronomy and Astrophysics*, 37:183–187, 1974.
- [3] E. Agullo, J. Demmel, J. Dongarra, B. Hadri, J. Kurzak, J. Langou, H. Ltaieff, P. Luszczek, and S. Tomov. Numerical linear algebra on emerging architectures: The plasma and magma projects. volume 180, page 012037. IOP Publishing, 2009.
- [4] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen. *LAPACK Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, third edition, 1999.
- [5] J. Barnes and P. Hut. A hierarchical O(nlogn) force-calculation algorithm. *nature*, 324:4, 1986.
- [6] R. Barriuso and A. Knies. SHMEM user's guide for C. Technical report, Cray Research Inc, 1994.
- [7] G. Bosilca, A. Bouteiller, A. Danalis, T. Herault, P. Lemarinier, and J. Dongarra. DAGuE: A generic distributed dag engine for high performance computing. *Parallel Computing*, 2011.
- [8] Chirag Dekate, Matthew Anderson, Maciej Brodowicz, Hartmut Kaiser, Bryce Adelstein-Lelbach, and Thomas Sterling. Improving the scalability of parallel n-body applications with an event-driven constraint-based execution model. *Int. J. High Perform. Comput. Appl.*, 26(3):319–332, August 2012.
- [9] James Dinan, Pavan Balaji, Ewing Lusk, P. Sadayappan, and Rajeev Thakur. Hybrid parallel programming with MPI and Unified Parallel C. In *Proceedings of the 7th ACM international conference on Computing Frontiers*, CF '10, pages 177–186, New York, NY, USA, 2010. ACM.
- [10] Truong Vinh Truong Duy, Katsuhiko Yamazaki, Kosai Ikegami, and Shigeru Oyanagi. Hybrid MPI-OpenMP paradigm on SMP clusters: MPEG-2 encoder and n-body simulation. *CoRR*, abs/1211.2292, 2012.
- [11] Mary Hall, Richard Lethin, Keshav Pingali, Dan Quinlan, Vivek Sarkar, John Shalf, Robert Lucas, Katherine Yelick, Pedro C Diniz, Alice Koniges, et al. ASCR programming challenges for exascale computing. 2011.
- [12] M. Herlihy and N. Shavit. *The art of multiprocessor programming*. Morgan Kaufmann, 2008.
- [13] International Organization for Standardization. Programming Languages – Fortran. ISO/IEC 1539-1:2010 Standard, 2010.
- [14] Youngjoon Jo and Milind Kulkarni. Enhancing locality for recursive traversals of recursive structures. In *Proceedings of the 2011 ACM international conference on Object oriented programming systems languages and applications*, OOPSLA '11, pages 463–482, New York, NY, USA, 2011. ACM.
- [15] Hartmut Kaiser, Maciek Brodowicz, and Thomas Sterling. Parallel an advanced parallel execution model for scaling-impaired applications. In *Proceedings of the 2009 International Conference on Parallel Processing Workshops*, ICPPW '09, pages 394–401, Washington, DC, USA, 2009. IEEE Computer Society.
- [16] Laxmikant Kale, Anshu Arya, Abhinav Bhatele, Abhishek Gupta, Nikhil Jain, Pritish Jetley, Jonathan Lifflander, Phil Miller, Yanhua Sun, Ramprasad Venkataraman, Lukasz Wesolowski, and Gengbin Zheng. Charm++ for productivity and performance: A submission to the 2011 HPC Class II challenge. Technical Report 11-49, Parallel Programming Laboratory, November 2011.
- [17] Laxmikant V. Kale and Sanjeev Krishnan. Charm++: A portable concurrent object oriented system based on c++. In *Proceedings of the eighth annual conference on Object-Oriented Programming Systems, Languages, and Applications*, OOPSLA '93, pages 91–108, New York, NY, USA, 1993. ACM.
- [18] J. Kurzak, P. Luszczek, A. YarKhan, M. Faverge, J. Langou, H. Bouwmeester, and J. Dongarra. *Multi- and Many-Core Technologies: Programming, Algorithms, & Applications*, chapter Multithreading in the PLASMA Library. Taylor & Francis, 2011.
- [19] P.J. Mucci, S. Browne, C. Deane, and G. Ho. PAPI: A portable interface to hardware performance counters. In *Proceedings of the Department of Defense HPCMP Users Group Conference*. Citeseer, 1999.
- [20] J. Nieplocha, B. Palmer, V. Tipparaju, M. Krishnan, H. Trease, and E. Aprà. Advances, applications and performance of the global arrays shared memory programming toolkit. *International Journal of High Performance Computing Applications*, 20(2):203–231, 2006.
- [21] OpenMP Architecture Review Board. OpenMP application program interface version 3.1, 2011.
- [22] H. Rein and S.-F. Liu. REBOUND: An open-source multi-purpose n-body code for collisional dynamics. *Astronomy and Astrophysics*, 537:128, 2012.
- [23] J. Reinders. *Intel threading building blocks: Outfitting C++ for multi-core processor parallelism*. O'Reilly Media, Incorporated, 2007.
- [24] John K. Salmon. *Parallel hierarchical N-body methods*. PhD thesis, California Institute of Technology, 1991.
- [25] H. Shan, B. Austin, N.J. Wright, E. Strohmaier, J. Shalf, and K. Yelick. Accelerating applications at scale using one-sided communication. In *The 6th Conference on Partitioned Global Address Space Programming Models*, 1993.
- [26] R. Stevens, A. White, et al. Architectures and technology for extreme scale computing. In *ASCR Scientific Grand Challenges Workshop Series, Tech. Rep.*, 2009.
- [27] UPC Consortium. UPC language specifications, v1.2. Technical Report LBNL-59208, Lawrence Berkeley National Lab, 2005.
- [28] M.S. Warren and J.K. Salmon. A parallel hashed oct-tree n-body algorithm. In *Proceedings of the 1993 ACM/IEEE conference on Supercomputing*, pages 12–21. ACM, 1993.
- [29] Mathias Winkel, Robert Speck, Helge Hübner, Lukas Arnold, Rolf Krause, and Paul Gibbon. A massively parallel, multi-disciplinary Barnes-Hut tree code for extreme-scale N-body simulations. *Computer Physics Communications*, 183:880–889, 2012.
- [30] S.C. Woo, M. Ohara, E. Torrie, J.P. Singh, and A. Gupta. Methodological considerations and characterization of the SPLASH-2 parallel application suite. In *Proceedings of the 22nd Annual International Symposium on Computer Architecture*, pages 24–36, 1995.
- [31] J. Zhang, B. Behzad, and M. Snir. Optimizing the Barnes-Hut algorithm in UPC. In *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis(SC)*, page 75. ACM, 2011.

The submitted manuscript has been created by UChicago Argonne, LLC, Operator of Argonne National Laboratory ("Argonne"). Argonne, a U.S. Department of Energy Office of Science laboratory, is operated under Contract No. DE-AC02-06CH11357. The U.S. Government retains for itself, and others acting on its behalf, a paid-up nonexclusive, irrevocable worldwide license in said article to reproduce, prepare derivative works, distribute copies to the public, and perform publicly and display publicly, by or on behalf of the Government.