

Class II Submission to the HPC Challenge Award Competition

Coarray Fortran 2.0

John Mellor-Crummey, Laksono Adhianto, Guohua Jin,
Mark Krentel, Karthik Murthy, William Scherer, and Chaoran Yang
Department of Computer Science
Rice University
Houston, TX, USA

1 General Description

We implemented four HPC Challenge benchmarks: Global HPL, Global RandomAccess, EP STREAM (Triad), and Global FFT in Coarray Fortran 2.0; in addition, we implemented an extra benchmark — Unbalanced Tree Search (UTS) [9]. Coarray Fortran 2.0 is an extension of Coarray Fortran developed at Rice University to improve expressiveness and performance by adding features such as teams, collectives, event-based synchronization, and asynchronous operations (including copy, collective operations, and function shipping) [7, 15]. In this submission, we have aimed to develop high performance solutions that are also clean and elegant. In particular, we accept additional code complexity in order to improve performance. Our use of a software routing algorithm for RandomAccess (see Section 3) exemplifies this design philosophy. Table 1 shows the total source lines of code and their breakdowns for each of our four HPCC benchmark implementations and UTS.

All benchmark codes were compiled using Rice’s Coarray Fortran 2.0 compiler and runtime system [13]. Our Coarray Fortran 2.0 compiler is a source-to-source compiler built using the Rose compiler infrastructure [12] developed by Lawrence Livermore National Laboratory. Our Coarray Fortran 2.0 runtime library is built on the UC Berkeley’s GASNet communication library [3]. We executed our benchmark codes on a range of Cray systems available to us over the past year. On the Cray XT4 and XT5, we used GASNet’s Portals and MPI conduits for communication; on the Cray XE6, we used GASNet’s new Gemini conduit, which is based on Cray’s GNI driver for Gemini. The version of GASNet and the conduit used for each set of experiments is specified in our table of results. For all of our experiments, we used Cray’s PGI programming environment and the Portland Group’s Fortran compiler for compiling the Fortran 90 codes generated by our CAF translator. Since our results were collected over time, different PGI compiler versions were used for various benchmarks; the version used for each set of experiments is specified in our table of results.

To evaluate the performance of our CAF 2.0 implementations of the FFT, and HPL benchmarks, we ran them on up to 4096 cores of Franklin, a Cray XT4 system at the National Energy Research Scientific Computing Center. Each node in Franklin contains a 2.3 GHz single socket quad-core AMD Opteron processor (Budapest) (theoretical peak performance of 9.2 GFlop/sec per core) and 2GB of memory per core. The memory speed is 800 MHz. Each node is connected to a dedicated SeaStar2 router through Hypertransport. The SeaStar2 interconnect is arranged as a 3D torus.

We evaluated our CAF 2.0 implementations of the STREAM and RandomAccess benchmarks on up to 4096 cores of Jaguar, a Cray XT4 system at Oak Ridge National Laboratory (ORNL). Jaguar contains 7,832 compute nodes. Each node contains a quad-core AMD Opteron 1354 (Budapest) processor running at 2.1 GHz. Some nodes use 8 GB of DDR2-800 memory, and others use DDR2-667 memory. Like Franklin, Jaguar’s nodes are also connected with a SeaStar2 router. We also experimented with STREAM on Hopper, a Cray XE6 that includes 6,384 nodes, each with two twelve-core AMD ‘MagnyCours’ 2.1-GHz processors and 32 GB DDR3 1333-MHz memory.

	STREAM	RandomAccess	FFT	HPL	UTS
Computation	32	188	180	536	267
Communication & synchronization	1	12	4	46	17
Declaration	17	118	103	109	151
Comments & spaces	13	91	163	95	109
Total	63	409	450	786	544

Table 1: Source lines of code for CAF 2.0 benchmarks.

# cores	STREAM (GByte/s)		RandomAccess (GUP/s)			FFT (GFLOP/s)	HPL (GFLOP/s)	UTS (MNode/s)
	Oct 2011 Jaguar(XT5)	Oct 2010 Franklin	Oct 2011 Hopper	Oct 2011 Hopper	Oct 2010 Jaguar(XT4)	Jul 2011 Franklin	Oct 2011 Franklin	Oct 2011 Jaguar(XT5)
	# cores/node							
	GASNet library							
	PGI(Gnu) compilers							
	12	4	16	4	4	4	4	12
	16.2	14.2	17.2	17.2	14.2	14.2	14.2	16.2
	10.9	10.0	11.7	11.7	10.0	10.0	10.0	10.9(4.4.3)
1								2.578
2		8.7				0.50		
4		8.5		0.0525	0.014	0.54	27.1	
8		17.0		0.0748	0.020	0.99		
16		34.0	0.119	0.117	0.031	1.77	95.9	
32		68.2	0.116	0.185	0.051	3.31		
64	82.0	137.6	0.150	0.303	0.084	6.69	363.5	163.1
128		272.5	0.212	0.434	0.14	11.90		325.8
256	327.0	544.4	0.339	0.691	0.24	22.82	1357.6	645.0
512		1089.8	0.456		0.44	38.61		1252
1024	1320	2179.2	0.746	1.84	0.64*	67.80	4990.7	2371
2048		4358.4	1.26	3.03	0.97*	97.49		4962
4096	7004	8733.4			1.67*	187.04	18333.2	7818
8192	14097				1.67*	357.80		12286

* These data for RandomAccess were measured on a larger version of the code that had many alternate implementations of verification. The routing algorithm is the same as that presented in the appendix and used to measure the other data points.

Table 2: Performance results of Coarray Fortran 2.0 implementation of the benchmarks on Cray XT.

Experiments with UTS were performed on the Cray XT5 partition of Jaguar, which consists of 18,688 compute nodes. Each compute node contains dual hex-core AMD Opteron 2435 (Istanbul) processors running at 2.6GHz, 16GB of DDR2-800 memory, and a SeaStar 2+ router.

Table 2 shows performance results of the four benchmarks running on up to 8192 cores of the aforementioned Cray systems.

2 EP STREAM Triad

The STREAM benchmark evaluates the extent to which a parallel system can deliver and sustain peak memory bandwidth by performing a simple vector operation that scales and adds two vectors:

$$a \leftarrow b + ac \tag{1}$$

Performance of the STREAM benchmark is measured in GByte/s, with the calculated performance defined as $24 \frac{m}{t_{min}} 10^{-9}$, where m is the size of the vectors, required to be at least a quarter of system memory; and

```

1 double precision, allocatable :: a(:)[*]
2 double precision, allocatable :: b(:)[*], c(:)[*]
3 ! allocate with the default team
4 allocate(a(ndim)[], b(ndim)[], c(ndim)[])
5 ...
6 do round = 1, rounds
7   do j = 1, rep
8     call triad(a,b,c,n,scalar)
9   end do
10  call team_barrier()
11 end do
12 ...
13 subroutine triad(a, b, c, n ,scalar)
14   double precision a(n), b(n), c(n), scalar
15   a = b + scalar * c
16 end subroutine triad

```

Figure 1: Implementation of STREAM benchmark.

t_{min} is the minimum execution time over at least 10 repetitions of the benchmark kernel. The STREAM benchmark is embarrassingly parallel; the work performed on any one node is independent of that performed on others.

Since the STREAM benchmark does not require communication between processes, the Coarray Fortran version is essentially identical to the sequential Fortran implementation, with the exception that all arrays are declared and allocated as coarrays. A sketch of the essence of our implementation is shown in Figure 1.

To deliver top performance, we outlined the STREAM triad calculation from the timing loop. Since our CAF 2.0 compiler presently represents coarray data using F90 pointers, having the triad computation in a separate routine enabled us to inform the compiler that the arrays were contiguous data by using explicit shape array declarations within triad. Eventually, we will either automatically outline coarray computations for performance, or use the Fortran 2008 CONTIGUOUS attribute to inform the back-end compiler about the contiguity of coarray data and avoid the need for outlining.

We initialized b and c with identical values (generated with the random number generator from RandomAccess), and set the scalar α to -1, so that the result should be 0. We considered the calculation verified if the maximum value of the difference between the computed value and 0 was less than 10^{-9} .

3 RandomAccess

The RandomAccess benchmark evaluates the rate at which a parallel system can apply updates to randomly indexed entries in a distributed table. Performance of the RandomAccess benchmark is measured in Giga Updates Per Second (GUP/s). GUP/s is calculated by identifying the number of table entries that can be randomly updated in one second, divided by 1 billion (1e9). The term “randomly” means that there is little relationship between one table index to be updated and the next. An update is a read-modify-write operation on a 64-bit word in the table. First, a table index is generated using a random number generator. Then, the table value at that index is combined with a literal value using an xor and the resulting value is written back to memory.

On distributed-memory parallel systems that lack hardware support for shared memory, fine-grain operations on remote data are expensive. To develop a high performance implementation of RandomAccess in CAF 2.0, we exploit the “1024 element look ahead and storage” allowed by the problem specification. First, each process image generates a batch of 1024 indices of table locations to be updated. Next, the code uses a hypercube-based pattern of bulk communication to route each update to the appropriate process image co-located with the table index being updated. Finally, each process image locally applies updates to its section of the distributed table.

Figure 2 shows a sketch of the hypercube routing algorithm used in our CAF 2.0 implementation of RandomAccess. Routing occurs in $\log P$ rounds. As each round begins, each processor has a set of elements. For each element, the processor must decide whether to communicate or retain in it the current round. Each pro-

```

1  event,allocatable :: delivered(:)[*],received(:)[*]
2  integer(8),allocatable :: fwd(:,:)[*]
3
4  do i = world_logsize-1, 0, -1
5    ...
6    call split(ret(:,last), retsizes(last), ret(:,current), retsizes(current), &
7              fwd(1:,out,i),fwd(0,out,i),bufsize,dist)
8
9    if (i < world_logsize-1) then
10   event_wait(delivered(i+1))
11   call split(fwd(1:,in,i+1), fwd(0,in,i+1), ret(:,current), retsizes(current), &
12             fwd(1:,out,i),fwd(0,out,i),bufsize,dist)
13   event_notify(received(i+1)[from])
14   endif
15
16   outgoing_size = fwd(0,out,i)
17   call event_wait(received(i))
18   fwd(0:outgoing_size,in,i)[partner] = fwd(0:outgoing_size,out,i)
19   call event_notify(delivered(i)[partner])
20   ...
21 end do

```

Figure 2: A sketch of hypercube routing for updates in CAF 2.0 RandomAccess.

cessor splits the updates in its possession at the time into two sets: those to retain and those to communicate in the current round. After the first round, a processor uses `event_wait` to wait for communication from the prior round to complete. Once a set of updates have been received from a communication partner, they are also split for communication or retention. Following the second split, each processor first waits for its communication partner’s buffer to be available by checking that its prior message to this communication partner has been already received. Each process finishes a round of routing by copying a vector of updates into its partner’s coarray and notifying its partner that they have been delivered. Our CAF 2.0 compiler translates this copy/`event_notify` pair into a non-blocking, asynchronous put-with-notify operation supported by the CAF 2.0 runtime.

