# An OpenSHMEM Implementation for the Adapteva Epiphany Coprocessor

James Ross[1] and David Richie[2]

[1] U.S. Army Research Laboratory, Aberdeen Proving Ground, MD 21005, USA
`james.a.ross176.civ@mail.mil`
[2] Brown Deer Technology, Forest Hill, MD 21050, USA
`drichie@browndeertechnology.com`

**Abstract.** This paper reports the implementation and performance evaluation of the OpenSHMEM 1.3 specification for the Adapteva Epiphany architecture within the Parallella single-board computer. The Epiphany architecture exhibits massive many-core scalability with a physically compact 2D array of RISC CPU cores and a fast network-on-chip (NoC). While fully capable of MPMD execution, the physical topology and memory-mapped capabilities of the core and network translate well to Partitioned Global Address Space (PGAS) programming models and SPMD execution with SHMEM.

**Keywords:** OpenSHMEM, Network-on-Chip (NoC), Single-Board Computer, Performance Evaluation

## 1 Introduction and Motivation

The OpenSHMEM communications library is designed for computer platforms using Partitioned Global Address Space (PGAS) programming models [1]. Historically, these were large Cray supercomputers, but now the OpenSHMEM interface may also be used on commodity clusters. The Adapteva Epiphany architecture represents a divergence in computer architectures typically used with OpenSHMEM and is just one of many emerging parallel architectures that present a challenge in identifying effective programming models to exploit them. While some researchers may be considering how the OpenSHMEM API may interact with coprocessors, the work presented here leverages the API for device-level operation. In some aspects, the Epiphany architecture resembles a symmetric multiprocessing (SMP) multi-core processor with a shared off-chip global memory pool. However, each core can directly address the private address space of neighboring cores across an on-chip 2D mesh network. Thus, the architecture also has the characteristics of a PGAS platform. Previous proof-of-concepts demonstrated that message passing protocols could achieve good application performance on the Epiphany architecture [2], [3]. However, it was unclear if the OpenSHMEM 1.3 standard could be fully implemented within the platform limitations and achieve high performance using a standard programming model without resorting to non-standard software extensions.

Existing open source OpenSHMEM implementations are inadequate within the constraints of the Epiphany architecture, so a new C language implementation named *ARL OpenSHMEM for Epiphany* was developed from scratch. The design emphasizes a reduced memory footprint, high performance, and simplicity, which are often competing goals. This paper discusses the Epiphany architecture in Sect. 2.1, the OpenSHMEM implementation and performance evaluation in Sect. 3, and a discussion of future work and potential standard extensions for embedded architectures in Sect. 4.

## 2   Background

The 16-core Epiphany-III coprocessor is included within the $99 ARM-based single-board computer and perhaps represents the low-cost end of programmable hardware suitable for SHMEM research and education. Many universities, students, and researchers have purchased the platform with over 10,000 sales to date. Despite this, programming the platform and achieving high performance or efficiency remain challenging for many users. Like GPUs, the Xeon Phi, and other coprocessors, typical applications comprise host code and device code. Only a minimal set of communication primitives exist within the non-standard Epiphany Hardware Utility Library (eLib) for multi-core barriers, locks, and data transfers [4]. The barrier and data transfer routines are not optimized for low latency. Other primitives within eLib use unconventional 2D row and column indexing, which cannot easily address arbitrary numbers of working cores or disabled cores. More complicated collectives, such as those in the OpenSHMEM specification, are left as an exercise for the application developer.

Although not discussed in detail in this paper, the CO-PRocessing Threads (COPRTHR) 2.0 SDK [5] further simplifies the execution model to the point where the host code is significantly simplified, supplemental, and even not required depending on the use case [6]. There are essentially two modes of possible execution. The first mode requires host code with explicit Epiphany coprocessor offload routines. The second mode uses a host-executable coprocessor program with the conventional main routine provided. The program automatically performs the coprocessor offload without host code. Combined with the work presented in this paper, the COPRTHR 2.0 SDK enables many OpenSHMEM applications to execute on the Epiphany coprocessor without any source code changes. Execution occurs as if the Epiphany coprocessor is the main processor driving computation. COPRTHR 1.6 was used to present the Threaded MPI model for Epiphany [2] as well as a number of applications [7], [8].

