



An OpenSHMEM Implementation for the Adapteva Epiphany Coprocessor

James Ross

james.a.ross176.civ@mail.mil
US Army Research Laboratory

David Richie

drichie@browndeertechnology.com
Brown Deer Technology



U.S. ARMY
RDECOM

Outline



- **Epiphany Architecture**
- **Programming Challenge**
- **Hardware/Software Setup**
- **OpenSHMEM Interface**
- **Benchmark Performance**
- **Conclusions & Future Work**



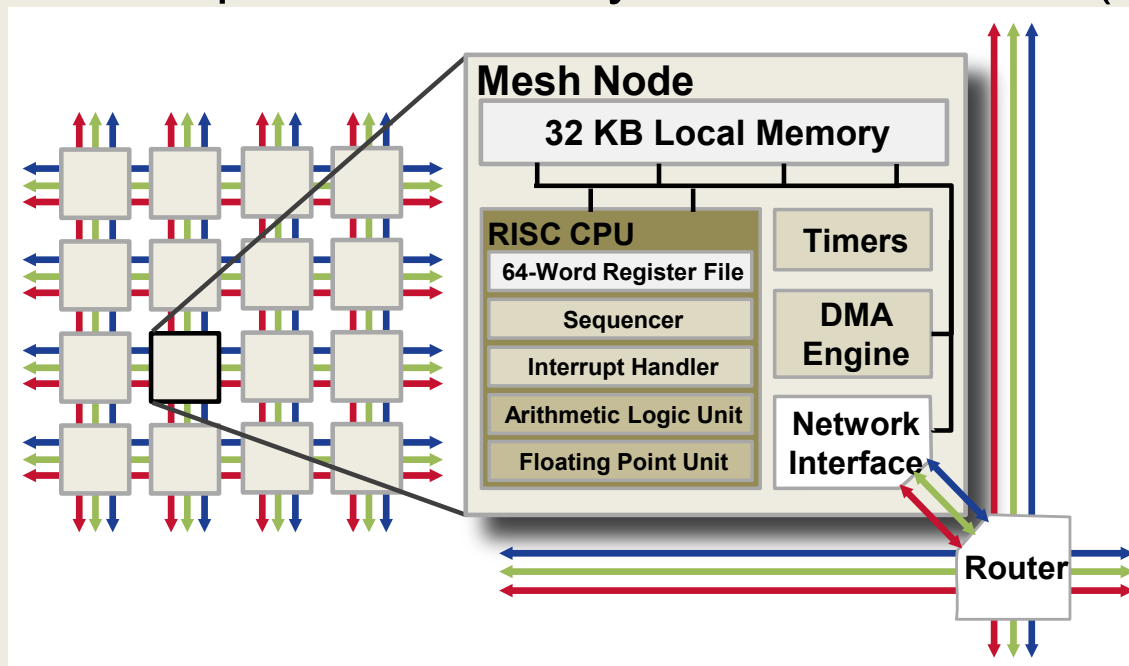
U.S. ARMY
RDECOM

UNCLASSIFIED

Adapteva Epiphany Architecture



- Design emphasizes simplicity, scalability, power-efficiency
- 2D array of RISC cores, 2D Network on Chip (NoC)
- 512 KB shared global scratch memory (32 KB/core, Epiphany-III)
- Fully divergent cores
- Minimal un-core functionality, e.g., no data or instruction cache
- Existing design scales to thousands of cores
- High performance/power efficiency ~50 GFLOPS/W (Epiphany-IV)





U.S. ARMY
RDECOM

UNCLASSIFIED

Hardware/Software Setup

ARL



Hardware

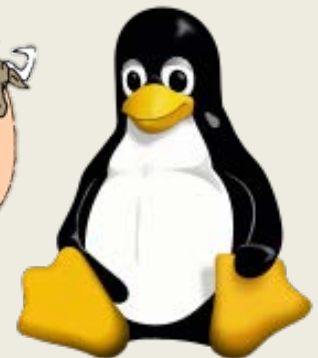
- Parallella development board (\$99 'Micro-Server')
- Dual-core ARM host processor
- 16-core Epiphany-III co-processor (@ 600 MHz)
 - 19.2 GFLOPS
 - 76.8 GB/s bandwidth

Software at time of publication

- OpenSHMEM 1.3 (our implementation)
- Parallella Linux image (2015.1)
 - Epiphany SDK 2015.1
 - GCC 4.8
- Brown Deer Technology
 - COPRTHR-2 Beta



Open
SHMEM



Brown Deer
Technology

adapteva

U.S. ARMY
RDECOM

Programming Challenge



Programming challenge:

- Can the OpenSHMEM API be efficiently used with embedded PGAS architectures like Epiphany?
- Distributed memory mapped cores with 32KB local memory per core
- Best viewed as a “distributed cluster on chip”
- Non-uniform memory access (NUMA) to mapped local memory
- Typical many-core programming models (OpenMP, OpenCL, CUDA) leverage common cache / memory hierarchy for SMP
- Here, there is no hardware data/instruction cache

Two key observations:

- Architecture resembles cluster on chip
- Inter-core data movement is key to performance

Proposition: Use SHMEM as device-level programming model

We are using the OpenSHMEM API to enable the device-level parallelism of a 2D RISC array processor