Similar software routing strategies have been used before, though never with CAF. Researchers at Sandia studied a different but related strategy for RandomAccess using all-to-all communication based on a hypercube communication pattern [11]. IBM also explored a software routing strategy for the RandomAccess benchmark on Blue Gene systems [4].

4 Global FFT

The HPC Challenge *FFT* (Fast Fourier Transform) benchmark measures the ability of a system to overlap computation and communication while calculating a very large Discrete Fourier Transform of size m with input vector z_i and output vector Z_i :

$$Z_k \leftarrow \sum_j^m z_j e^{-2\pi i \frac{jk}{m}}; 1 \leq k \leq m$$

Performance of the FFT benchmark is measured in GFLOP/s, with calculated performance defined as $5 \frac{m \log_2 m}{t} 10^{-9}$, where m is the size of the DFT and t is the execution time (in seconds). The number of processors for this benchmark may be implementation-specific; in particular, it is allowed to be an integral power of 2. Parallel FFT algorithms has been well studied in the past [16, 2, 6]. The reference FFT implementation of the HPC Challenge benchmarks uses a 1D algorithm based on [16].

Our CAF 2.0 FFT implementation uses a radix 2 binary exchange formulation that consists of five parts: permutation of data to move each source element to the position that is its binary bit reversal; local (in-core) FFT computation for levels that use data co-located with a single CAF process image; transposition of the data from block to cyclic layout; FFT computation for the remaining layers (which now requires only local data); and reverse transposition to restore data from cyclic to block layout. This is shown in Figure 3.

```

1      complex, allocatable :: c(:)[*]
2      ...
3      ! permute local data
4      ...
5      ! calculate twiddle table
6      ...
7      ! compute levels of FFT using the data co-located with a single CAF process image
8      ...
9      ! transpose data from block to cyclic distribution
10     ...
11     ! compute remaining levels
12     ...
13     ! transpose data back from cyclic to block distribution

```

Figure 3: Sketch of FFT computation.

5 Global HPL

In our Coarray Fortran 2.0 implementation of HPL, the main matrix is declared as a coarray and is distributed in block-cyclic fashion in both dimensions. We use a single block size for both dimensions. However, given a processor core topology, the number of matrix blocks mapped to each processor in different dimensions can be different. For simplicity, we also use the same block size of the distribution as the width of the panel in the LU panel factorization with row partial pivoting. Choosing the right block size for the block-cyclic matrix distribution is very important. The block size determines the efficiency of core computation, the load balance between processors, and the communication latency exposed. This is machine and algorithm dependent. Our implementation takes an optional block size as an input parameter. For the final runs of the experiments we used 92 double precision elements as the block size.

We use Coarray Fortran 2.0 teams to represent processor subsets. We create column teams for processors that own the same matrix column and row teams for processors that own the same matrix row. All collective operations are performed within pre-arranged teams for high efficiency. This is particularly important for achieving scalable high performance on large machines. To hide communication latency, we used a team-based asynchronous broadcast operation `team.broadcast_async` to broadcast panels. This operation is one of the asynchronous collective operations that we have designed into Coarray Fortran 2.0. We used a user-defined all-reduce MAX operation for implementing the row partial pivoting and broadcasting the pivot elements. We used Cray Scientific Libraries package, LibSci mainly for computing matrix-vector multiplication (`dger`) and matrix-matrix multiplication (`dgemm`).

Our implementation allocates a matrix with $12K \times 12K$ double precision elements on each processor, more than half of the 2GBytes system memory. Other than the main matrix, coarrays are also used as buffers for panel broadcast and communication. A double panel buffer is used in order to perform panel factorization and update of trialing matrix in parallel. Random numbers are generated for the matrix during its initialization phases before the main computation and during verification. We verified and used 1.11×10^{-16} as machine precision for 64-bit floating-point values. We verified our results by checking the scaled residual formulated in [5].

6 Unbalanced Tree Search (UTS)

The UTS benchmark involves building and counting the number of nodes in an unbalanced n-ary tree. The tree is based on a geometric/binomial distribution. Each node in this tree is characterized by a 20-byte descriptor that is used to determine the number of children for the node and (indirectly) for all of its children. Descriptors are calculated from an SHA1 hash of the parent descriptor and the child's index; a more detailed explanation may be found in the UTS v1.1 benchmark [1]). Since each node's children are completely determined by its descriptor, the parent-child links need not be explicitly maintained: The UTS tree is virtual.

Our implementation of UTS is based on the T1WL UTS benchmark. We implement a UTS tree based

```

1      !while there is work to do
2  do while(queue_count .gt. 0)
3      delete_queue_end(descriptor, depth)
4      call process_work_item(descriptor, depth)
5      ...
6      !check if someone needs work
7      if ((incoming_lifeline .ne. 0) .and. (queue_count .ge. lifeline_threshold)) then
8          call push_work()
9      endif
10     enddo
11
12     ! attempting to steal work from another image
13     steal_from_img = get_random_image(my_rank, my_rank)
14     spawn steal_work_spawn(my_rank, 0)[steal_from_img]
15
16     ! set up lifelines
17     neighbor_index = 0
18     do while (neighbor_index .lt. max_neighbor_index)
19         next_neighbor = mod(my_rank+(2**neighbor_index), world_size)
20         spawn set_lifelines(my_rank, neighbor_index)[next_neighbor]
21         neighbor_index = neighbor_index + 1
22     enddo

```

Figure 4: Outline of the UTS benchmark.

on a geometric distribution with expected children size per node 4, maximum tree depth 18, and initial seed for the root descriptor 19.

The highlights of our implementation are as follows:

1. Initial work sharing

Process 0 builds the initial few levels of the UTS tree and then distributes these nodes to different processes. Each process then starts working on building the rest of the tree rooted at these nodes. Process 0 does this by spawning the function *copy_item_and_activate* on the processes.

2. Random work stealing

Each process, after finishing its quota of work, tries to steal work (i.e. nodes whose subtree has not been constructed) from a random process. If a process is unsuccessful in n steal attempts (set to 1 in these experiments) it then quiesces by setting up *lifelines* ([14]). Line 12 in Figure 4 depicts an attempt to steal from a remote node.

3. Work sharing via lifelines

Process A sets up a *lifeline* on another image B to indicate that A is available for any excess work that B obtains in the future. In our implementation, a process sets up lifelines on each of its hypercube neighbors (the processes at offsets $2^0, 2^1, \dots, 2^{\log_2(\text{num-processes})}$). Process A uses function shipping to establish lifelines on the remote nodes. Since the spawn is executed on the process on which the lifeline has to be set up, the communication overhead for this operation is reduced to a single round trip.

6.1 Challenges posed by the benchmark

The geometric distribution of nodes in the tree causes there to be a great deal of imbalanced work in this application. Work stealing in general is a good strategy, but stealing alone cannot succeed for the following reasons:

- Amount to steal

The stolen work might yield only a small subtree before petering out. In effect, the time spent in a network transaction to check if a process has work, lock the process's work queue, and then steal the work might not be fruitful. Few papers ([8], [10]) papers have looked at this in detail. Our steals are limited by the *ActiveMessageMediumPacket* size (a max of 9 items).

- Timing of the steal
The steal attempt from a process A might land at process B when B is working on a few nodes. Concluding that B has a small amount of work would be erroneous if those nodes give rise to a large subtree.

References

- [1] Uts 1.1 benchmark. <http://barista.cse.ohio-state.edu/wiki/index.php/UTS>.
- [2] R. C. Agarwal, F. G. Gustavson, and M. Zubair. A high performance parallel algorithm for 1-D FFT. In *Supercomputing '94: Proceedings of the 1994 conference on Supercomputing*, pages 34–40, Los Alamitos, CA, USA, 1994. IEEE Computer Society Press.
- [3] Dan Bonachea. GASNet specification, v1.8. Technical report, University of California at Berkeley, Berkeley, CA, USA, 2006.
- [4] Rahul Garg and Yogish Sabharwal. Software routing and aggregation of messages to optimize the performance of HPCC randomaccess benchmark. In *SC '06: Proceedings of the 2006 ACM/IEEE conference on Supercomputing*, New York, NY, USA, 2006. ACM.
- [5] HPC challenge awards: Class 2 specification. <http://www.hpcchallenge.org/class2specs.pdf>, June 2005.
- [6] S. Lennart Johnsson and Robert L. Krawitz. Cooley-Tukey FFT on the Connection Machine. *Parallel Computing*, 18:1201–1221, 1991.
- [7] John Mellor-Crummey, Laksono Adhianto, William N. Scherer, III, and Guohua Jin. A new vision for Coarray Fortran. In *PGAS '09: Proceedings of the Third Conference on Partitioned Global Address Space Programming Models*, pages 1–9, New York, NY, USA, 2009. ACM.
- [8] Min, Iancu Seung-Jai, Yelick Costin, and Katherine. Hierarchical work stealing on manycore clusters. In *Proceedings of Fifth Conference on Partitioned Global Address Space Programming Models*, PGAS'11, 2011.
- [9] Stephen Olivier, Jun Huan, Jinze Liu, Jan Prins, James Dinan, P. Sadayappan, and Chau-Wen Tseng. UTS: An unbalanced tree search benchmark. *Lecture Notes in Computer Science*, 4382/2007:235–250, 2007.
- [10] Stephen Olivier, Jun Huan, Jinze Liu, Jan Prins, James Dinan, P. Sadayappan, and Chau-Wen Tseng. Uts: an unbalanced tree search benchmark. In *Proceedings of the 19th international conference on Languages and compilers for parallel computing*, LCPC'06, pages 235–250, Berlin, Heidelberg, 2007. Springer-Verlag.
- [11] Steven J. Plimpton, Ron Brightwell, Courtenay Vaughan, Keith D. Underwood, and Mike Davis. A simple synchronous distributed-memory algorithm for the HPCC RandomAccess benchmark. In *CLUSTER*. IEEE, 2006.
- [12] Daniel J. Quinlan. Rose: Compiler support for object-oriented frameworks. *Parallel Processing Letters*, 10(2/3):215–226, 2000.
- [13] Rice Coarray Fortran 2.0. <http://caf.rice.edu>, 2010.
- [14] Vijay A. Saraswat, Prabhanjan Kambadur, Sreedhar Kodali, David Grove, and Sriram Krishnamoorthy. Lifeline-based global load balancing. In *Proceedings of the 16th ACM symposium on Principles and practice of parallel programming*, PPOPP '11, pages 201–212, New York, NY, USA, 2011. ACM.