### 2.1   Epiphany Architecture

Many modern computer architectures address the "memory wall problem" by including increasingly complex cache hierarchies and core complexity, wider memory buses, memory stacking, and complex packaging to maintain the SMP hardware and software architecture. The Epiphany architecture unwinds decades of

these types of changes – it is a cache-less, 2D array of RISC cores with a fast network-on-chip (NoC) that can an be simply described as a "cluster on a chip". Each core within the Epiphany-III architecture contains 32 KB of SRAM which is shared between instructions and local data. The Epiphany architecture can scale to one megabyte of SRAM per core, but there is a linear design trade-off between the number of cores and available memory for a fixed die space. The core local memory is memory-mapped, and each core may directly access the local memory of any core within the mesh network. Each core has shared memory access to off-chip global DRAM, although this access is significantly slower than local memory or non-uniform memory access (NUMA) to neighboring core memory. The highest performance and most energy-efficient applications leverage inter-core communication and on-chip data reuse. Like many high performance computing (HPC) clusters, the inter-core communication is generally explicit in order to achieve highest performance. The architecture is also scalable by tiling multiple chips without additional "glue logic". The tight coupling between the core logic and the on-chip mesh network enables very low-latency operation of OpenSHMEM routines. An architectural overview appears in Figure 1. Unlike most application programming interfaces for communication, there is no additional software layer to handle networking for hardware abstraction. As we will discuss in further detail, the OpenSHMEM implementation for Epiphany performs network operations directly.
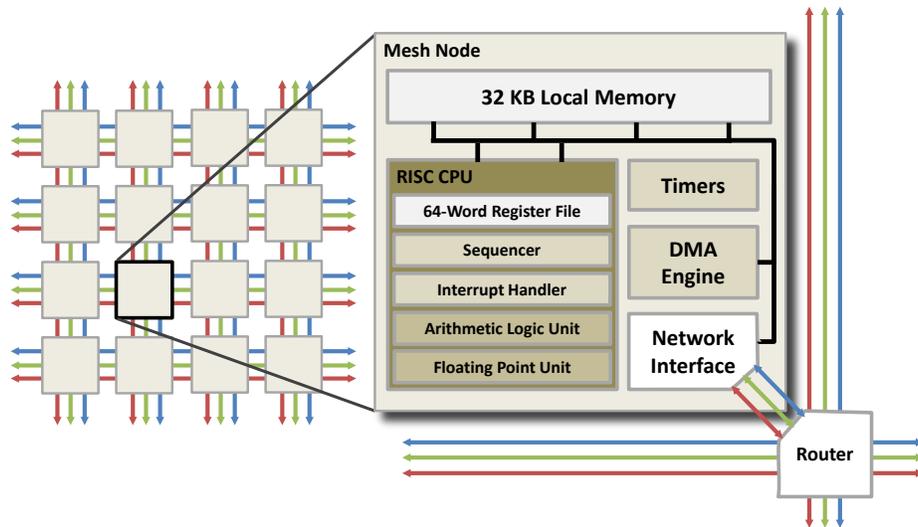


**Fig. 1.** The 16-core Epiphany-III architecture is a 2D array of RISC CPU cores. It contains a 64-word register file, sequencer, interrupt handler, integer and floating point units, timers, and DMA engines for the fast network-on-chip

## 3 Implementation and Performance Evaluation

Due to the tight memory constraints of the Epiphany memory and availability of specialized hardware features, the OpenSHMEM reference implementation built on GASNet was not suitable for deployment on the Epiphany cores. As a credit to the OpenSHMEM specification and the Adapteva Epiphany architecture documentation, the full OpenSHMEM 1.3 implementation was written and optimized over a period of a few weeks. The entire library, including the optional extensions described in detail later, is approximately 1800 lines of code and does not require additional software. The software directly targets the underlying hardware features and was designed to be extremely lightweight in order to compile to small binaries expected with embedded architectures.

Linear scaling algorithms were avoided, and many of the collective routines use dissemination or recursive doubling algorithms, optimized for low-latency on the Epiphany network. The remote memory access routines, `shmem_TYPE_put` and `shmem_TYPE_get`, use hand-tuned memory-mapped load and store primitives with a hardware loop feature specific to the Epiphany architecture. The non-blocking remote memory access routines use the dual-channel Direct Memory Access (DMA) engine on each processor network node. The distributed locking and atomic routines leverage an atomic `TESTSET` instruction that performs an atomic "test-if-not-zero" and conditional write. An optional hardware barrier implementation was also developed for a specialized `shmem_barrier_all` for extremely low-latency global barriers. An optional inter-processor interrupt and corresponding interrupt service routine (ISR) enable faster `shmem_TYPE_get` operations by interrupting the remote core to use the optimized `shmem_TYPE_put`.