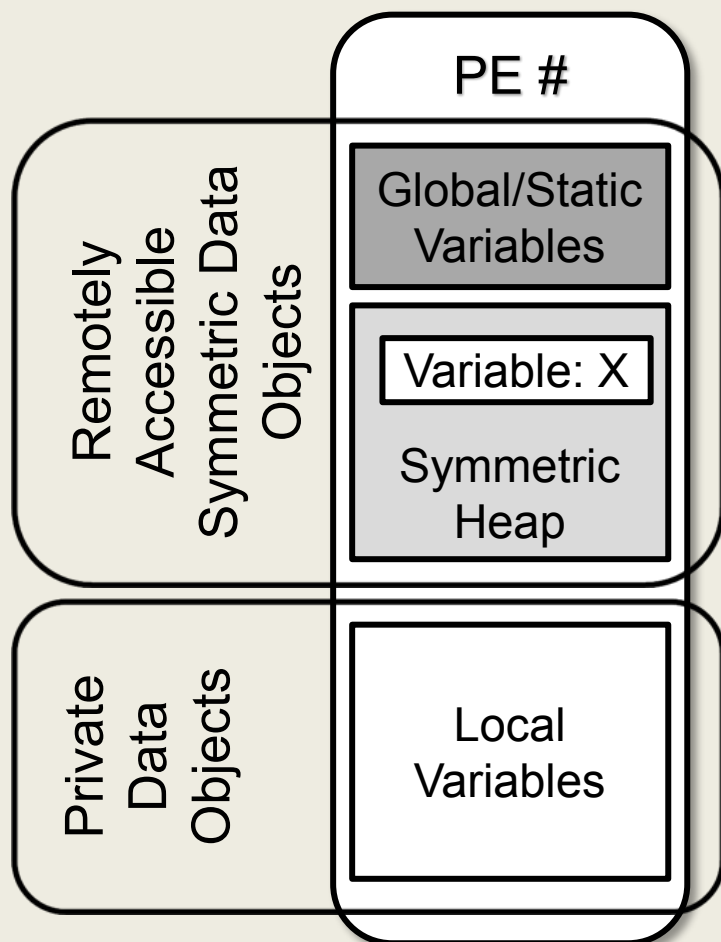
U.S. ARMY
RDECOM

UNCLASSIFIED

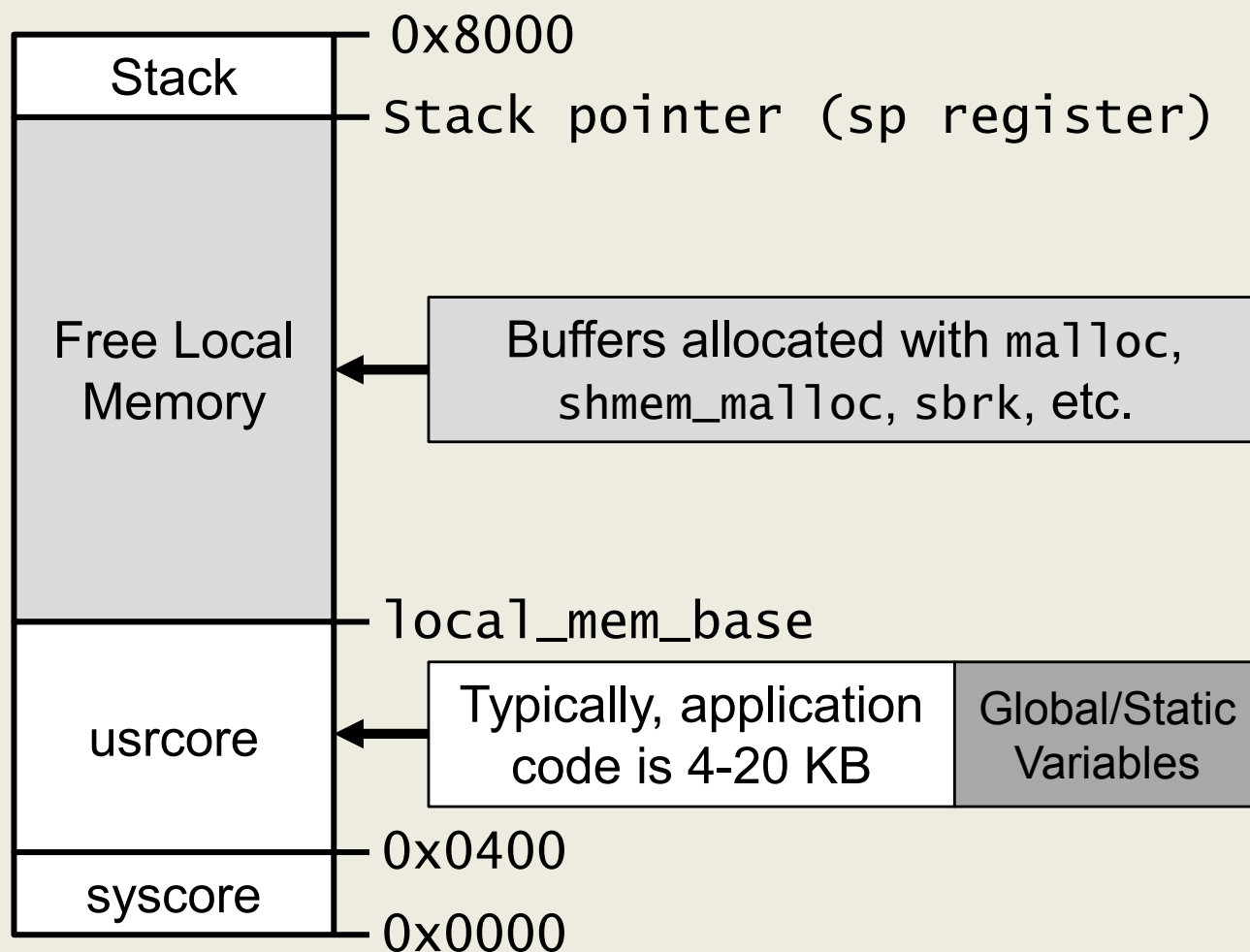
OpenSHMEM



PGAS Memory Model



Epiphany-III Core Memory Layout (Typical)



U.S. ARMY
RDECOM**OpenSHMEM****ARL**

Why OpenSHMEM for PGAS?

- Unified Parallel C, Split-C, Fortress, Coarray Fortran, Chapel, X10 require building/modifying a compiler instead of a library

So why not MPI?

- MPI substantially larger library with higher-level abstraction and two-sided communication means more code (source and binary)
- OpenSHMEM has improved data referencing semantics and reduced interface complexity
 - No message tags, no status, no special types, no I/O
 - MPI makes no assumptions for symmetric memory allocation. Requires correct remote address calculation

Why not eSDK (e-lib)? (See Supplemental Slides)

- Non-standard, less portable code
- No abstraction between physical row/column and virtual process
- Memory management, collectives, reductions, atomics missing

U.S. ARMY
RDECOM

OpenSHMEM Design



Objectives:

- Maintain/conform to OpenSHMEM API
- Complete implementation of API
- Design must lead to efficient implementation for architecture
 - Tree-based algorithms, high performance, small code size

Challenges:

- Significant local memory constraint (32KB per core)
 - Typical OpenSHMEM uses significantly larger amounts of memory
 - Entire library is 200+ routines, typical applications use just several

Boons:

- Extremely low latency on-chip inter-core communication
- No middleware (transparent network interface for memory access)
- Simple DMA engine for asynchronous communication



U.S. ARMY
RDECOM

Implementation



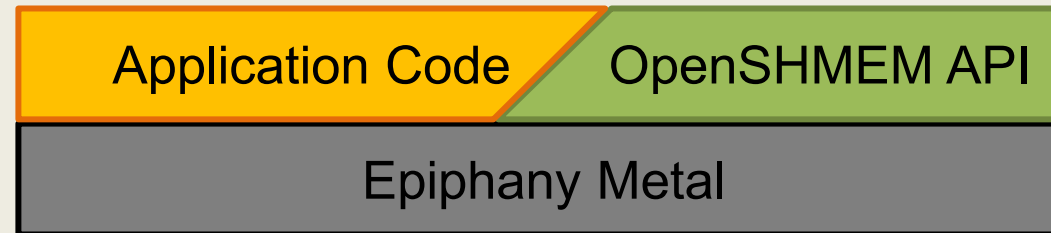
- Currently implemented in header-only library (~1800 LOC C/inline asm)
 - Compiler optimizes for constant arguments, reduces application code size
 - Data type differentiation handled by macros
- No software stack layers for low-level communication primitives
 - Typical implementations based on additional software stacks (GASNet, network interfaces, other software layers)
- Optimized for low latency performance, efficient hierarchical communication (dissemination barrier, recursive-doubling, etc)

U.S. ARMY
RDECOM

Implementation



OpenSHMEM / Application / Epiphany software stack:



- OpenSHMEM library replaces device-side eSDK (e-lib) code.
- No consideration for coprocessor offloading (e-hal, libcoprthr)
- Message passing between cores literally corresponds to load or store instruction

...not this:





U.S. ARMY
RDECOM

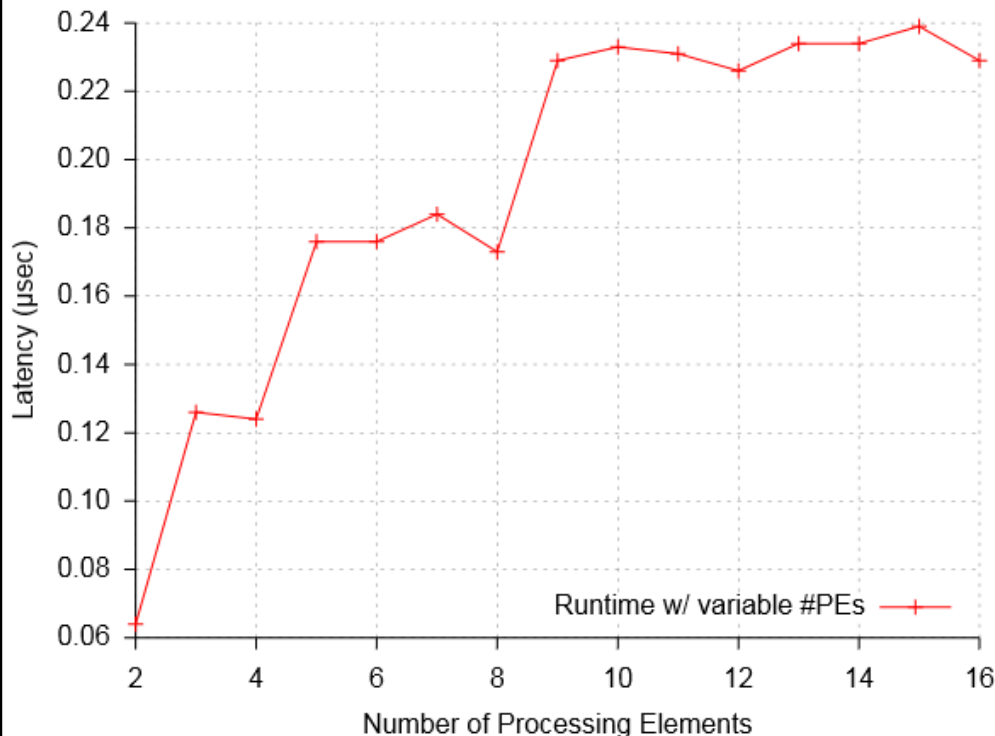
UNCLASSIFIED

Performance Results



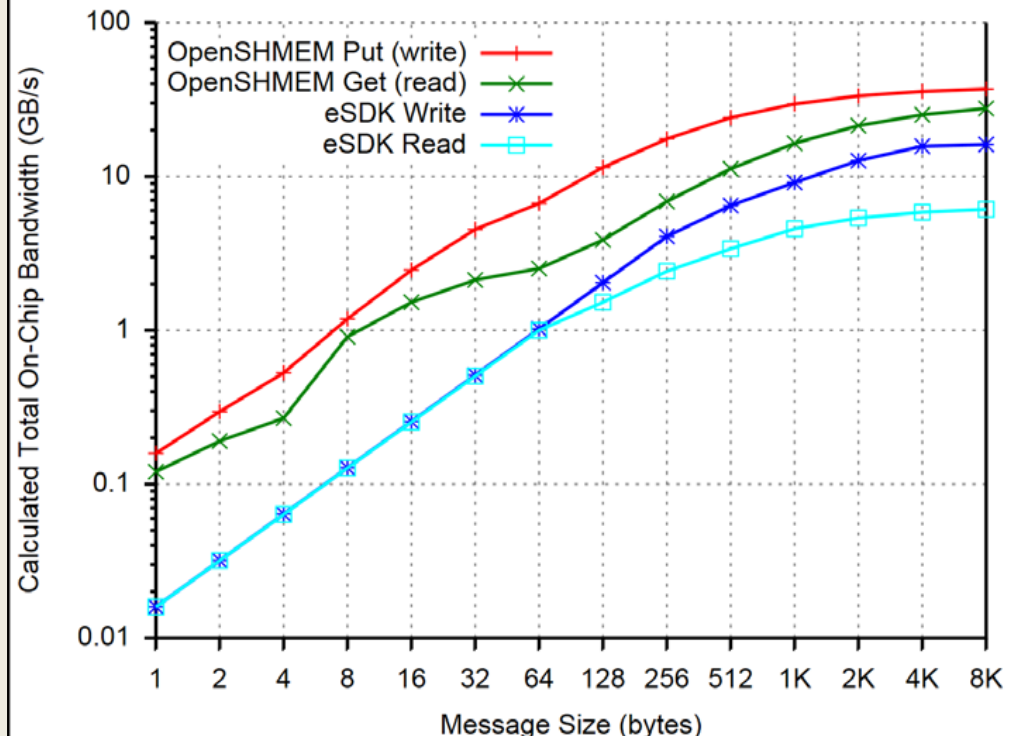
- Few directly comparable routines between eSDK (e-lib) and OpenSHMEM

Epiphany-III OpenSHMEM Barrier Performance (Dissemination)



- >9.1x speedup** for software barrier
- 20x speedup** for fixed 16 core case (hardware WAND barrier = 0.1 µsec)
- eSDK barrier = **2.0 µsec**

Memory Copying Performance (#PEs=16)



- 2.1-9.9x speedup** for all message sizes
- Peak Put bandwidth = **2.4 GB/s** per core
- Corresponds to alternating ldrd/strd



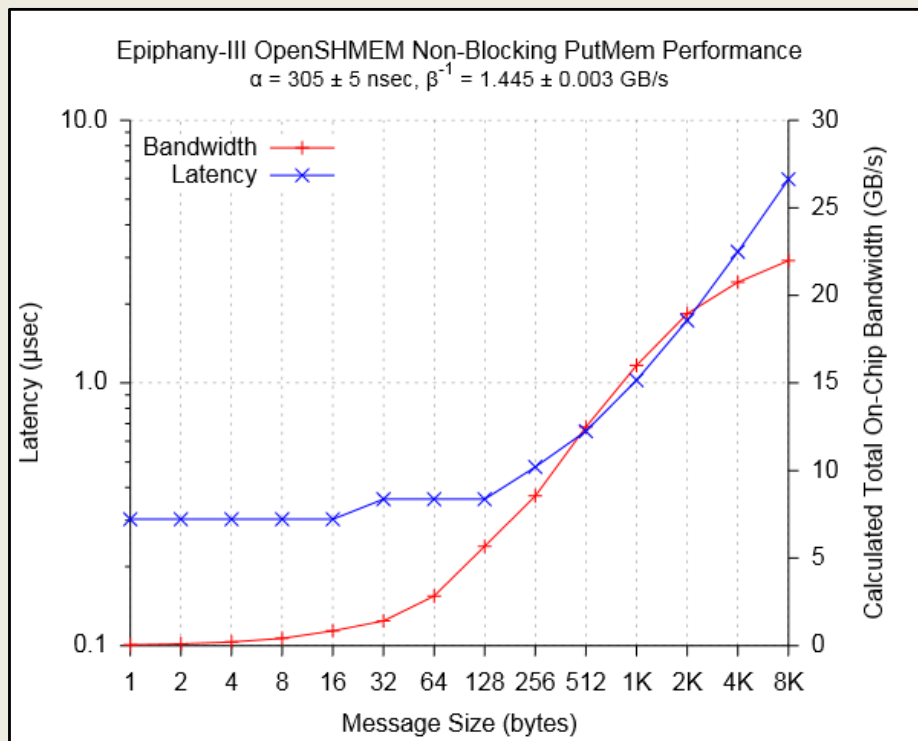
U.S. ARMY
RDECOM

UNCLASSIFIED

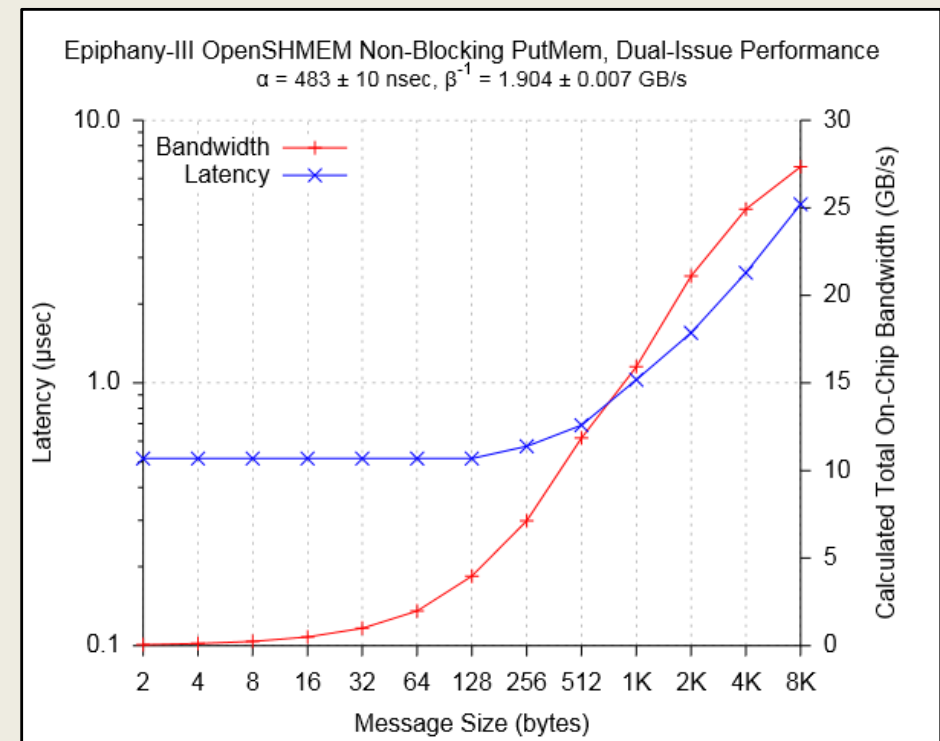
Performance Results



- Non-blocking DMA performance (automatic dual-issue DMA scheduling)
 - Universally underperforms blocking message passing
 - Dual-issue DMA increases latency for marginal gain for large transfers
 - Hardware errata on E3 prevents full DWORD/clock (4.8 GB/s)