- [15] William N. Scherer, III, Laksono Adhianto, Guohua Jin, John Mellor-Crummey, and Chaoran Yang. Hiding latency in coarray fortran 2.0. In *PGAS '10: Proceedings of the Third Conference on Partitioned Global Address Space Programing Models*, New York, NY, USA, 2010.
- [16] Daisuke Takahashi and Yasumasa Kanada. High-performance radix-2, 3 and 5 parallel 1-D complex FFT algorithms for distributed-memory parallel computers. *J. Supercomput.*, 15(2):207–228, 2000.

A EP STREAM Triad code

A.1 stream.caf

```
1  program stream
2  use module_triad
3  implicit none
4  integer(8), parameter :: POLY = 7, offset=0, rounds=10, rep=10
5  integer(8), parameter :: n=(256 * 1024 * 1024 / 18)
6
7  double precision, allocatable :: a(:)[*], b(:)[*], c(:)[*]
8  double precision :: tmin, scalar
9
10 integer(8):: i, j, round, ndim, me, start_time, times(rounds), rate,ran
11 real(8) rnd
12
13 ndim = n + offset
14 allocate(a(ndim)[], b(ndim)[], c(ndim)[])
15
16 ! initialize
17 me = team_rank()
18 call random_number(rnd) ! the range of rnd is 0 ... 1
19 ran = rnd *(me+1)*n      ! the range of ran is 0 ... n(n+1)
20 scalar = -1.0D0
21
22 do i = 1, n
23   ran = xor(ishft(ran,1), iand(-ishft(ran,-63),POLY))
24   b(i) = ran
25 end do
26 c = b
27
28 times = 0.0D0
29
30 ! compute the triad
31 do round = 1, rounds
32   call system_clock(start_time, rate)
33   do j = 1, rep
34     call triad(a,b,c,n,scalar)
35   end do
36   call system_clock(times(round), rate)
37   times(round) = times(round) - start_time
38   call team_barrier()
39 end do
40 tmin = minval(times)
41
42 if (me == 0) then
43   print *, 'min_time=',tmin
44   print *, 'n=',ndim,"total_mem_usage:",ndim*24
45   print *, 'BW/PE(MB/s,WC)=' ,24.0D0*n*rep/tmin*1D-6*rate
46 end if
47
48 ! verification
49 a = abs(a)
50 if (maxval(a) > 1e-9) then
51   print *, 'Verification failed'
52 end if
53
54 deallocate(a, b, c)
55 end program stream
```

A.2 module_triad.caf

```
1  module module_triad
2  contains
3  subroutine triad(a, b, c, n, scalar)
4     integer n
5     double precision a(n), b(n), c(n), scalar
6     a = b + scalar * c
7  end subroutine triad
8  end module module_triad
```

B RandomAccess code

B.1 randomaccess.caf

```
1      subroutine print_problem_stats(nupdates_local)
2          use module_team
3          use module_table
4
5          implicit none
6          integer(8) :: nupdates_local
7
8          if (world_rank == 0) then
9              write(*,'(A)') "randomaccess_benchmark"
10             write(*,'(A,I18,A,I18,A)') "global_table_size=", local_table_size * 8 * world_size, &
11                 "bytes," , local_table_size * world_size, "elements"
12             write(*,'(A,I18,A,I18,A)') "local_table_size=", local_table_size * 8, &
13                 "bytes," , local_table_size, "elements"
14             write(*,'(/,A,I18)') "global_updates=", nupdates_local * world_size
15             write(*,'(A,I18)') "local_updates=", nupdates_local
16         endif
17     end subroutine print_table
18
19     subroutine print_gups(start_time, end_time, rate, updates, style)
20         use module_team
21
22         implicit none
23         real, parameter :: billion = 1e9
24         integer(8)      :: updates
25         integer         :: start_time, end_time, ticks, rate
26         real            :: GUPS
27         character(*)    :: style
28
29         ticks = end_time - start_time
30         if (world_rank == 0) then
31             GUPS = (updates * 1.0 * world_size * rate) / (ticks * billion)
32             write(*,'(/,A,I18,A)') "clock_rate=", rate, "ticks_per_second"
33             write(*,'(A,I18)') "ticks=", ticks
34             write(*,'(A,F18.2,A)') "elapsed_time=", ticks/(1.0 * rate), "seconds"
35             write(*,'(A,A,E18.4)') style, ":GUPS=", GUPS
36         endif
37     end subroutine
38
39     subroutine perform_updates(nupdates_local, bunch_size, style)
40         use module_team
41         use module_route
42         use module_random_seq
43
44         implicit none
45         integer         :: start_time, end_time, rate, bunch_size
46         integer(8)      :: nupdates_local, start_pos, round
47         character(*)    :: style
48
49         dropped_updates = 0
50
51         start_pos = world_rank * nupdates_local
52         call random_seq_set_position(start_pos)
53
54         call team_barrier()
55         call system_clock(start_time, rate)
56
57         do round = 1, nupdates_local - bunch_size + 1, bunch_size
58             call random_seq_get_bunch(retain(1:bunch_size,0))
59             call route
60         enddo
61
62         call team_barrier()
63         call system_clock(end_time)
64
65         call print_gups(start_time, end_time, rate, nupdates_local, style)
66     end subroutine perform_updates
67
68
69     program randomaccess
70         use module_team
71         use module_table
72         use module_split
```

```

73     use module_route
74
75     implicit none
76     integer, parameter :: bunch_size = 1024
77     character(len=32)  :: arg
78     integer           :: local_table_bits
79     integer(8)       :: nupdates_local, error_count, error_bound
80
81     call team_init()
82
83     ! arg 1 specifies the number of elements in the local table as a power of 2
84     local_table_bits = 27
85     if (command_argument_count() > 0) then
86         call get_command_argument(1, arg)
87         if (len_trim(arg) .ne. 0) then
88             read(arg, '(I2)') local_table_bits
89         endif
90     endif
91
92     call table_init(local_table_bits)
93     call route_init(bunch_size)
94
95     nupdates_local = ishft(1_8, local_table_bits + 2)
96
97     call print_problem_stats(nupdates_local)
98
99     ! timed updates
100    call perform_updates(nupdates_local, bunch_size, "routing")
101
102    ! verification step: apply updates again (gratuitously timed)
103    call perform_updates(nupdates_local, bunch_size, "routing_verification")
104
105    call count_update_errors(error_count, error_bound, dropped_updates)
106
107    if (world_rank == 0) then
108        write(*, '(/,A,I18)')      "errors_found", error_count
109        write(*, '(/,A,I18)')      "error_upper_bound", error_bound
110    endif
111
112    end program randomaccess

```

B.2 module_route.caf

```

1  module module_route
2      use module_team
3      use module_split
4      integer(8), allocatable, dimension(:,:) :: retain
5      integer(8), dimension(0:1) :: retain_sizes
6
7      event, allocatable, dimension(:) :: delivered[*]
8      event, allocatable, dimension(:) :: received[*]
9      integer(8), allocatable, dimension(:,:,) :: fwd[*]
10
11     integer :: buffer_size, nelements_per_bunch
12     integer :: last, current
13     contains
14
15     subroutine route_init(nelements_per_bunch_)
16         nelements_per_bunch = nelements_per_bunch_
17         buffer_size = max(2000000, 2 * nelements_per_bunch)
18         allocate(retain(buffer_size, 0:1))
19         allocate(fwd(0:buffer_size, 0:1, 0:world_logsize-1)[])
20         allocate(delivered(0:world_logsize-1)[])
21         allocate(received(0:world_logsize-1)[])
22         call event_init(delivered, world_logsize)
23         call event_init(received, world_logsize)
24         do i = world_logsize-1, 0, -1
25             call event_notify(received(i))
26         enddo
27     end subroutine route_init
28
29     subroutine route
30         use module_bits
31         integer, parameter :: out = 0, in = 1
32         integer :: distance, partner, i, from
33         integer :: outgoing_size, itemp, n, iii, iik, advance, next_slot

```

```

34
35     advance(n) = 1 - n ! statement function to advance indices
36
37     last = 0
38     current = 1
39
40     retain_sizes(last) = nelements_per_bunch
41
42     do i = world_logsize-1, 0, -1
43         distance = ishft(1, i)
44
45         partner = mod(world_rank + distance + world_size, world_size)
46
47         retain_sizes(current) = 0
48         fwd(0,out,i) = 0 ! zero count of elements in outgoing buffer
49
50         !-----
51         ! partition elements retained in last routine into ones to retain
52         ! or forward in the current round
53         !-----
54         call split(retain(:,last), retain_sizes(last), &
55                  retain(:,current), retain_sizes(current), &
56                  fwd(1:,out,i), fwd(0,out,i), buffer_size, distance)
57
58         !-----
59         ! partition elements forwarded in last routine into ones to retain
60         ! or forward in the current round
61         !-----
62         if (i < world_logsize-1) then
63             call event_wait(delivered(i+1))
64             call split(fwd(1:,in,i+1), fwd(0,in,i+1), &
65                      retain(:,current), retain_sizes(current), &
66                      fwd(1:,out,i), fwd(0,out,i), buffer_size, distance)
67             call event_notify(received(i+1)[from])
68         endif
69
70         call event_wait(received(i))
71         outgoing_size = fwd(0,out,i)
72         fwd(0:outgoing_size,in,i)[partner]= fwd(0:outgoing_size,out,i)
73         call event_notify(delivered(i)[partner])
74
75         from = mod(world_rank - distance + world_size, world_size)
76
77         last = advance(last)
78         current = advance(current)
79     end do
80
81     call apply_updates(retain(:,last), retain_sizes(last))
82
83     call event_wait(delivered(0))
84     call apply_updates(fwd(1:,in,0), fwd(0,in,0))
85     from = mod(world_rank - 1 + world_size, world_size)
86     call event_notify(received(0)[from])
87
88     end subroutine route
89
90 end module route