Many of the OpenSHMEM routines have some component that is hardware accelerated on the Epiphany architecture such as zero-overhead hardware loops for copying data, memory-mapped loads and stores, the `TESTSET` instruction for remote locks and atomics, a wait on AND (`WAND`) instruction for a low-latency `shmem_barrier_all`. The `MULTICAST` experimental feature would enable energy-efficient, low-latency broadcasts but is presently unused. The point-to-point synchronization routines are among the simplest to implement and do not have a section dedicated to discussion. Generally, they spin-wait on local values until they meet the criteria defined by the routine. The memory ordering routines need only verify that both DMA engines have an idle status by spin-waiting on the relevant special register. There are no intermediate data copies in this implementation.

The performance evaluation of the Epiphany OpenSHMEM implementation began with the OpenSHMEM micro-benchmark codes. The timing code had to be modified because the `gettimeofday` routine is only accurate to a microsecond, and many of the operations operate in the sub-microsecond regime.

Many of the communication routines in the performance evaluation include the parameters $\alpha$ and $\beta^{-1}$ in the figure subtitle along with their standard deviations. These two parameters are from the "$\alpha$-$\beta$ model" for communication in HPC. They neatly summarize the communication time ($T_c$) to include the la-

tency ($\alpha$) and marginal cost ($\beta$) to transfer a message (of size $L$) in equation 1. The $\beta^{-1}$ parameter is the peak effective core bandwidth for the routine.

$$T_c = \alpha + \beta \cdot L \qquad (1)$$

## 3.1   Library Setup, Exit, Query Routines

The `shmem_init` routine retrieves or calculates the local processing element (PE) number (for `shmem_my_pe`) and number of PEs (for `shmem_n_pes`), configures the optimized hardware barrier or collective dissemination barrier arrays, obtains the SHMEM heap memory offset, and precalculates a few other addresses for improved runtime performance. The `shmem_ptr` routine can directly calculate remote memory locations using simple logical shift and bitwise operations.

## 3.2   Memory Management Routines

Memory management on the Epiphany processor is atypical. Each Epiphany-III core has a flat 32 KB local memory map from address `0x0000` to `0x7fff`. Programs are typically loaded starting at `0x0100` if extremely constrained for memory, or `0x0400` if using the COPRTHR 2 interface. The stack pointer typically moves downward from the high address. Data used for the application, including the SHMEM data heap, begins directly after the program space. Figure 2 shows the typical memory layout of an Epiphany-III core using the COPRTHR 2 interface as it relates to the PGAS model. The static or global variables that are typically defined within the application appear below the free local memory address within the symmetric heap. They are still symmetrical across all Epiphany cores as the program binary is identical.

Due to the tight memory constraints, a more modern memory allocator was not addressed in this work. The basic memory management system calls `brk` and `sbrk` are more suited for controlling the amount of memory allocated from the SHMEM data heap for each process element because there is no virtual address abstraction. Instead, there is a local base memory tracking pointer that stores the current free memory base address and incremented with each allocation. The memory management routines build on these calls, but care must be taken to adhere to the following rules:

1. `shmem_free` must be called in the reverse order of allocation if making subsequent allocations
2. `shmem_realloc` can only be used on the last (re)allocated pointer
3. `shmem_align` alignment must be a power of 2 greater than 8 (default is 8)