- **1.45 GB/s** peak performance per core



- **1.91 GB/s** per core peak performance
- Code performs *shmem_putmem_nbi* twice with half total data



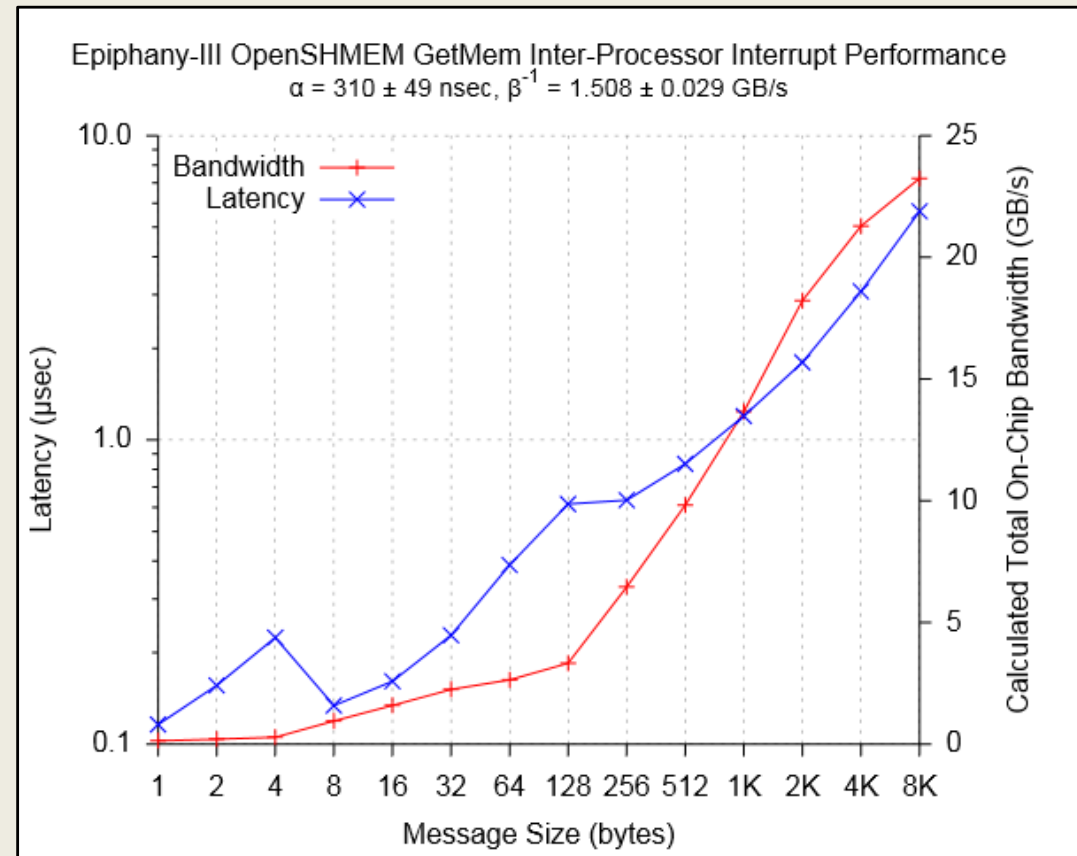
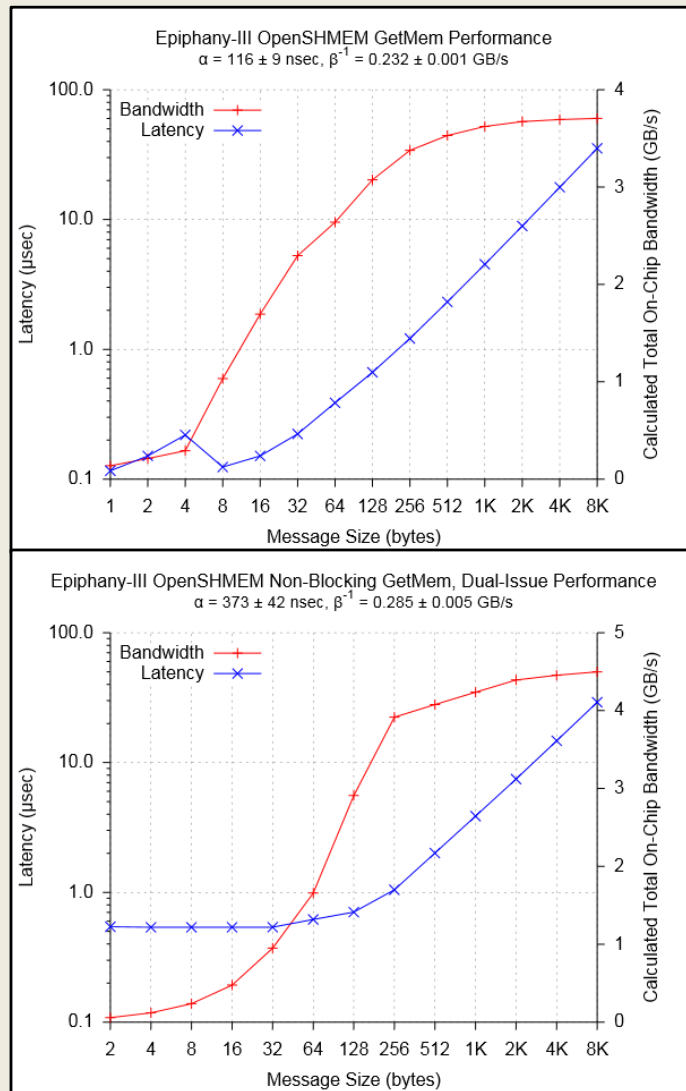
U.S. ARMY
RDECOM

UNCLASSIFIED

Performance Results



- Get performance worse than Put because read network requires response/stall
- Use Inter-Processor Interrupt (experimental feature) to cause remote Put



- **200-300 MB/s** per core peak Get performance
- **1.5 GB/s** peak performance with IPI
 - Crossover at $N \geq 64$ bytes