```

B.3 module_split.caf

```

1     module module_split
2         use module_table
3         use module_team
4         integer(8) :: dropped_updates
5     contains
6
7         subroutine split(in, in_s, keep, keep_s, fwd, fwd_s, buffersize, distance)
8
9             implicit none
10            integer(8), dimension(:) :: in, keep, fwd
11            integer(8) :: in_s, keep_s, fwd_s
12            integer :: buffersize, distance, partner_rank, i, target_rank
13
14            partner_rank = world_rank + distance
15
16            do i = 1, in_s

```

```

17         target_rank = ishft(in(i), -local_table_logsize)
18         target_rank = iand(target_rank, world_size - 1)
19
20         if (target_rank < world_rank) then
21             target_rank = target_rank + world_size
22         endif
23         if (target_rank - world_rank < distance) then
24             if (keep_s < buffersize) then
25                 keep_s = keep_s + 1
26                 keep(keep_s) = in(i)
27             else
28                 dropped_updates = dropped_updates + 1
29             endif
30         else
31             if (fwd_s < buffersize) then
32                 fwd_s = fwd_s + 1
33                 fwd(fwd_s) = in(i)
34             else
35                 dropped_updates = dropped_updates + 1
36             endif
37         endif
38     enddo
39 end subroutine split
40 end module module_split

```

B.4 module_table.caf

```

1  module module_table
2      use module_team
3      integer(8), allocatable :: table(:)[*]
4      integer(8) :: local_table_size, local_table_logsize
5  contains
6
7      subroutine table_init(local_n)
8          implicit none
9          integer :: local_n
10         integer(8) :: start_index, i
11
12         local_table_size = ishft(1_8, local_n)
13         local_table_logsize = local_n
14         allocate(table(0:local_table_size-1)[])
15         start_index = world_rank * local_table_size
16         do i = 0, local_table_size-1
17             table(i) = start_index + i
18         end do
19     end subroutine table_init
20
21     subroutine apply_updates(buffer, size)
22         implicit none
23         integer(8) :: buffer(:)
24         integer(8) :: size, index, i
25         do i = 1, size
26             index = iand(buffer(i), local_table_size - 1)
27             table(index) = ieor(table(index), buffer(i))
28         end do
29     end subroutine apply_updates
30
31     subroutine reduce_sum(in1, in2, out, nbytes)
32         integer(8) :: in1(2), in2(2), out(2)
33         integer :: nbytes
34         out = in1 + in2
35     end subroutine reduce_sum
36
37     subroutine count_update_errors(error_count, error_bound, dropped_updates)
38         integer(8) :: error_count, error_bound, dropped_updates, i, start_index
39         integer(8) :: errors(2)
40         errors = 0
41         start_index = world_rank * local_table_size
42         do i = 0, local_table_size-1
43             if (table(i) .ne. start_index + i) then
44                 errors = errors + 1
45             end if
46         end do
47
48         errors(2) = errors(2) + dropped_updates
49         call team_reduce(errors, errors, 0, reduce_sum)

```

```

50
51     error_count = errors(1)
52     error_bound = errors(2)
53     end subroutine count_update_errors
54
55 end module module_table

```

B.5 module_random_seq.caf

```

1  module module_random_seq
2      integer(8) :: POLY
3      integer(8) :: PERIOD
4      integer(8) :: ran      ! next random number in the sequence
5  contains
6
7      !-----
8      ! generate the random number following val in a sequence specified
9      ! by the generator
10     !-----
11     integer(8) function random_seq_get_next(val)
12         integer(8) :: val
13         random_seq_get_next = ieor(ishft(val,1), iand(-ishft(val,-63), POLY))
14     end function
15
16     !-----
17     ! fast forward to the nth random number that would be returned by
18     ! get_next_random in a sequence beginning with the value 1.
19     ! this code is a Fortran version of the C function urng() found
20     ! in the HPC Challenge benchmark reference implementation.
21     !-----
22     subroutine random_seq_set_position(nth)
23         implicit none
24         integer(8) :: nth, n, temp
25         integer(8) :: m2(0:63)
26         integer :: i,j
27
28         POLY = 7
29         PERIOD = 1317624576693539401_8
30
31         n = nth
32         do while (n < 0)
33             n = n + PERIOD
34         enddo
35
36         do while (n > PERIOD)
37             n = n - PERIOD
38         enddo
39
40         if (n == 0) then
41             ran = 1
42         else
43             temp = 1
44             do i = 0, 63
45                 m2(i) = temp
46                 temp = random_seq_get_next(temp)
47                 temp = random_seq_get_next(temp)
48             enddo
49
50             do i = 62, 0, -1
51                 if (btest(n, i)) exit
52             enddo
53
54             ran = 2
55             do while (i > 0)
56                 temp = 0
57                 do j = 0, 63
58                     if (btest(ran, j)) temp = ieor(temp, m2(j))
59                 enddo
60                 ran = temp
61                 i = i - 1
62                 if (btest(n,i)) ran = random_seq_get_next(ran)
63             enddo
64         endif
65     end subroutine random_seq_set_position
66
67     subroutine random_seq_get_bunch(rand)

```

```

68     integer(8), dimension(:) :: rand
69     do i = 1, size(rand)
70         rand(i) = ran
71         ran = random_seq_get_next(ran)
72     enddo
73 end subroutine random_seq_get_bunch
74
75 end module module_random_seq

```

B.6 module_team.caf

```

1  module module_team
2  use module_bits
3  integer :: world_rank, world_size, world_logsize
4  contains
5
6  subroutine team_init()
7      world_rank = team_rank()
8      world_size = team_size()
9      world_logsize = number_of_bits(world_size - 1)
10 end subroutine team_init
11
12 end module module_team

```

B.7 module_bits.caf

```

1  module module_bits
2  contains
3  !-----
4  ! number_of_bits: compute the index of the leftmost non-zero bit in i
5  !-----
6  function number_of_bits(i) result(num)
7      num = 0
8      itmp = i
9      do while (itmp .gt. 0)
10         num = num + 1
11         itmp = ishft(itmp, -1)
12     end do
13 end function number_of_bits
14
15 end module module_bits

```

C Global FFT

C.1 fft.caf

```

1  program fft_driver
2  use module_fft
3  use module_reverse
4  implicit none
5
6  integer(8) :: n, two_n, world_size, world_logsize
7  integer(8) :: local_size, local_logsize
8  integer(8) :: tstart, tend, rate, mystart, rank
9  double precision :: tdelta, gflops, tsec
10 character(len=32) :: arg
11
12 world_size = team_size()
13 world_logsize = number_of_bits(world_size)-1
14 rank = team_rank()
15
16 ! by default the local core will have 1M of elements
17 local_logsize = 21
18
19 !-----
20 ! get the command line
21 !-----
22 if (0 < command_argument_count()) then
23     call get_command_argument(1, arg)
24     if (len_trim(arg) .ne. 0) then
25         read(arg, '(I2)') local_logsize
26     endif
27 endif

```

```

28
29
30 !-----
31 ! INIT: compute local size
32 !-----
33 ! from hpcc spec: the number of elements should be at lest
34 !  $2^5 * m \leq 0.25 * \text{local\_mem\_size}$ ,
35 ! since  $m = 2^{\text{local\_logsize}}$ , then:
36 !  $2^{\text{local\_logsize}} \leq 2^{-7} * \text{local\_memsize}$ 
37 !  $\text{local\_logsize} \leq \log_2(\text{local\_memsize}) - 7$ 
38 local_size = ishft(1,local_logsize)
39
40 !-----
41 ! INIT: compute my part of computation
42 !-----
43 call bitinit()
44 mystart = rank * local_size
45 call fft_init(local_size, mystart, rank)
46
47 !-----
48 ! computation
49 !-----
50 call system_clock(tstart, rate)
51 call fft_inner(local_size, 1_8)
52 call system_clock(tend, rate)
53 call fft_verif(world_size, local_size, mystart)
54
55 !-----
56 ! timing output
57 !-----
58 if ((0 .eq. rank) .and. (local_logsize > 12)) then
59     n = local_logsize + world_logsize
60     two_n = ishft(1,n)
61
62     tsec = real(tend-tstart) / rate
63     write (*,'(A,F8.4)') "Elapsed_time: ", tsec
64     gflops = ((5 * n * two_n) / tsec) * 1D-9
65     write (*, '(A,I7,A,I9,A,F8.4)') "Num_PEs: ", world_size, &
66                                     "; Local_size: ", local_size, "; GFlops: ", gflops
67
68 endif
69 call team_barrier()
70
71 end program fft_driver

```

C.2 module_fft.caf

```

1 module module_fft
2     use module_permute
3     use module_bits
4     use module_transpose
5
6     complex(8), allocatable, dimension (:) :: c[*] ! data array
7     complex(8), allocatable, dimension (:) :: spare[*] ! scratch space
8     complex(8), allocatable, dimension (:) :: twiddles
9
10     ! controls whether the fft calculated is a known pattern
11     ! with a known result
12     logical, parameter :: debug_mode = .false.
13 contains
14
15 !-----
16 ! Problem initialization
17 !-----
18 subroutine initialize_data_array(n_local_size, local_start, rank, buffer)
19     implicit none
20     integer(8) :: n_local_size, local_start, rank, i, j
21     integer :: n
22     complex(8) :: buffer(0:n_local_size - 1)
23     integer(4), parameter :: SEED = 314159265_4
24     real(8) :: h, h2
25     integer, dimension(:), allocatable :: rseeds
26
27     ! Initialize random number generator
28     call random_seed(size=n)
29     allocate(rseeds(n))
30     do j = 1, n
31         rseeds(j) = (rank + 1) * SEED * j