This is a pragmatic approach that we feel is reasonable and won't even be noticed on most codes. Calling `shmem_free` moves the local base memory tracking pointer to the address in the function argument so most routines only need to call it once for the first allocated buffer in a series if freeing all memory. The `shmem_realloc` routine could be designed to copy the contents of the old
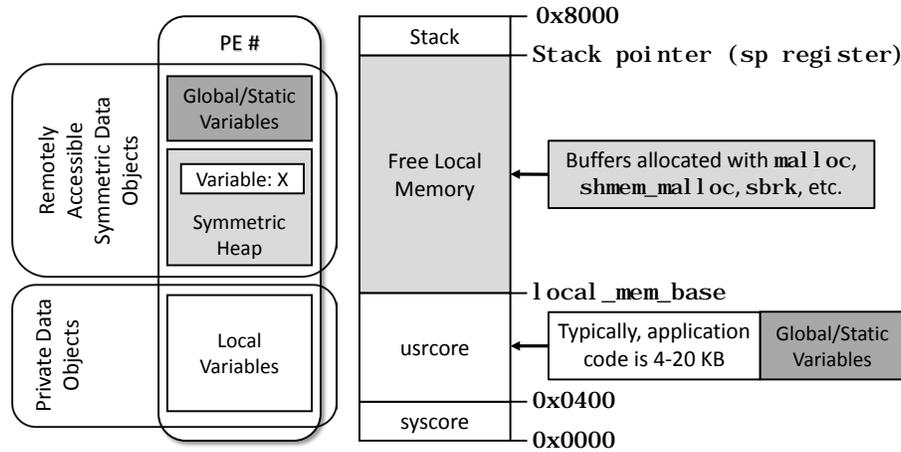
PE #

Remotely Accessible Symmetric Data Objects

Global/Static Variables

Variable: X

Symmetric Heap

Private Data Objects

Local Variables

Stack

0x8000

Stack pointer (sp register)

Free Local Memory

Buffers allocated with `malloc`, `shmem_malloc`, `sbrk`, etc.

local_mem_base

usrcore

Typically, application code is 4-20 KB

Global/Static Variables

0x0400

syscore

0x0000

**Fig. 2.** The PGAS memory model (left) and the equivalent typical memory layout on an Epiphany-III core (right)

buffer to the new buffer, however, this would waste the memory space in the original allocation (a precious commodity on the Epiphany architecture). Future developments with COPRTHR 2 may address these deficiencies by exporting the COPRTHR host-side memory management to the coprocessor threads.

### 3.3 Remote Memory Access Routines

Inter-process memory copying on the Epiphany is trivial, and a simple loop over incrementing source and destination arrays can be done in C code. However, like many optimized `memcpy` routines, high-performance copies are non-trivial. A high-performance inter-processor memory copy routine does not appear to be in the eLib library. So after quite some time of hand-tuning in assembly, a put-optimized method was written that makes use of a "zero-overhead" hardware loop and four-way unrolled staggered double-word loads and (remote) stores. A specialization for the edge case of unaligned memory is also included since the Epiphany architecture requires loads and stores to be memory aligned to the data size. Assuming the fast path is taken, the core can transfer a double-word (8 bytes) per clock cycle. However, since the 8 byte load operation requires an additional cycle, the effective peak network copy is 8 bytes every two clocks. For a clock rate of 600 MHz, peak contiguous network transfers may achieve up to 2.4 GB/sec. Having the NoC and core clocks pinned ensures that application communication performance scales with the chip clock speed. The same put-optimized memory copy subroutine is used for get operations. This is sub-optimal, but remote read operations will never be as high-performance as remote write operations on the Epiphany architecture, so they should generally be avoided. Remote direct read operations are slower than equivalent remote direct

write operations because the read request must first traverse the network to the receiving core network interface, then the data must traverse the network back to the requesting core. Unlike a remote direct write operation which can issue store instructions without a response, the read operation stalls the requesting core until the load instruction returns data to a register. Issuing multiple requests does little to mitigate this performance issue, thus, the throughput of the optimized `shmem_put` is approximately an order of magnitude greater than `shmem_get` as shown in Figure 3.



**Fig. 3.** Performance of optimized `shmem_put` (top left) and `shmem_get` (top right) for contiguous data exchange operations for 16 processing elements, speedup comparison with eLib (bottom left), and experimental inter-processor user interrupt for high-performance `shmem_get` (bottom right).

In order to address this performance disparity with contiguous remote reads, an inter-processor interrupt is configured and signaled by the receiving core, causing an equivalent fast write to be executed. The receiving core is then signaled to continue upon completion of the inter-processor ISR. This is an experimental feature because it uses the user interrupt and must be enabled with `SHMEM_USE_IPI_GET` during compilation. It has the greatest performance impact for large transfers. The method has a turnover point for buffers larger than 64

bytes so that smaller transfers are read directly and larger transfers use the inter-processor interrupt. All results for contiguous block transfers and a performance comparison with the equivalent eLib interface in the eSDK are shown in Figure 3.