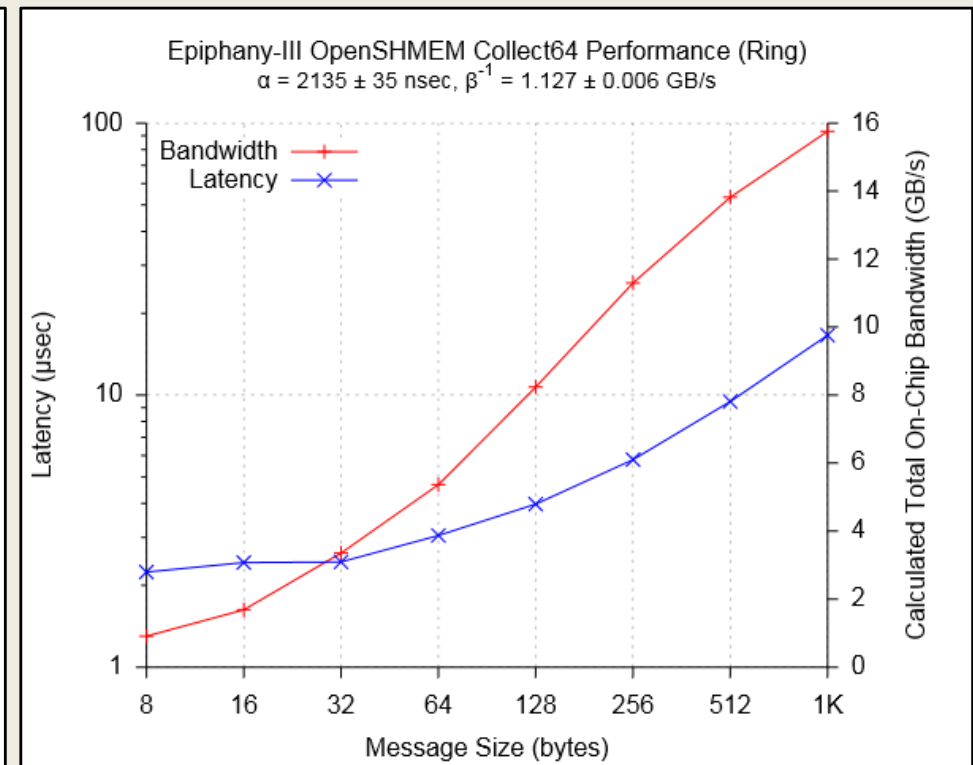
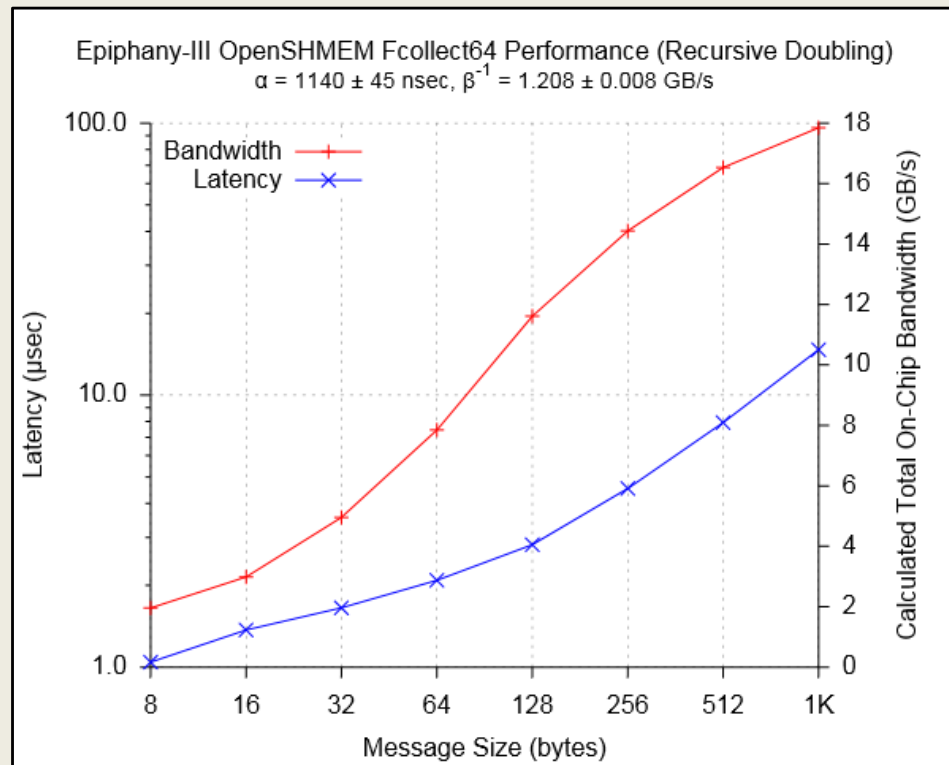
U.S. ARMY
RDECOM

UNCLASSIFIED

Performance Results



- Collect peak bandwidth good, but with relatively high latency for small sizes
 - Concatenates blocks of data from multiple PEs to all PEs



- **1.21 GB/s** peak performance per core
- **1.13 GB/s** peak performance per core



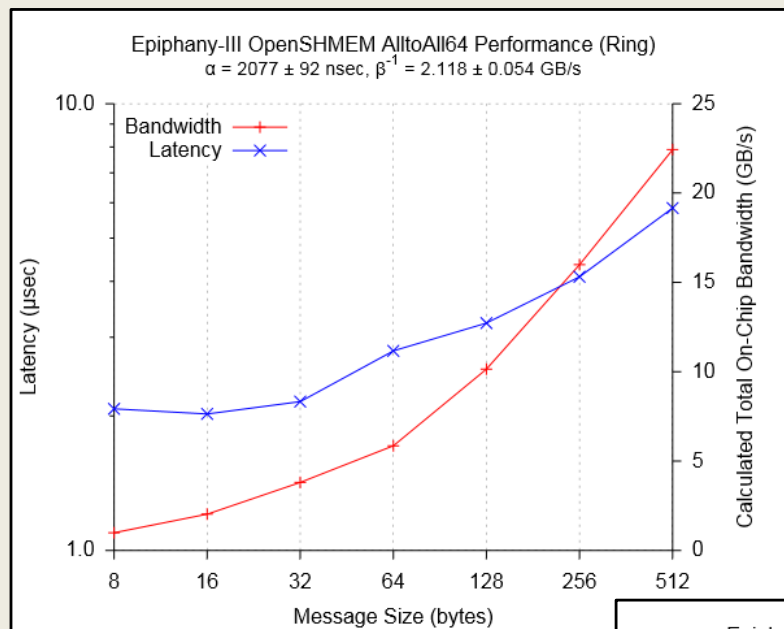
U.S. ARMY
RDECOM

UNCLASSIFIED

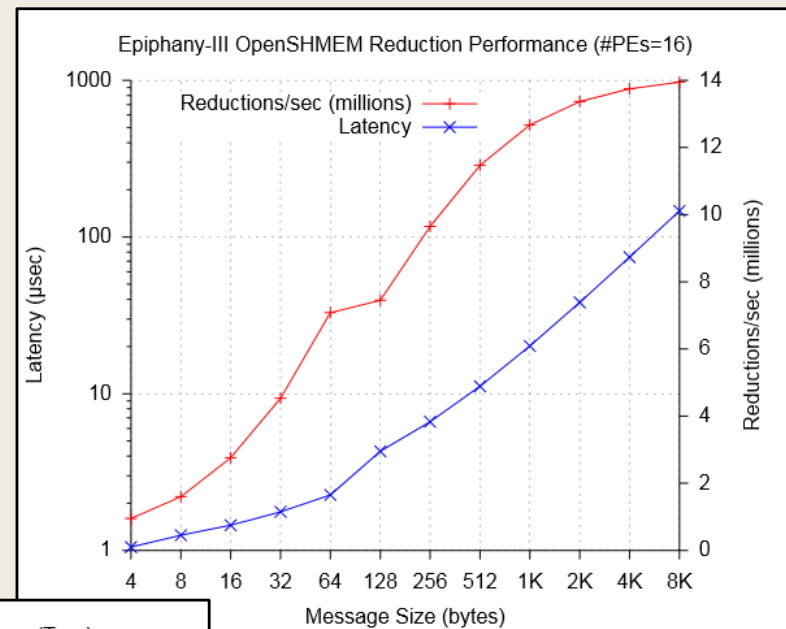
Performance Results



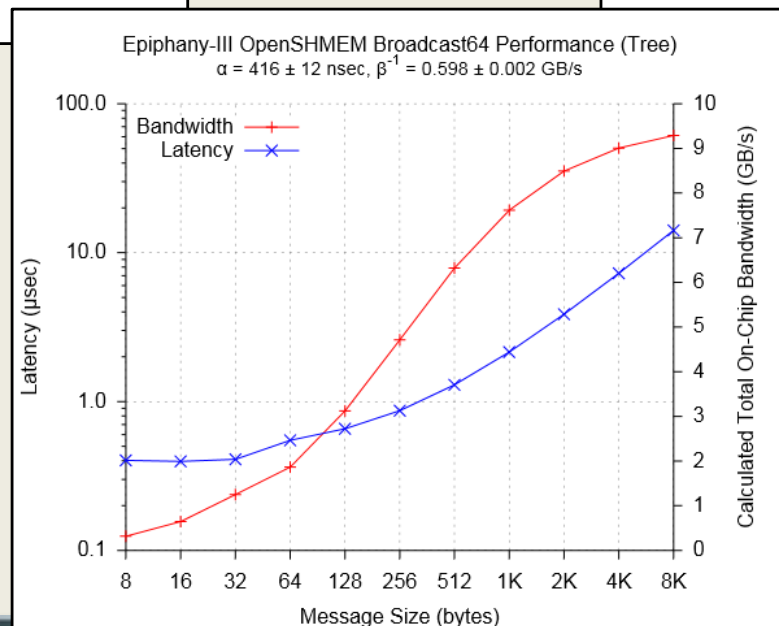
- Collective operations achieve high performance