```



```

32     end do
33     call random_seed(PUT = rseeds)
34     deallocate(rseeds)
35
36     ! Initialize data with either random numbers or an ascending sequence
37     do i = 0, n_local_size - 1
38         if (.not. debug_mode) then
39             call random_number(h)
40             call random_number(h2)
41             buffer(i) = cmplx(h, h2)
42         else
43             buffer(i) = cmplx(i + local_start + 1)
44         endif
45     enddo
46 end subroutine initialize_data_array
47
48
49 subroutine fft_init(n_local_size, local_start, rank)
50     implicit none
51     integer(8) :: n_local_size, i, local_start, rank
52
53     ! Allocate memory
54     allocate(twiddles(0:n_local_size/2 - 1))
55     allocate(c(0:n_local_size - 1)[])
56     allocate(spare(0:n_local_size - 1)[])
57
58     ! Initialize data (in a separate routine so we can reinitialize
59     ! at verification time)
60     call initialize_data_array(n_local_size, local_start, rank, c)
61 end subroutine
62
63
64 !-----
65 ! verification
66 !-----
67 subroutine fft_verif(world_size, n_local_size, local_start)
68     implicit none
69     integer(8) :: i, j, rank, world_size, n_local_size, local_start, mei
70     real(8)      :: norm, error, max_error, residue, logm
71     real(8), parameter :: Epsilon = 1.1E-16
72     real(8), parameter :: max_residue = 16.0
73
74
75     ! Perform inverse fft
76     ! c = c / (world_size * n_local_size)
77     do i = 0, n_local_size - 1
78         c(i) = c(i) / (world_size * n_local_size)
79     end do
80     call fft_inner(n_local_size, -1.8)
81
82     ! Regenerate source data
83     rank = team_rank()
84     call initialize_data_array(n_local_size, local_start, rank, spare)
85
86
87     ! Find max error for infinity norm
88     max_error = -1.0; mei = -1
89     do i = 0, n_local_size - 1
90         error = abs(c(i) - spare(i))
91         if (error .gt. max_error) then; mei = i; endif
92         max_error = max(max_error, error)
93     end do
94
95     ! Calculate residue
96     logm = number_of_bits(world_size) - 1 + number_of_bits(n_local_size) - 1
97     residue = (max_error / Epsilon) / logm
98     if (residue .lt. max_residue .and. 0 .eq. rank) then
99         write (*,'(A)') "Verification successful"
100     else if (residue .ge. max_residue) then
101         write (*,'(A,ES10.4,A)') "Verification failed (residue=", residue, ")"
102         write (*,'(A,ES11.4)') "Max error:", max_error
103         write (*,'(A,2ES13.4,A,2ES11.4,A)') "In:", c(mei), "; Out:", spare(mei), ")"
104     endif
105 end subroutine
106
107 !-----
108 ! fft: compute the fft of complex c. c is of length n_local_size

```

```

109      ! ATTN: using local view
110      !-----
111      subroutine fft_inner(n_local_size, direction)
112          implicit none
113          integer(8)          :: n_local_size, direction, world_size, n_world_size
114          integer(8)          :: rank, lo, hi, lstride, i, j, k
115          integer(8)          :: levels, l, loc_comm, m, m2
116          complex(8)         :: ce, cr, cl
117          double precision    :: two_pi, angle_base
118
119      !-----
120      ! prep for computation
121      !-----
122      rank = team_rank()
123      world_size = team_size()
124      n_world_size = world_size * n_local_size
125      two_pi = 2.0d0 * acos(-1.0d0) * direction
126
127      ! we assume the number of procs and the problem size are power of two
128      levels = number_of_bits(n_world_size) - 1
129      loc_comm = number_of_bits(n_local_size)
130
131      ! Permute local data
132      call permute(c, n_world_size, spare)
133
134      ! Generate twiddle factor table. We only use the first half, since the
135      ! second half is just -1 times the first.
136      do i = 0, n_local_size/2 - 1
137          spare(i) = exp(cmplx(0.0,(i) * ((-two_pi)/n_local_size)))
138      enddo
139
140      !-----
141      ! phase 1 computation
142      !-----
143      do l = 1, loc_comm-1          ! --- for each local level in the FFT
144          m = ishft(1, l)
145          m2 = ishft(m, -1)
146          lstride = ishft(n_local_size, -1)
147          twiddles(0:m2 - 1) = spare(0:n_local_size/2-1:lstride)
148          do k = 0, n_local_size-1, m          ! --- for each butterfly in a level
149              do j = k, k + m2 - 1          ! --- for each point in a half bfly
150                  ce = twiddles(j - k)
151                  cr = ce * c(j + m2)
152                  cl = c(j)
153                  c(j) = cl + cr
154                  c(j + m2) = cl - cr
155              end do
156          end do
157      enddo
158
159      ! Get ready for the second phase -- transpose to cyclic layout.
160      ! We can't fit in memory a twiddle table for this phase, alas.
161      call transpose(c, n_world_size, world_size, spare, .true.)
162
163      !-----
164      ! phase 2 computation
165      !-----
166      do l = loc_comm, levels          ! --- for each level in the fft
167          m = ishft(1, l) / world_size          ! --- local elts/bfly
168          m2 = ishft(m, -1)
169          angle_base = (-two_pi) / real(ishft(1,l))
170          do k = 0, n_local_size-1, m          ! --- for each butterfly in a level
171              do j = k, k + m2 - 1          ! --- for each point in a half bfly
172                  ce = exp(cmplx(0.0,((j - k) * world_size + rank) * angle_base))
173                  cr = ce * c(j + m2)
174                  cl = c(j)
175                  c(j) = cl + cr
176                  c(j + m2) = cl - cr
177              end do
178          end do
179      enddo
180
181      ! restore to original layout
182      call transpose(c, n_world_size, n_local_size/world_size, spare, .false.)
183      end subroutine fft_inner
184
185      end module

```

C.3 module_bits.caf

```
1 module module_bits
2 contains
3
4 !-----
5 ! number_of_bits: compute the index of the leftmost non-zero bit in i
6 !-----
7 function number_of_bits(i) result(num)
8 implicit none
9 integer(8) :: i, itmp, num
10 num = 0
11 itmp = i
12 do while (itmp .gt. 0)
13 num = num + 1
14 itmp = ishft(itmp, -1)
15 end do
16 end function number_of_bits
17
18 end module module_bits
```

C.4 module_permute.caf

```
1 module module_permute
2
3 use module_bits
4 use module_reverse
5
6 contains
7
8 subroutine packf(input, output, n, p, n_b, cp_b, npadding, do_bitreverse)
9 implicit none
10
11 ! dummy arguments
12 integer(8) :: n, n_b, p, npadding, p_b, cp_b
13 complex(8) :: output(0:n-1 + p * npadding)
14 complex(8) :: input(0:n-1)
15 logical :: do_bitreverse
16
17 ! local variables
18 integer(8) :: pe_bufstart(0:p-1)
19 integer(8) :: pe_bufflen, buf, p_bits, ii, i, jj, j, ooffset, ioffset
20
21 !-----
22 ! copy the constant from the caller into local variable
23 ! without the copy, it will have runtime-crash on gfortran44
24 !-----
25 p_b = cp_b
26
27 !-----
28 ! compute a vector that will enable us to use the low log_2(p) bits
29 ! of an index to find the start of the buffer for a processor based
30 ! on the bitreverse of those bits
31 !-----
32 p_bits = number_of_bits(p-1)
33 pe_bufflen = n/p + npadding
34 if (do_bitreverse) then
35 do i = 0, p - 1
36 pe_bufstart(i) = i_bitreverse(i, p_bits) * pe_bufflen
37 enddo
38 else
39 do i = 0, p - 1
40 pe_bufstart(i) = i * pe_bufflen
41 enddo
42 endif
43
44 !print *, "n=", n, "n_b = ", n_b, "p =", p, "p_b =", p_b, "npadding=", npadding
45
46 ! special case handling for small numbers of processors
47 if (p_b .gt. p) then; p_b = p; endif
48
49 do jj = 0, p-1, p_b
50 ooffset = 0
51 do ii = 0, n-1, p * n_b
52 do j = jj, jj + p_b - 1
53 buf = pe_bufstart(j)
```

```

54         ioffset = ooffset
55         do i = ii, min(ii + p * (n_b - 1), n-1), p
56             output(buf + ioffset) = input(i+j)
57             ioffset = ioffset + 1
58         enddo
59     enddo
60     ooffset = ooffset + n_b
61 enddo
62 enddo
63 end subroutine packf
64
65
66 subroutine permute_locally(dest, src, n, cn_b)
67     implicit none
68
69     complex(8), dimension(0:) :: src, dest
70     integer(8) :: n_b, cn_b
71     integer(8) :: i, j, n, n_bits
72     n_b = cn_b
73
74     n_bits = number_of_bits(n - 1)
75
76     if (n_b .gt. n) then; n_b = n; endif
77
78     do j = 0, n_b
79         do i = j, n - 1, n_b
80             dest(i_bitreverse(i, n_bits)) = src(i)
81         end do
82     end do
83 end subroutine permute_locally
84
85 subroutine permute(c, n, scratch)
86     implicit none
87
88     ! parameters
89     integer(8), parameter :: SIZEOF_COMPLEX = 16
90     complex(8), dimension(0:) :: c[*]
91     complex(8), dimension(0:) :: scratch[*]
92     integer(8) :: n, world_size, local_n, block_size
93
94     world_size = team_size()
95     local_n = n / world_size
96     block_size = local_n / world_size
97
98     call packf(c, scratch, local_n, world_size, 32_8, 1024_8, 0_8, .true.)
99     call team_alltoall(scratch, c, block_size)
100    call permute_locally(c, scratch, local_n, ishft(1_8,22))
101
102 end subroutine permute
103
104
105 end module module_permute

```

C.5 module_reverse.caf

```

1  module module_reverse
2
3      integer(8) :: mask64, mask32, mask16, mask8, mask4, mask2, mask1
4
5  contains
6
7  !-----
8  ! i_bitreverse: return the bitreverse of the n-bit integer i
9  !-----
10     function i_bitreverse(i, n)
11
12         integer(8) :: i_bitreverse, i, n
13         integer(8) :: ival, itmp, imask, ishft
14
15         interchange(ival, imask, ishft) = &
16             ior(ishft(iand(ival, imask), -ishft), &
17                 ishft(iand(ival, not(imask)), ishft))
18         itmp = interchange(i, mask32, 32_8)
19         itmp = interchange(itmp, mask16, 16_8)
20         itmp = interchange(itmp, mask8, 8_8)
21         itmp = interchange(itmp, mask4, 4_8)