### 3.4 Non-blocking Remote Memory Access Routines

The set of non-blocking remote memory access routines (`shmem_put_nbi` and `shmem_get_nbi`) makes use of the on-chip DMA engine. The DMA engine has two independent DMA channels per processor node so that two non-blocking transfers may execute concurrently. Each channel has a separate DMA specification of the source and destination address configuration. The configuration is capable of 2D DMA operations with flexible stride sizes. This could support an extension to the OpenSHMEM 1.3 standard for non-blocking strided remote memory access routines if needed. The performance results for the non-blocking remote memory access routines appear in Figure 4.
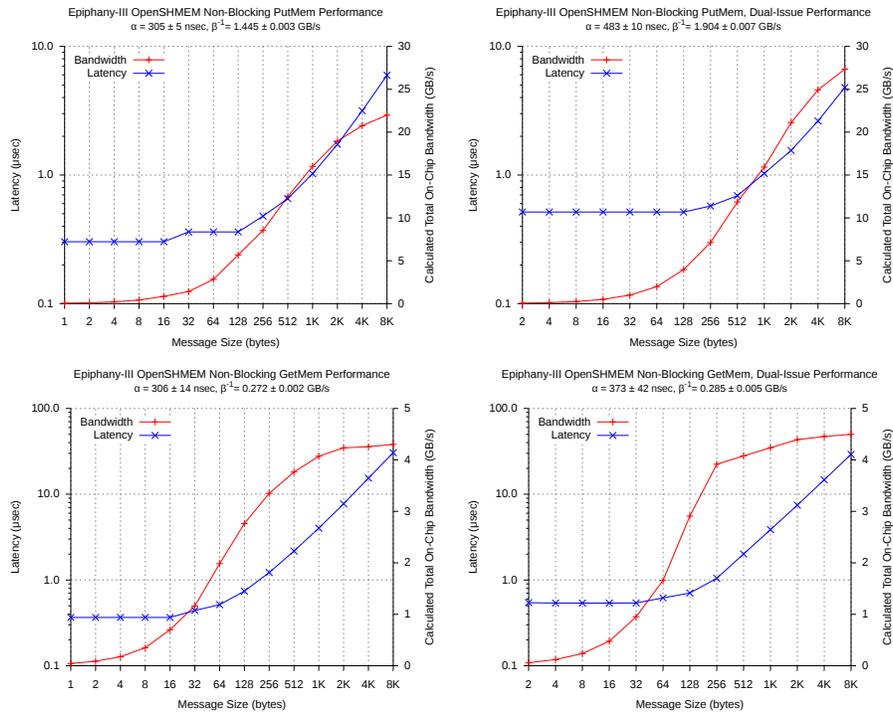


**Fig. 4.** Performance of non-blocking remote memory access routines.

Application performance improvement may be realized for large non-blocking transfers by splitting transfers into two portions and calling two non-blocking

transfers, however, the performance benefit is marginal and often worse. Due to hardware errata in the Epiphany-III, the DMA engine is throttled to less than half of the peak bandwidth of 8 bytes per clock, or 4.8 GB/sec [9]. If fully enabled, as expected in future chips, the DMA engine would be used for the blocking remote memory access routines rather than remote load/store instructions. In general, it may be faster to use blocking transfers because the DMA engine setup overhead is relatively high, and there are often bank conflicts with the concurrent computation and DMA engine access, hindering fully overlapped communication and computation. The blocking operation, `shmem_quiet`, spin-waits on the DMA status register. Alternatively, a DMA ISR could be used to continue the `shmem_quiet` operation, but it is not clear how this could be higher performance.

### 3.5 Atomic Memory Operations

The Epiphany-III ISA does not have support for atomic instructions, but the `TESTSET` instruction used for remote locks may be used to define other atomic operations in software. With the current code design, it is trivial to extend to other atomic operations with a single line of code if additional atomic operations are defined by the OpenSHMEM specification in the future. At a core level, memory access for both fetch and set operations completes in a single clock cycle and is therefore implicitly atomic. The fetch operation still must traverse the network to the remote core and return the result. Each data type specialization uses a different lock on the remote core as per the specification. The performance results for the 32-bit integer atomic routines appear in Figure 5.
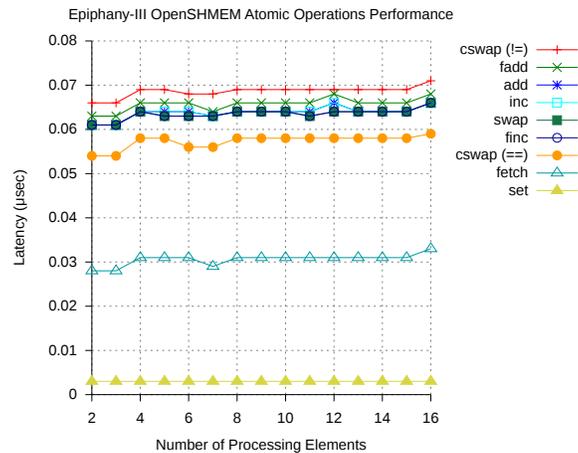