AlltoAll Peak = **2.12 GB/s**



14 MegaReductions/s



Effective Broadcast = Peak
Put / $\log_2(\text{NPES}) = \mathbf{600}$
MB/s



Example

```

1  #include <shmem.h>
2  #define N 3040
3
4  long pSync[SHMEM_REDUCE_SYNC_SIZE] = { SHMEM_SYNC_VALUE };
5  int pwrk[SHMEM_REDUCE_MIN_WRKDATA_SIZE];
6  extern float dot_product(float* a, float* b, int n); // optimized dot product routine**
7  float c = 0.0f; // collective result
8
9  int main(void)
10 {
11     int nd8m1 = (N >> 3) - 1; // N/8 - 1
12     shmem_init();
13
14     int me = shmem_my_pe();
15     int npes = shmem_n_pes();
16
17     // allocation
18     float* a = (float*)shmem_malloc(N*sizeof(float));
19     float* b = (float*)shmem_malloc(N*sizeof(float));
20
21     // initialization...
22     for (int i = 0; i < N; i++) a[i] = b[i] = 1.0f;
23
24     shmem_barrier_all();
25
26     c = dot_product(a, b, nd8m1); // assembly optimized routine
27     shmem_float_sum_to_all(&c, &c, 1, 0, 0, npes, pwrk, pSync);
28
29     shmem_free(b);
30     shmem_free(a);
31
32     shmem_finalize();
33 }

```

| code | size (bytes) |
|-------------|--------------|
| syscore | 1024 |
| dot_product | 208 |
| main* | 1916 |

**Combined, this code achieves 87%
peak performance for 16 cores:
16.8 GFLOPS
67.3 GB/s
Nelements = 48,640 (16*3040)**

* Includes inlined *shmem_(init, my_pe, n_pes, malloc, barrier_all, float_sum_to_all, free, finalize)* routines

** Assembly optimized *dot_product* available in supplemental slides

**U.S. ARMY
RDECOM**

Conclusions & Future Work



- This work reported on full OpenSHMEM 1.3 implementation for embedded PGAS architecture, Adapteva's Epiphany
- OpenSHMEM provides an effective programming model for this class of architecture
- Header-only implementation enables compiler optimizations for program size and performance
- Performance exceeds eSDK (e-lib) in speed, code size
- Performance/latencies in another performance class than traditional implementations due to on-chip network – a favorable comparison
- Distribution for community input (on ARL github account)
- Improve memory management with COPRTHR-2
- Recommendations to OpenSHMEM standard committee for embedded PGAS architectures



U.S. ARMY
RDECOM

Resources



- US Army Research Laboratory GitHub
 - <https://github.com/USArmyResearchLab>
 - Intend to release ARL OpenSHMEM for Epiphany with example codes and benchmarks as free and open source software
- COPRTHR 2.0 Resources/Download
 - http://www.browndeertechnology.com/resources_epiphany_developer_coprthr2.htm
- Parallella (Epiphany-III host platform)
 - <http://parallella.org/>