```

```

22     itmp = interchange(itmp, mask2, 2_8)
23     itmp = interchange(itmp, mask1, 1_8)
24
25     i_bitreverse = ishft(itmp, n - 64_8) ! provide result as n-bit value
26 end function i_bitreverse
27
28 subroutine bitinit
29     mask64 = not(0_8)
30     mask32 = ishft(mask64,32_8)
31     mask16 = ieor(ishft(mask32,-16_8),mask32)
32     mask8 = ieor(ishft(mask16,-8_8),mask16)
33     mask4 = ieor(ishft(mask8,-4_8),mask8)
34     mask2 = ieor(ishft(mask4,-2_8),mask4)
35     mask1 = ieor(ishft(mask2,-1_8),mask2)
36 end subroutine bitinit
37
38
39 end module module_reverse

```

D Global HPL

D.1 hpl.caf

```

1     program HPL
2     use support
3     use HPLmod
4
5     integer :: m, n, i, k
6     integer, allocatable :: seed(:)
7     double precision :: start_time, end_time, cputime, gflops, normA, normx, resid, norm_r, norm_c
8     double precision, allocatable :: normb(:)[*]
9     double precision, pointer :: a(:, :)
10    double precision, pointer :: b(:)
11    character(len=32) :: arg
12
13    ! compute my location in a 2D processor grid
14    nprocs = team_size()
15    me = team_rank()
16    mycol = me / NPROW
17    myrow = me - mycol * NPROW
18
19    BLKSIZE = min(PROBLEMSIZE / NPCOL, 16)
20    if (0 .lt. command_argument_count()) then
21        call get_command_argument(1, arg)
22        if (len_trim(arg) .ne. 0) then
23            read(arg, '(I10)') BLKSIZE
24        end if
25    end if
26
27    if (nprocs .ne. NPCOL * NPROW) then
28        if (me .eq. 0) print *, 'execution_needs_', NPCOL * NPROW, 'processors'
29        stop
30    end if
31
32    i = NPROW
33    rowp2 = 1
34    rowlog2 = 0
35    do while (i .gt. 1)
36        i = i / 2
37        rowp2 = rowp2 * 2
38        rowlog2 = rowlog2 + 1
39    end do
40
41    ! compute my part of the matrix and allocate
42    m = localsize(PROBLEMSIZE, 0, BLKSIZE, NPROW, myrow)
43    n = localsize(PROBLEMSIZE, 0, BLKSIZE, NPCOL, mycol)
44    allocate(ab(m,n+1)[], x(n)[])
45    a => ab(1:m,1:n)
46    b => ab(1:m,n+1)
47
48    call team_split(team_world, myrow, mycol, rteam, myrow, ierr)
49    call team_split(team_world, myrow, mycol, rteam_, myrow, ierr)
50    call team_split(team_world, mycol, myrow, cteam, mycol, ierr)
51
52    call random_seed(size = k)

```

```

53 allocate(seed(k))
54 call random_seed(get = seed)
55 call init(m, n)
56
57 ! collect time for the main body of work
58 call barrier()
59 call cpu_time(start_time)
60 call lup(m, n)
61 call backsolve(m, n)
62 call barrier()
63 call cpu_time(end_time)
64 cputime = end_time - start_time
65
66 gflops = (2.0/3.0*PROBLEMSIZE + 1.5) * (PROBLEMSIZE/1.0e9) * (PROBLEMSIZE/cputime)
67 if (me .eq. 0) then
68   print *, 'size,nprocs,blksize,time,gflops:',PROBLEMSIZE,',',nprocs,',',BLKSIZE,',',cputime,',',gflops
69 end if
70
71 ! verify results
72 call random_seed(put = seed)
73 call init(m,n)
74 normA = norm_r(a(1,1),m,n,0)
75 normx = norm_c(x(1),n)
76 allocate(normb(1)[])
77 if (mycol .eq. mod(PROBLEMSIZE/BLKSIZE, NPCOL)) then
78   normb(1) = maxval(abs(b))
79   call team_reduce(normb, normb, 0, REDUCE_MAX, cteam)
80   call team_broadcast(normb, 0, cteam)
81   if (n .gt. 0 .and. m .gt. 0) b = b - matmul(a,x)
82 else
83   b = 0.0
84   if (n .gt. 0 .and. m .gt. 0) b = - matmul(a,x)
85 end if
86 call team_broadcast(normb, mod(PROBLEMSIZE/BLKSIZE, NPCOL), rteam)
87 call team_reduce(b(1:m), b(1:m), 0, REDUCE_SUM, rteam)
88 resid = norm_r(b(1),m,1,1)
89 resid = resid / (epsilon * (normA * normx + normb(1)) * PROBLEMSIZE)
90 if (me .eq. 0 .and. resid .lt. THRESHOLD) print *, "result_VALID_scaled_residual:", resid
91 if (me .eq. 0 .and. resid .ge. THRESHOLD) print *, "result_INVALID_scaled_residual:", resid
92
93 deallocate(ab,x,seed,normb)
94 end program HPL
95
96 function norm_r(a,s,t,is_vect) result(r)
97 use support
98 use HPLmod
99 integer :: s,t,i,is_vect
100 double precision :: r, a(s,t)
101 double precision, allocatable :: w(:)[*], nval(:)[*]
102 allocate (w(s)[], nval(1)[])
103 w = 0.0
104 do i = 1, s
105   w(i) = sum(abs(a(i,:)))
106 end do
107 if (is_vect .ne. 1) call team_reduce(w(1:s), w(1:s), 0, REDUCE_SUM, rteam)
108 if (mycol .eq. 0) then
109   nval(1) = 0.0
110   if (s .gt. 0) nval(1) = maxval(w)
111   call team_select(nval, nval, 0, SELECT_MAX, cteam)
112 end if
113 call team_broadcast(nval, 0)
114 r = nval(1)
115 deallocate(w, nval)
116 end function norm_r
117
118 function norm_c(v,n) result(r)
119 use support
120 use HPLmod
121 double precision :: r, v(n), w(n)
122 double precision, allocatable :: nval(:)[*]
123 allocate (nval(1)[])
124 nval(1) = 0.0
125 if (n .gt. 0) then
126   w = abs(v)
127   nval(1) = maxval(w)
128 end if
129 call team_select(nval, nval, 0, SELECT_MAX, rteam)

```

```

130   call team_broadcast(nval, 0)
131   r = nval(1)
132   deallocate(nval)
133   end function norm_c

```

D.2 module_hpl.caf

```

1   module HPLmod
2   integer, parameter :: PROBLEMSIZE = 64*12*1024, NUMPANELS = 2, NPCOL = 64, NPROW = 64
3   double precision, parameter :: EPSILON = 1.11e-16, THRESHOLD = 16.0
4   integer :: nprocs, rowp2, rowlog2, me, mycol, myrow, BLKSIZE
5   team :: rteam, rteam_, cteam, subcteam
6
7   ! panelinfo: i, j, iloc, jloc, mloc, nloc, rproc, cproc, permute
8   integer, target, allocatable :: panelinfo_1(:)[*], panelinfo_2(:)[*]
9   double precision, target, allocatable :: panelbuff_1(:)[*], panelbuff_2(:)[*]
10  double precision, target, allocatable :: ab(:,:)[*]
11  double precision, allocatable :: x(:)[*]
12  event, allocatable, dimension(:) :: delivered[*]
13  integer :: isz = 0, bsz = 0
14
15  type :: panelptr
16     integer, pointer :: info(:)
17     double precision, pointer :: buff(:)
18  end type
19
20  type (panelptr) :: panels(1:NUMPANELS)
21
22  end module HPLmod

```

D.3 module_init.caf

```

1   subroutine init(m, n)
2   use HPLmod
3   integer :: m, n
4   do i = 0, NPROW * NPCOL - 1 - me
5     do j = 1, n+1
6       call random_number(ab(:,j))
7     end do
8   end do
9   do j = 1, n+1
10    ab(:,j) = ab(:,j) * 2.0 - 1.0
11  end do
12  end subroutine init

```

D.4 module_panel.caf

```

1   ! initialize a panel starting at global index (pp, pp)
2   subroutine initpanel(p, pp)
3   use HPLmod
4   use support
5   integer :: p, pp, blk, rproc, cproc, mloc, nloc, iloc, jloc      ! pp is 0-based
6
7   blk = pp / BLKSIZE
8   rproc = blk - blk / NPROW * NPROW
9   cproc = blk - blk / NPCOL * NPCOL
10  mloc = localsize(PROBLEMSIZE - pp, pp, BLKSIZE, NPROW, myrow)
11  nloc = localsize(PROBLEMSIZE + 1 - pp, pp, BLKSIZE, NPCOL, mycol)
12  isz = 8 * BLKSIZE + 2 * NPROW + 10
13  bsz = (mloc + nloc + BLKSIZE + 1) * BLKSIZE
14  if (p .eq. 1) then
15    allocate(panelinfo_1(isz)[], panelbuff_1(bsz)[])
16    panels(1)%info(1:isz) => panelinfo_1
17    panels(1)%buff(1:bsz) => panelbuff_1
18  else
19    allocate(panelinfo_2(isz)[], panelbuff_2(bsz)[])
20    panels(2)%info(1:isz) => panelinfo_2
21    panels(2)%buff(1:bsz) => panelbuff_2
22  end if
23  iloc = localsize(pp, 0, BLKSIZE, NPROW, myrow)
24  jloc = localsize(pp, 0, BLKSIZE, NPCOL, mycol)
25  panels(p)%info(1:9) = (/pp+1, pp+1, iloc+1, jloc+1, mloc, nloc, rproc, cproc, 0/) ! 0 to 1-based
26  end subroutine initpanel