**Fig. 5.** Performance of OpenSHMEM atomic operations for 32-bit integers and a variable number of processing elements. Atomic operations are performed in a tight loop on the next neighboring processing element.

### 3.6 Collective Routines

Multi-core barriers are critical to performance for many parallel applications. The Epiphany-III includes hardware support for a fully collective barrier with the WAND instruction and corresponding ISR. This hardware support is included as an experimental feature within the OpenSHMEM library and must be enabled by specifying SHMEM_USE_WAND_BARRIER at compile time. After several implementations of barrier algorithms, it was determined that a dissemination barrier was the highest-performing software barrier method. It is not clear if this algorithm will continue to achieve the highest performance on chip designs with a larger number of cores; alternative tree algorithms may be needed. The eLib interface in the eSDK uses a counter-based collective barrier and requires a linearly increasing amount of memory with the number of cores. The dissemination barrier requires $8 \cdot log_2(N)$ bytes of memory, where $N$ is the number of processing elements within the barrier. The use of this synchronization array mitigates the need for signaling by locks at each stage of the barrier. The collective eLib barrier completes in 2.0 $\mu$sec while the WAND barrier completes in 0.1 $\mu$sec. The performance for group barriers for a subset of the total processing elements is shown in Figure 6. The latency of the dissemination barrier increases logarithmically with the number of cores so that more than eight cores take approximately 0.23 $\mu$sec.
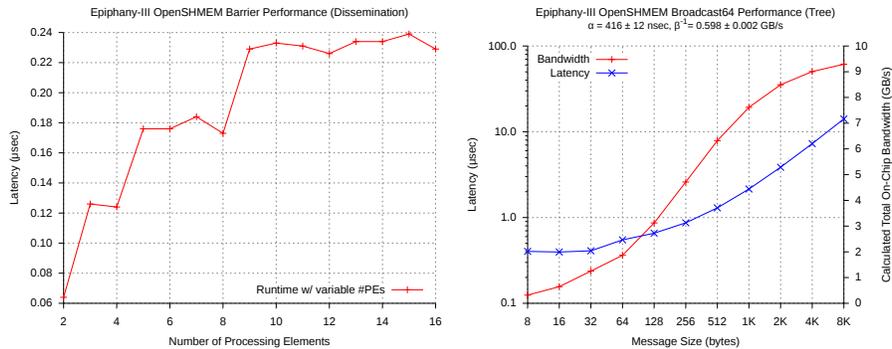


**Fig. 6.** Performance of shmem_barrier for variable number of processing elements (left) and the performance of shmem_broadcast64 for variable message sizes (right)

Broadcasts are important in the context of the Epiphany application development in order to limit the replication of off-chip memory accesses to common memory. It is faster to retrieve off-chip data once and disseminate it to other processing elements in an algorithmic manner than for each processing element to fetch the same off-chip data. The data are distributed with a logical network tree, moving the data the farthest distance first in order to prevent subsequent stages increasing on-chip network congestion. The broadcast routines use the same

high-performance memory copying subroutine as the contiguous data transfers. Effective core bandwidth approaches the theoretical peak performance for this algorithm and is approximately $2.4/log_2(N)$ GB/sec. Figure 6 shows collective broadcast performance for variable message sizes.

The `shmem_collect` and `shmem_fcollect` routines use ring and recursive doubling algorithms for concatenating blocks of data from multiple processing elements. Each uses the optimized contiguous memory copying routine. There is likely room for improvement with these routines; the measured performance appears in Figure 7.
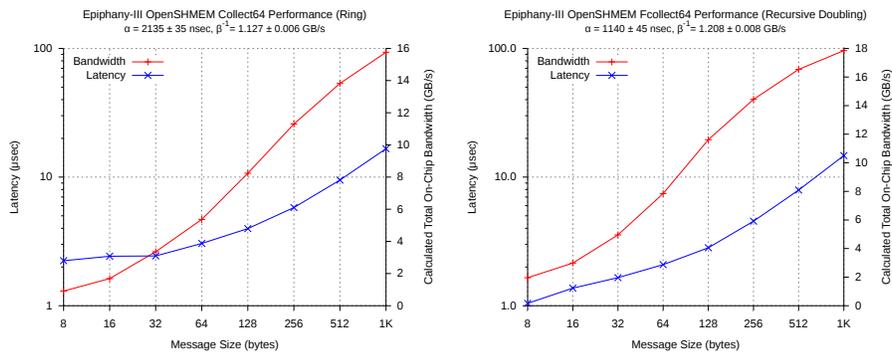