Questions: james.a.ross176.civ@mail.mil



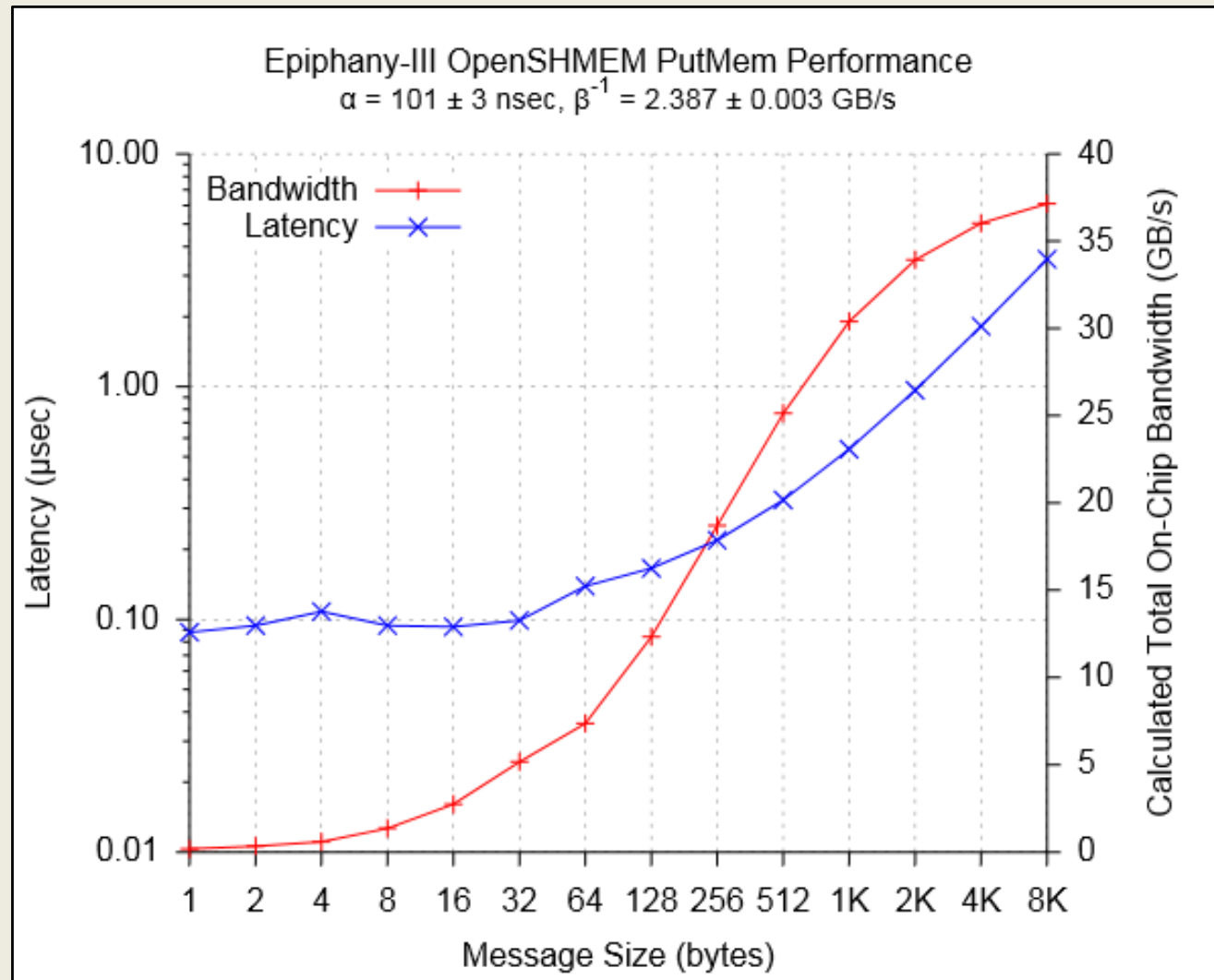
U.S. ARMY
RDECOM

ARL

Supplemental Slides

U.S. ARMY
RDECOM

More Figures

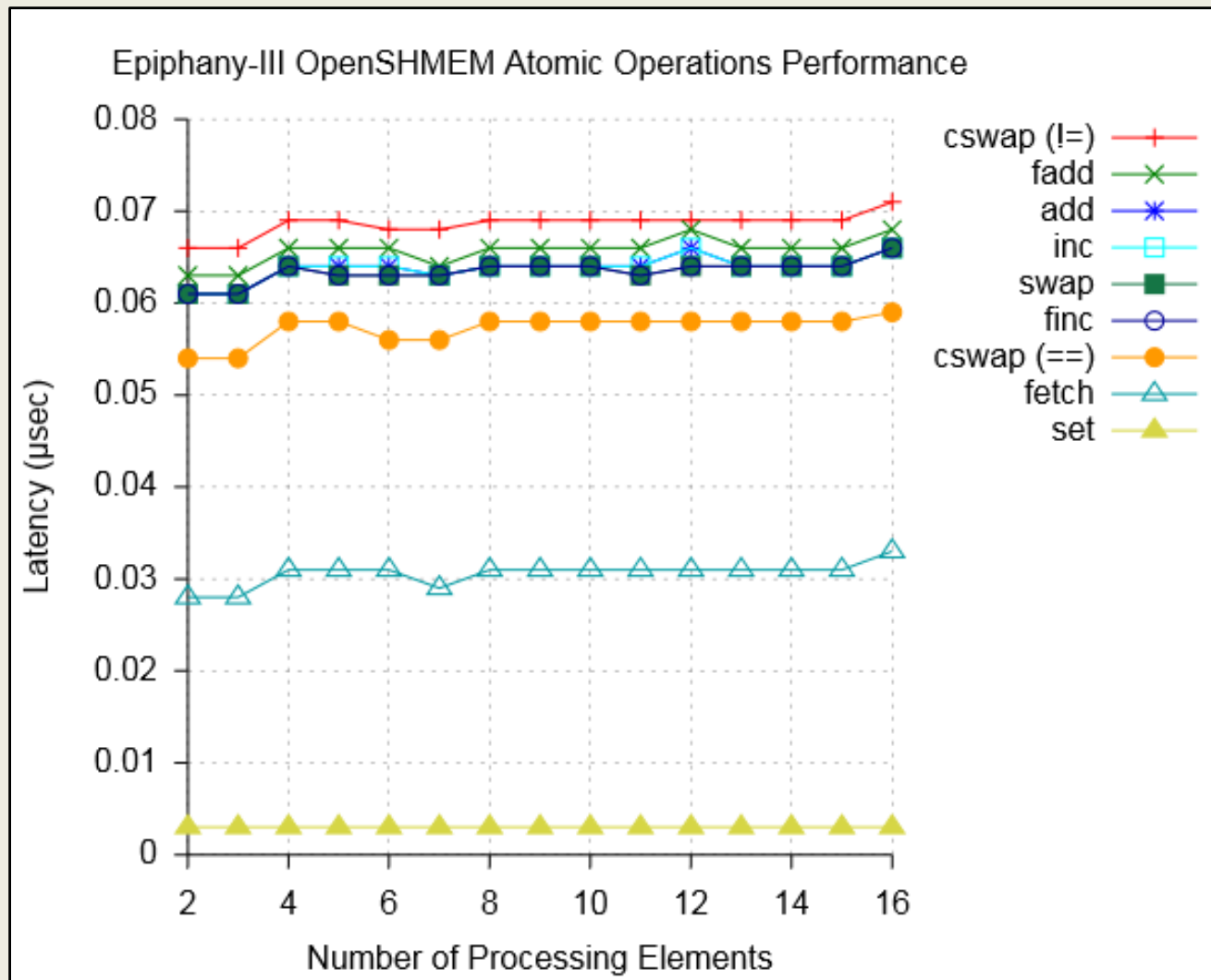
ARL



U.S. ARMY
RDECOM

More Figures

ARL





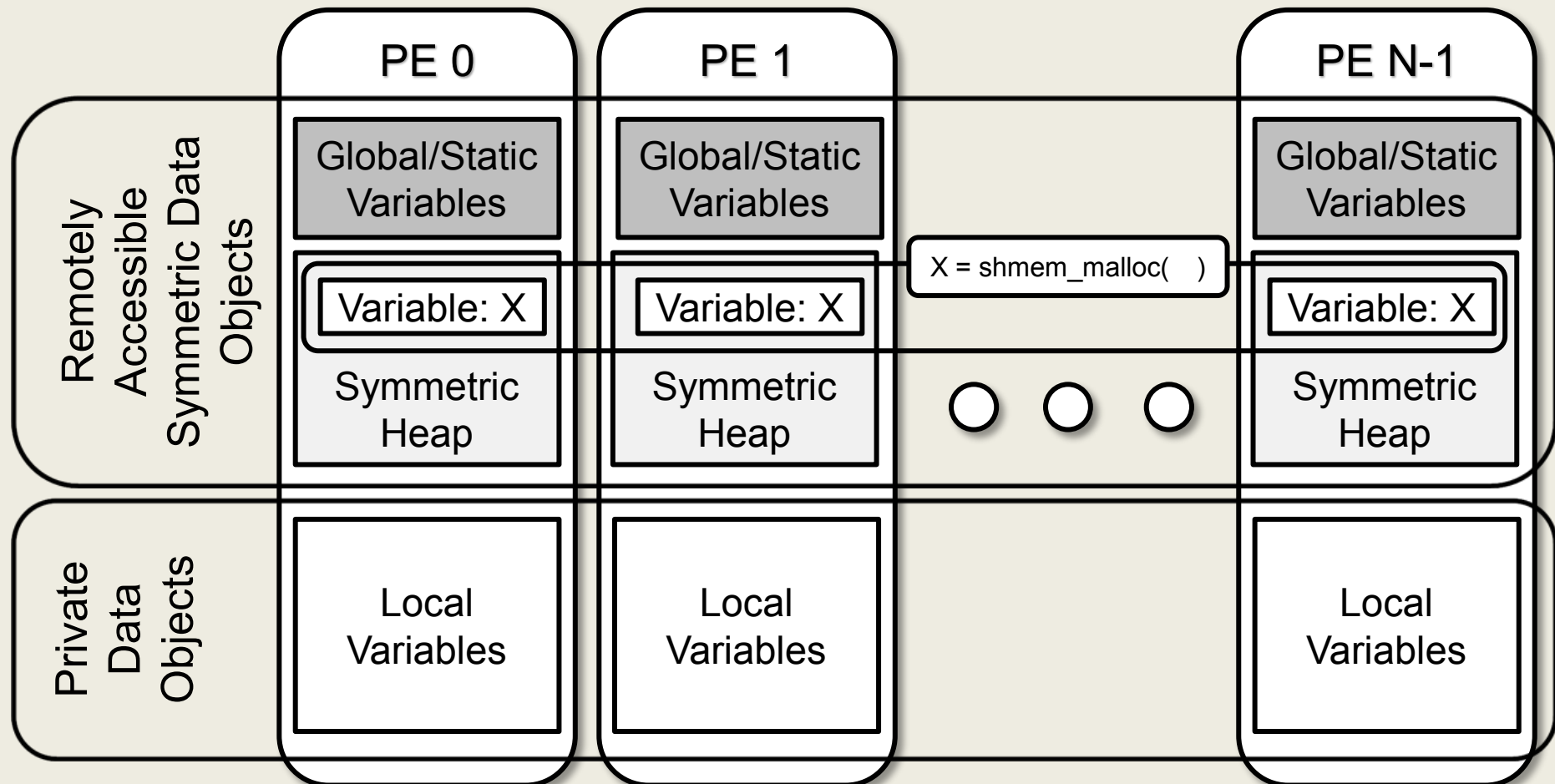
U.S. ARMY
RDECOM

UNCLASSIFIED

OpenSHMEM



- Open standard API, portable across many platforms
- Library using Partitioned Global Address Space (PGAS) programming model
- Processing Element (Epiphany core) is an OpenSHMEM process
- Symmetric Objects have same address on all PEs, global shared memory





U.S. ARMY
RDECOM

UNCLASSIFIED

Dot Product



```
1  .global _dot_product
2  _dot_product:
3      movts lc, r2
4      mov r2, %low(.Lstart)
5      movts ls, r2
6      mov r2, %low(.Lend-4)
7      movts le, r2
8      .balignw 8,0x01a2
9      mov r24, #0
10     mov r25, #0
11     mov r26, #0
12     mov r27, #0
13     ldrd r48, [r0], #1
14     fsub r44, r24, r24
15     ldrd r50, [r1], #1
16     fsub r45, r24, r24
17     ldrd r52, [r0], #1
18     fsub r46, r24, r24
19     ldrd r54, [r1], #1
20     fsub r47, r24, r24
21     ldrd r56, [r0], #1
22     fmadd r24, r48, r50
23     ldrd r58, [r1], #1
24     fmadd r25, r49, r51
25     ldrd r60, [r0], #1
26     fmadd r26, r52, r54
27     ldrd r62, [r1], #1
28     fmadd r27, r53, r55
29     .Lstart:
30     ldrd r48, [r0], #1
31     fmadd r44, r56, r58
32     ldrd r50, [r1], #1
33     fmadd r45, r57, r59
34     ldrd r52, [r0], #1
35     fmadd r46, r60, r62
36     ldrd r54, [r1], #1
37     fmadd r47, r61, r63
38     ldrd r56, [r0], #1
39     fmadd r24, r48, r50
40     ldrd r58, [r1], #1
41     fmadd r25, r49, r51
42     ldrd r60, [r0], #1
43     fmadd r26, r52, r54
44     ldrd r62, [r1], #1
45     fmadd r27, r53, r55
46     .Lend:
47     fmadd r44, r56, r58
48     fmadd r45, r57, r59
49     fmadd r46, r60, r62
50     fmadd r47, r61, r63
51     fadd r24, r24, r25
52     fadd r26, r26, r27
53     fadd r44, r44, r45
54     fadd r46, r46, r47
55     fadd r24, r24, r26
56     fadd r44, r44, r46
57     fadd r0, r24, r44
58     rts
```

```
1  float dot_product(const float* a, const float* b, int nd8m1) {
2      float c = 0.0f;
3      int n = (nd8m1 + 1)*8;
4      for (int i = 0; i < n; i++) {
5          c += a[i] * b[i];
6      }
7      return c;
8  }
```

The assembly-optimized dot product (left) uses a hardware loop, 8-way loop unrolling, dual-issue pipelined loads/stores, fused multiply-adds. It is functionally identical to the C code (above). The inner loop achieves peak performance of the processor in both core bandwidth and computation. For N = 3040, that's an overall 96% peak performance with the call overhead, header, and footer code.