```

D.5 module_lup.caf

```
1      subroutine lup(m, n)
2      use HPLmod
3      use support
4      integer :: m, n, ni, nn, p, pp = 0, cp = 1, pproc, cproc, ub, colstart
5      ni = localsize(PROBLEMSIZE+1, 0, BLKSIZE, NPCOL, mycol)
6      allocate(delivered(1:NUMPANELS)[@team_world])
7      call event_init(delivered, NUMPANELS)
8
9      nn = ni
10     ! build panels to fill the panel buffer, factorize and update them
11     call initpanel(1, pp)
12     pp = pp + BLKSIZE
13     if (mycol .eq. panels(1)%info(8)) nn = nn - BLKSIZE
14     call initpanel(2, pp)
15
16     if (mycol .eq. panels(1)%info(8)) call fact(m, n, 1)
17     mloc = panels(1)%info(5)
18     pproc = panels(1)%info(8)
19     call team_broadcast_async(panelbuff_1(1:(mloc+BLKSIZE+1)*BLKSIZE), panels(1)%info(8), delivered(1), rteam)
20
21     ! factorize rest panels and finish all updates
22     cproc = panels(2)%info(8)
23
24     do j = pp, PROBLEMSIZE - 1, BLKSIZE
25
26         cp = j / BLKSIZE + 1
27         cp = mod(cp - 1, 2) + 1
28         if (cp .eq. 1) deallocate(panelinfo_1, panelbuff_1)
29         if (cp .eq. 2) deallocate(panelinfo_2, panelbuff_2)
30         call initpanel(cp, j)
31         call event_wait(delivered(3-cp))
32         numcol = 0
33         colstart = 0
34
35         if (mycol .eq. cproc) then
36             numcol = localsize(min(BLKSIZE, PROBLEMSIZE-j), j, BLKSIZE, NPCOL, mycol)
37             if (numcol .gt. 0) then
38                 if (NPCOL .eq. 1) call update(m, n, BLKSIZE, numcol, 3 - cp)
39                 if (NPCOL .ne. 1) call update(m, n, 0, numcol, 3 - cp)
40             end if
41             call fact(m, n, cp)
42             colstart = colstart + numcol
43         end if
44         if (mycol .eq. pproc) colstart = colstart + BLKSIZE
45         ub = (panels(cp)%info(5)+BLKSIZE+1)*BLKSIZE
46         if (cp .eq. 1) call team_broadcast_async(panelbuff_1(1:ub), panelinfo_1(8), delivered(1), rteam)
47         if (cp .eq. 2) call team_broadcast_async(panelbuff_2(1:ub), panelinfo_2(8), delivered(2), rteam_)
48
49         if (nn-numcol .gt. 0) call update(m,n, colstart, nn-numcol, 3 - cp)
50
51         if (mycol .eq. cproc) nn = nn - BLKSIZE
52         pproc = cproc
53         cproc = mod(cproc+1, NPCOL)
54
55     end do
56     numcol = localsize(1, PROBLEMSIZE, BLKSIZE, NPCOL, mycol)
57     colstart = ni - panels(cp)%info(4)
58     if (numcol .gt. 0) call update(m, n, colstart, numcol, cp)
59
60     deallocate(panelinfo_1, panelinfo_2, panelbuff_1, panelbuff_2)
61
62     end subroutine lup
```

D.6 module_fact.caf

```
1      ! perform factorization of a panel
2      subroutine fact(m, n, p)
3      use support
4      use HPLmod
5      double precision, pointer, DIMENSION(:,:) :: a
6      double precision, pointer, DIMENSION(:) :: b, c, piv, l
7      integer :: m, n, p, clb, rlb, rub, rproc, pivl, pivg, pivproc, mloc, numcols, pivltmp(1)
8      double precision, allocatable :: w(:)[*]
9
```



```

10 with team cteam
11 allocate(w(2*BLKSIZE+4)[])
12
13 mloc = panels(p)%info(5)
14 a => ab(1:m,1:n)
15 b => ab(1:m,n+1)
16 l => panels(p)%buff(1:mloc*BLKSIZE)
17 c => panels(p)%buff(mloc*BLKSIZE+1:(mloc+BLKSIZE)*BLKSIZE)
18 piv => panels(p)%buff((mloc+BLKSIZE)*BLKSIZE+1:(mloc+BLKSIZE+1)*BLKSIZE)
19 clb = panels(p)%info(4)
20 rlb = panels(p)%info(3)
21 rub = rlb + mloc - 1
22 rproc = panels(p)%info(7)
23 numcols = min(BLKSIZE, n - clb + 1)
24 do i = 1, numcols
25     ! compute pivots, switch rows, and update the panel
26     w(1:BLKSIZE+3) = 0.0
27     w(BLKSIZE+4) = nrow
28     if (rub .ge. rlb) then
29         pivltmp = maxloc(abs(a(rlb:rub,clb+i-1)))
30         pivl = pivltmp(1) ! 1-based
31         pivg = pivl + rlb - 2 ! 0-based
32         pivg = pivg + BLKSIZE * ((NPROW - 1) * (pivg / BLKSIZE) + myrow) + 1 ! 1-based
33         w(BLKSIZE+1) = a(rlb+pivl-1,clb+i-1)
34         w(BLKSIZE+2:BLKSIZE+4) = (/pivl, pivg, myrow/)
35         w(1:numcols) = a(rlb+pivl-1, clb:clb+numcols-1)
36     end if
37     call team_allselect(w(1:BLKSIZE+4), w(1:BLKSIZE+4), SELECT_MAX, cteam)
38     piv(i) = w(BLKSIZE+3)
39
40     pivl = w(BLKSIZE+2)
41     pivproc = w(BLKSIZE+4)
42     if (myrow .eq. rproc .and. myrow .ne. pivproc) then
43         w(BLKSIZE+5:BLKSIZE+numcols+4)[pivproc@cteam] = a(rlb,clb:clb+numcols-1)
44     end if
45     call barrier(cteam)
46
47     c(i::BLKSIZE) = w(1:numcols)
48     if (myrow .eq. rproc) then
49         if (myrow .ne. pivproc) then
50             a(rlb, clb:clb+numcols-1) = w(1:numcols)
51         else if (myrow .eq. pivproc .and. pivl .ne. 1) then
52             a(pivl+rlb-1, clb:clb+numcols-1) = a(rlb,clb:clb+numcols-1)
53             a(rlb, clb:clb+numcols-1) = w(1:numcols)
54         end if
55     else if (myrow .eq. pivproc) then
56         a(pivl+rlb-1, clb:clb+numcols-1) = w(BLKSIZE+5:BLKSIZE+numcols+4)
57     end if
58
59     if (myrow .eq. rproc) rlb = rlb + 1
60     a(rlb:rub,clb+i-1) = a(rlb:rub,clb+i-1) / w(BLKSIZE+1)
61     if (i < BLKSIZE) call dger(rub-rlb+1,numcols-i,-1.0d0,a(rlb,clb+i-1),1,w(i+1),1,a(rlb,clb+i),m)
62
63 end do
64
65 rlb = panels(p)%info(3)
66 do i = clb, clb + numcols - 1
67     l((i-clb)*mloc+1:(i-clb+1)*mloc) = a(rlb:rlb+mloc-1,i)
68 end do
69 deallocate(w)
70
71 end with team
72 end subroutine fact

```

D.7 module_update.caf

```

1 ! use currpanel to update numcol columns of the trailing matrix of a
2 subroutine update(m, n, coldiff, numcol, p)
3 use HPLmod
4 use support
5 integer :: m, n, coldiff, numcol, lb, ub, pos, ierr, k, jlb, mykey, iu, ipiv, nn, lb_p, tsize
6 double precision, pointer, DIMENSION(:,:) :: a
7 double precision, pointer, DIMENSION(:) :: c, piv, l, u
8 integer, pointer, DIMENSION(:) :: pairs, aloc, uloc, len1, len2, permute
9 integer, pointer :: flag, numpairs
10 double precision, allocatable :: w(:,:)[*]

```

```

11 integer, allocatable :: wptr(:)[*]
12 integer :: i, j, p, ii, mloc, nloc, iloc, jloc, mloc_p, src, dst, srcloc, dstloc, numRows
13 integer :: fnds, fndd, rproc, srcprow, dstprow, wrows, wlb, mydist, partner, lessip2
14
15 with team cteam
16
17 iloc = panels(p)%info(3)
18 jloc = panels(p)%info(4)
19 mloc = panels(p)%info(5)
20 nloc = panels(p)%info(6)
21 rproc = panels(p)%info(7)
22 a => ab(1:m,1:n)
23 l => panels(p)%buff(1:mloc*BLKSIZE)
24 c => panels(p)%buff(mloc*BLKSIZE+1:(mloc+BLKSIZE)*BLKSIZE)
25 piv => panels(p)%buff((mloc+BLKSIZE)*BLKSIZE+1:(mloc+BLKSIZE+1)*BLKSIZE)
26 u => panels(p)%buff((mloc+BLKSIZE+1)*BLKSIZE+1:(mloc+nloc+BLKSIZE+1)*BLKSIZE)
27 permute => panels(p)%info(9:8*BLKSIZE+2*NPROW+10)
28 nn = min(numcol, PROBLEMSIZE - panels(p)%info(2))
29 allocate(w(nn+1, BLKSIZE)[], wptr(1)[])
30 flag => permute(1)
31 numpairs => permute(2)
32 pairs => permute(3:4*BLKSIZE+2)
33 aloc => permute(4*BLKSIZE+3:6*BLKSIZE+2)
34 uloc => permute(6*BLKSIZE+3:8*BLKSIZE+2)
35 len1 => permute(8*BLKSIZE+3:8*BLKSIZE+NPROW+2)
36 len2 => permute(8*BLKSIZE+NPROW+3:8*BLKSIZE+2*NPROW+2)
37
38 ! compute row permutation from pivot information
39 if (flag .eq. 0) then
40     ipiv = piv(1)
41     pairs(1:2) = (/ipiv, panels(p)%info(1)/)
42     numpairs = 1
43     if (ipiv .ne. panels(p)%info(1)) then
44         pairs(3:4) = (/panels(p)%info(1), ipiv/)
45         numpairs = 2
46     end if
47     do i = 2, min(BLKSIZE, PROBLEMSIZE-panels(p)%info(2)+1)
48         src = panels(p)%info(1) + i - 1
49         dst = piv(i)
50         fnds = 0
51         fndd = 0
52         do j = 1, numpairs
53             if (fnds .eq. 0 .or. src .ne. dst .and. fndd .eq. 0) then
54                 if (src .eq. pairs((j-1)*2+2)) fnds = j
55                 if (src .ne. dst .and. dst .eq. pairs((j-1)*2+2)) fndd = j
56             end if
57         end do
58
59         if (fnds .eq. 0) then
60             pairs(numpairs*2+1:numpairs*2+2) = (/src, dst/)
61             numpairs = numpairs + 1
62             pos = numpairs
63         else
64             pairs((fnds-1)*2+2) = dst
65             pos = fnds
66         end if
67
68         if (src .ne. dst) then
69             if (fndd .eq. 0) then
70                 pairs(numpairs*2+1:numpairs*2+2) = (/dst, src/)
71                 numpairs = numpairs + 1
72                 pos = numpairs
73             else
74                 pairs((fndd-1)*2+2) = src
75                 pos = fndd
76             end if
77         end if
78         if (i .ne. pos) then
79             src = pairs((i-1)*2+1)
80             dst = pairs((i-1)*2+2)
81             pairs((i-1)*2+1:(i-1)*2+2) = pairs((pos-1)*2+1:(pos-1)*2+2)
82             pairs((pos-1)*2+1:(pos-1)*2+2) = (/src, dst/)
83         end if
84     end do
85
86     ! compute the locations in a and u where rows should be placed
87     ii = 1