**Fig. 7.** Performance of linear scaling `shmem_collect64` and recursive doubling `shmem_fcollect64` for variable message sizes on 16 processing elements

The `shmem_TYPE_OP_to_all` reduction routines are important for many multi-core applications. The routines use different algorithms depending on the number of processing elements. A ring algorithm is used for processing elements that number in non-powers of two and a dissemination algorithm for powers of two. The symmetric work array is used for temporary storage and the symmetric synchronization array is used for multi-core locks and signaling. The performance of `shmem_int_sum_to_all` appears in Figure 8. Other routines vary marginally in performance due to data types and the arithmetic operation used. Reductions that fit within the symmetric work array have improved latency as seen in the figure.

The performance of the contiguous all-to-all data exchange, `shmem_alltoall`, appears in Figure 9. This routine has a relatively high overhead latency compared to other collectives.

### 3.7 Distributed Locking Routines

The distributed locking routines, `shmem_set_lock` and `shmem_test_lock`, are easily supported by the atomic `TESTSET` instruction. The actual lock address is defined in the implementation to be on the first processing element. These
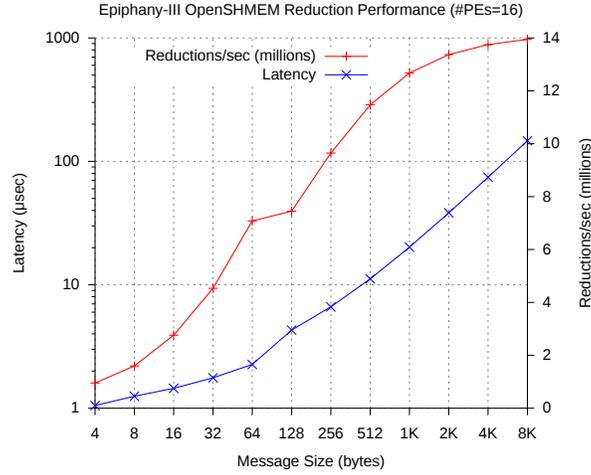
**Fig. 8.** Reduction performance for `shmem_int_sum_to_all` for all 16 processing elements. The latency and the number of collective reductions per second are shown. The effect of the minimum symmetric work array size for reductions, defined as `SHMEM_REDUCE_MIN_WRKDATA_SIZE` per the OpenSHMEM specification, is apparent for small reductions
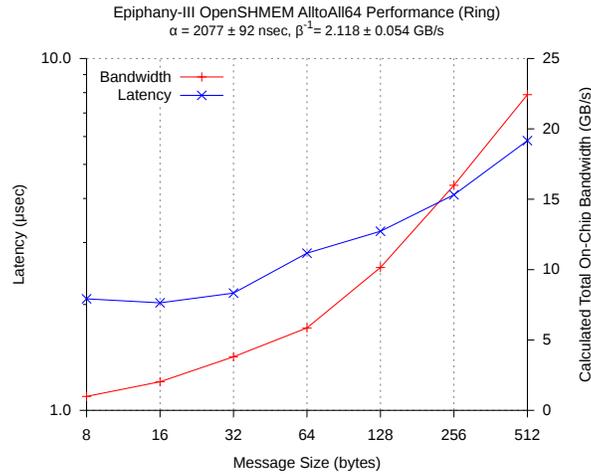


**Fig. 9.** Performance of the new (to version 1.3) contiguous all-to-all data exchange operation, `shmem_alltoall`, for 16 processing elements

locking mechanisms are also the basis for the atomic operations detailed in Sect. 3.5 but for multiple processing elements. The `shmem_clear_lock` routine is a simple remote write to free the lock. Although this scheme works well for the 16 processing elements on the Epiphany-III, the performance bottleneck will likely be a problem scaling to much larger core counts. Application developers should avoid using these global locks.