U.S. ARMY
RDECOM

Supplemental SDK Comparison



| Core/Device-Level Feature | OpenSHMEM 1.3 | eSDK (e-lib) |
|--|-------------------|-----------------------------------|
| Query Routines | Yes | Yes |
| Memory Management Routines | Yes | None |
| Remote Memory Access Routines | Yes | Yes |
| Non-Blocking Remote Memory Access Routines | Yes | Yes (needs better abstraction) |
| Atomic Memory Operations | Yes | None |
| Collective Routines | Yes | Barrier only |
| Point-To-Point Synchronization Routines | Yes | None |
| Memory Ordering Routines | Yes | Partial |
| Distributed Locking Routines | Yes | Partial |
| Signal/Interrupt Configuration Routines | Application code | Yes |
| Read/Write Special Registers | Application code | Yes |
| Timing Routines | Library Interface | Yes |



U.S. ARMY
RDECOM

UNCLASSIFIED

Supplemental SDK Comparison



| eSDK (e-lib interface) | OpenSHMEM |
|---|--|
| <code>e_coreid_t e_get_coreid(void)</code> | <code>int shmem_my_pe(void)</code> <code>int shmem_n_pes(void)</code> |
| <code>void *e_get_global_address(unsigned row, unsigned col, const void *ptr)</code> | <code>void *shmem_ptr(const void *target, int pe)</code> |
| <code>e_coreid_t e_coreid_from_coords(unsigned row, unsigned col)</code> <code>void e_coords_from_coreid(e_coreid_t coreid, unsigned *row, unsigned *col)</code> <code>e_bool_t e_is_on_core(const void *ptr)</code> <code>void e_neighbor_id(e_coreid_wrap_t dir, e_coreid_wrap_t wrap, unsigned *row, unsigned *col)</code> | <code>int shmem_pe_accessible(int pe)</code> <code>int shmem_addr_accessible(const void *addr, int pe)</code> |
| <code>unsigned e_ctimer_get(e_ctimer_id_t timer)</code> <code>unsigned e_ctimer_set(e_ctimer_id_t timer, unsigned int val)</code> <code>unsigned e_ctimer_start(e_ctimer_id_t timer, e_ctimer_config_t config)</code> <code>unsigned e_ctimer_stop(e_ctimer_id_t timer)</code> <code>void e_wait(e_ctimer_id_t timer, unsigned int clicks)</code> | No standard equivalent, but included in extended/experimental interface |
| <code>void *e_read(const void *remote, void *dst, unsigned row, unsigned col, const void *src, size_t n)</code> <code>void *e_write(const void *remote, const void *src, unsigned row, unsigned col, void *dst, size_t n)</code> <code>int e_dma_start(e_dma_desc_t *descriptor, e_dma_id_t chan)</code> <code>int e_dma_busy(e_dma_id_t chan)</code> <code>void e_dma_wait(e_dma_id_t chan)</code> <code>int e_dma_copy(void *dst, void *src, size_t n)</code> <code>void e_dma_set_desc(e_dma_id_t chan, unsigned config, e_dma_desc_t *next_desc, unsigned strd_i_src, unsigned strd_i_dst, unsigned count_i, unsigned count_o, unsigned strd_o_src, unsigned strd_o_dst, void *addr_src, void *addr_dst, e_dma_desc_t *desc)</code> | <code>shmem_put</code> <code>shmem_iput</code> <code>shmem_p</code> <code>shmem_[TYPE]_put</code> <code>shmem_[TYPE]_iput</code> <code>shmem_[TYPE]_p</code> <code>shmem_put[SIZE]</code> <code>shmem_iput[SIZE]</code> <code>shmem_putmem</code> <code>shmem_put_nbi</code> <code>shmem_[TYPE]_put_nbi</code> <code>shmem_put[SIZE]_nbi</code> <code>shmem_putmem_nbi</code> <code>shmem_get</code> <code>shmem_iget</code> <code>shmem_g</code> <code>shmem_[TYPE]_get</code> <code>shmem_[TYPE]_iget</code> <code>shmem_[TYPE]_g</code> <code>shmem_get[SIZE]</code> <code>shmem_iget[SIZE]</code> <code>shmem_getmem</code> <code>shmem_get_nbi</code> <code>shmem_[TYPE]_get_nbi</code> <code>shmem_get[SIZE]_nbi</code> <code>shmem_getmem_nbi</code> |
| <code>void e_irq_attach(e_irq_type_t irq, sighandler_t handler)</code> <code>void e_irq_set(unsigned row, unsigned col, e_irq_type_t irq)</code> <code>void e_irq_clear(unsigned row, unsigned col, e_irq_type_t irq)</code> <code>void e_irq_global_mask(e_bool_t state)</code> <code>void e_irq_mask(e_irq_type_t irq, e_bool_t state)</code> | No equivalent in SHMEM (but some of this occurs transparently within library, optionally) |
| <code>void e_mutex_init(unsigned row, unsigned col, e_mutex_t *mutex, e_mutexattr_t *attr)</code> <code>void e_mutex_lock(unsigned row, unsigned col, e_mutex_t *mutex)</code> <code>unsigned e_mutex_trylock(unsigned row, unsigned col, e_mutex_t *mutex)</code> <code>void e_mutex_unlock(unsigned row, unsigned col, e_mutex_t *mutex)</code> <code>void e_barrier_init(volatile e_barrier_t bar_array[], volatile e_barrier_t *tgt_bar_array[])</code> <code>void e_barrier(volatile e_barrier_t *bar_array, volatile e_barrier_t *tgt_bar_array[])</code> | <code>void shmem_barrier_all(void);</code> <code>void shmem_barrier(int PE_start, int logPE_stride, int PE_size, long *psync);</code> <code>void shmem_fence(void);</code> <code>void shmem_quiet(void);</code> <code>void shmem_set_lock(long *lock);</code> <code>void shmem_clear_lock(long *lock);</code> <code>int shmem_test_lock(long *lock);</code> |
| No equivalent on-chip memory management interface | <code>void *shmem_malloc(size_t size)</code> <code>void *shmem_align(size_t alignment, size_t size)</code> <code>void shmem_free(void *ptr)</code> <code>void *shmem_realloc(void *ptr, size_t size)</code> |



U.S. ARMY
RDECOM

Supplemental

ARL

| eSDK (e-lib interface) | OpenSHMEM |
|---|--|
| No equivalent atomic operations | shmem_[TYPE]_finc shmem_[TYPE]_inc shmem_[TYPE]_fadd shmem_[TYPE]_add shmem_[TYPE]_cswap shmem_[TYPE]_swap shmem_[TYPE]_fetch shmem_[TYPE]_set |
| No equivalent collective routines | shmem_broadcast[SIZE] shmem_fcollect[SIZE] shmem_collect[SIZE] shmem_alltoall[SIZE] shmem_alltoalls[SIZE] shmem_[TYPE]_sum_to_all shmem_[TYPE]_prod_to_all shmem_[TYPE]_min_to_all shmem_[TYPE]_max_to_all shmem_[TYPE]_and_to_all shmem_[TYPE]_or_to_all shmem_[TYPE]_xor_to_all |
| No equivalent point-to-point synchronization routines | shmem_[TYPE]_wait shmem_[TYPE]_wait_until |