## 4    Future Work and Discussion of Extensions for Embedded Architectures

It is our intention to release *ARL OpenSHMEM for Epiphany*, as well as the performance evaluation codes and benchmarks used in this paper, as open source software through the U.S. Army Research Laboratory's GitHub account [10] for Parallella community input and further development. The Epiphany architecture may also be updated in the future to add more hardware support for many of the existing OpenSHMEM routines. Many of the currently proposed OpenSHMEM extensions and updates should be addressable. A non-blocking strided remote memory access routine could be supported with the existing DMA engine as mentioned in Sect. 3.4. Some other extensions do not make sense for the architecture. For example, Epiphany is not a multithreaded architecture and, although it can be performed via software, is not the ideal mechanism for improving performance. The OpenSHMEM standard should remain sufficiently lightweight to address low-level operations without relying on specific architectural features.

One of the more challenging portions of the OpenSHMEM standard for the Epiphany architecture and other embedded architectures are the memory management routines. It makes some sense for some platforms to have a pre-allocated symmetric heap from which memory allocations will be made. Within an Epiphany local core, there is no memory virtualization between the physical address and the memory address returned by the allocation routines as available memory is linearly removed from the symmetric heap. The limitations of the available core space make it challenging to introduce a Linux-like abstract model of virtual memory. As OpenSHMEM is a low-level interface and application developers are already accustomed to explicitly managing memory, it may make some sense to improve memory management interfaces, such as those discussed in Sect. 3.2, for embedded architectures.

## 5    Conclusion

OpenSHMEM provided an effective and pragmatic programming model for the Epiphany architecture. The header-only implementation enabled compiler optimizations for program size and application performance that is difficult to achieve using a standard pre-compiled library. We demonstrated improved performance and many useful features compared to the current eLib library despite the additional software abstraction with the OpenSHMEM interface. The *ARL*

*OpenSHMEM for Epiphany* demonstrated high-performance execution while approaching hardware theoretical networking limits, and low-latency operation for many of the OpenSHMEM routines.

## References

1. Barbara Chapman, Tony Curtis, Swaroop Pophale, Stephen Poole, Jeff Kuehn, Chuck Koelbel, and Lauren Smith. Introducing OpenSHMEM: SHMEM for the PGAS Community. In *Proceedings of the Fourth Conference on Partitioned Global Address Space Programming Model*, PGAS '10, pages 2:1–2:3, New York, NY, USA, 2010. ACM.
2. David Richie, James Ross, Song Park, and Dale Shires. Threaded MPI Programming Model for the Epiphany RISC Array Processor. *Journal of Computational Science*, 9:94 – 100, 2015. Computational Science at the Gates of Nature.
3. James Ross and David Richie. Implementing OpenSHMEM for the Adapteva Epiphany RISC Array Processor. *Procedia Computer Science*, 80:2353 – 2356, 2016. International Conference on Computational Science 2016, ICCS 2016, 6-8 June 2016, San Diego, California, USA.
4. GitHub - adapteva/epiphany-libs: Epiphany runtime libraries and utilities. `https://github.com/adapteva/epiphany-libs`. Accessed: 2016-05-24.
5. COPRTHR-2 Epiphany/Parallella Developer Resources. `http://www.browndeertechnology.com/resources_epiphany_developer_coprthr2.htm`. Accessed: 2016-07-01.
6. David Richie and James Ross. Advances in Run-Time Performance and Interoperability for the Adapteva Epiphany Coprocessor. *Procedia Computer Science*, 80:1531 – 1541, 2016. International Conference on Computational Science 2016, ICCS 2016, 6-8 June 2016, San Diego, California, USA.
7. James A. Ross, David A. Richie, and Song J. Park. Implementing Image Processing Algorithms for the Epiphany Many-Core Coprocessor with Threaded MPI. IEEE, September 2015.
8. James A. Ross, David A. Richie, Song J. Park, and Dale R. Shires. Parallel Programming Model for the Epiphany Many-Core Coprocessor Using Threaded MPI. *Microprocessors and Microsystems*, 2016.
9. Adapteva, Inc. *E16G301 Epiphany$^{TM}$ 16-Core Microprocessor Datasheet*, June 2013. Rev 14.03.11.
10. US Army Research Laboratory - GitHub. `https://github.com/USArmyResearchLab`. Accessed: 2016-05-24.