```

```

88     len2 = 0
89     do i = 1, numpairs
90         src = pairs((i-1)*2+1)
91         srcprow = (src - 1) / BLKSIZE
92         srcprow = mod(srcprow, NPROW)
93         len2(srcprow+1) = len2(srcprow+1) + 1      ! 1-based, procid also 1-based
94         if (myrow .eq. srcprow) then
95             srcloc = localsize(src, 0, BLKSIZE, NPROW, myrow)
96             aloc(ii) = srcloc - iloc
97             dst = pairs((i-1)*2+2)
98             if (myrow .eq. rproc) then
99                 dstprow = (dst - 1) / BLKSIZE
100                dstprow = mod(dstprow, NPROW)
101                if (dstprow .eq. rproc) then
102                    uloc(ii) = dst - panels(p)%info(1)
103                    if (dst - panels(p)%info(1) .ge. BLKSIZE) then
104                        dstloc = localsize(dst, 0, BLKSIZE, NPROW, myrow)
105                        uloc(ii) = iloc - dstloc
106                    end if
107                else
108                    fndd = 0
109                    do j = 1, numpairs
110                        if (fndd .eq. 0 .and. dst .eq. pairs((j-1)*2+1)) fndd = j
111                    end do
112                    uloc(ii) = pairs((fndd-1)*2+2) - panels(p)%info(1)
113                end if
114            else
115                uloc(ii) = dst - panels(p)%info(1)
116            end if
117            ii = ii + 1
118        end if
119    end do
120    flag = 1
121 end if

122
123 ! permute rows and broadcast u
124 w = 0
125 numrows = len2(myrow+1)
126 len1(:) = len2(:)
127 if (myrow .eq. rproc) then
128     do j = 0, nn - 1
129         do i = 1, numrows
130             if (uloc(i) .ge. 0) u(uloc(i)+j*BLKSIZE+1) = a(aloc(i)+iloc, jloc+coldiff+j)
131             if (uloc(i) .lt. 0) a((-uloc(i))+iloc, jloc+coldiff+j) = a(aloc(i)+iloc, jloc+coldiff+j)
132         end do
133     end do
134 else
135     w(1,1:numrows) = uloc(1:numrows)
136     do k = 0, nn - 1
137         do i = 1, numrows
138             w(k+2,i) = a(aloc(i)+iloc, jloc+coldiff+k)
139         end do
140     end do
141 end if
142 if (myrow .eq. rproc) then
143     len1(myrow+1) = 0
144     numrows = 0
145 end if

146
147 wptr(1) = numrows
148 mydist = myrow - rproc
149 if (mydist < 0) mydist = mydist + NPROW
150 partner = ieor(mydist, rowp2)
151 call barrier()
152 if (NPROW - rowp2 .ne. 0 .and. partner .lt. NPROW) then
153     partner = mod(partner + rproc, NPROW)
154     if (mydist .eq. 0) then
155         lb = (mloc+BLKSIZE+1)*BLKSIZE+1
156         ub = (mloc+nloc+BLKSIZE+1)*BLKSIZE
157         if (p .eq. 1) then
158             mloc_p = panelinfo_1(5)[partner]
159             lb_p = (mloc_p+BLKSIZE+1)*BLKSIZE+1
160             panelbuff_1(lb_p : lb_p+ub-lb)[partner] = panelbuff_1(lb : ub)
161         else
162             mloc_p = panelinfo_2(5)[partner]
163             lb_p = (mloc_p+BLKSIZE+1)*BLKSIZE+1
164             panelbuff_2(lb_p : lb_p+ub-lb)[partner] = panelbuff_2(lb : ub)

```

```

165     end if
166   else if (mydist .eq. rowp2) then
167     w(:, 1 : len1(myrow+1))[partner] = w(:, 1 : len1(myrow+1))
168   else if (iand(mydist, rowp2) .ne. 0) then
169     wlb = wptr(1)[partner]
170     w(:, wlb+1 : wlb+len1(myrow+1))[partner] = w(:, 1 : len1(myrow+1))
171   else if (len1(partner+1) .gt. 0) then
172     wptr(1) = wptr(1) + len1(partner+1)
173   end if
174 end if
175
176 call barrier()
177 partner = ieor(mydist, rowp2)
178 if (NPROW - rowp2 .ne. 0 .and. partner .lt. NPROW) then
179   partner = mod(partner + rproc, NPROW)
180   if (mydist .eq. 0 .and. len1(partner+1) .gt. 0) then
181     do i = 1, len1(partner+1)
182       iu = w(1,i) + 1 ! w(1,i) is 0-based
183       u(iu:iu+(nn-1)*BLKSIZE:BLKSIZE) = w(2:nn+1,i)
184     end do
185   end if
186 end if
187
188 do i = 1, NPROW-rowp2-1
189   lessip2 = mod(rproc+i, NPROW)
190   len1(lessip2+1) = len1(lessip2+1) + len1(mod(lessip2+rowp2, NPROW)+1)
191 end do
192
193 mykey = 0
194 wrows = wptr(1)
195 if (mydist .lt. rowp2) mykey = 1
196 if (mydist .ge. rowp2 .and. NPROW - rowp2 .ne. 0) mykey = 2
197 call team_split(cteam, mykey, myrow, subcteam, mykey, ierr)
198
199 with team subcteam
200
201 if (mydist .lt. rowp2) then
202   do i = 1, rowlog2
203     partner = ieor(mydist, ibset(0, i-1))
204     partner = mod(rproc + partner, NPROW)
205     if (mydist .lt. ibset(0, i) .and. mydist .lt. ibset(0, i-1)) then
206       lb = (mloc+BLKSIZE+1)*BLKSIZE+1
207       ub = (mloc+nloc+BLKSIZE+1)*BLKSIZE
208       if (p .eq. 1) then
209         mloc_p = panelinfo_1(5)[partner]
210         lb_p = (mloc_p+BLKSIZE+1)*BLKSIZE+1
211         panelbuff_1(lb_p : lb_p+ub-lb)[partner] = panelbuff_1(lb : ub)
212       else
213         mloc_p = panelinfo_2(5)[partner]
214         lb_p = (mloc_p+BLKSIZE+1)*BLKSIZE+1
215         panelbuff_2(lb_p : lb_p+ub-lb)[partner] = panelbuff_2(lb : ub)
216       end if
217       if (len1(partner+1) > 0) w(:,wrows+1:wrows+len1(partner+1)) = w(:,1:len1(partner+1))[partner]
218     else if (mydist .ge. ibset(0, i)) then
219       w(:,wrows+1:wrows+len1(partner+1)) = w(:,1:len1(partner+1))[partner]
220     end if
221     call barrier()
222     if (mydist .lt. ibset(0, i) .and. mydist .lt. ibset(0, i-1)) then
223       do k = 0, nn - 1
224         do j = 1, len1(partner+1)
225           iu = w(1,wrows+j) + 1 ! w(1,i) is 0-based
226           u(iu+k*BLKSIZE) = w(k+2,wrows+j)
227         end do
228       end do
229     end if
230   end if
231   if (mydist .lt. ibset(0, i) .and. mydist .ge. ibset(0, i-1)) then
232     do k = 0, nn - 1
233       do j = 1, numrows
234         a(aoloc(j)+illoc, jloc+coldiff+k) = u(uloc(j)+1+k*BLKSIZE)
235         u(uloc(j)+1+k*BLKSIZE) = w(k+2, j)
236       end do
237     end do
238     do j = numrows + 1, len1(myrow+1)
239       iu = w(1,j) + 1
240       u(iu:iu+(nn-1)*BLKSIZE:BLKSIZE) = w(2:, j)
241     end do

```

```

242     end if
243     if (mydist .ge. ibset(0, i) .or. mydist .lt. ibset(0, i-1)) wrows = wrows + len1(partner+1)
244
245     do j = 0, rowp2 - 1
246         partner = ieor(j, ibset(0, i-1))
247         if (partner .gt. j) then
248             partner = mod(rproc + partner, NPROW)
249             len1(mod(rproc+j, NPROW)+1) = len1(mod(rproc+j, NPROW)+1) + len1(partner+1)
250             len1(partner+1) = len1(mod(rproc+j, NPROW)+1)
251         end if
252     end do
253 end do
254 else if (NPROW - rowp2 .gt. 1) then
255     call team_broadcast(u, 0)
256     do i = 1, numrows
257         a(aloc(i)+iloc, jloc+coldiff:) = u(uloc(i)+1::BLKSIZE)
258     end do
259 end if
260 call team_free(subcteam)
261
262 end with team
263
264 if (NPROW - rowp2 .ne. 0 .and. ieor(mydist, rowp2) .lt. NPROW) then
265     partner = mod(rproc + ieor(mydist, rowp2), NPROW)
266     if (iand(mydist, rowp2) .eq. 0) then
267         lb = (mloc+BLKSIZE+1)*BLKSIZE+1
268         ub = (mloc+nloc+BLKSIZE+1)*BLKSIZE
269         if (p .eq. 1) then
270             mloc_p = panelinfo_1(5)[partner]
271             lb_p = (mloc_p+BLKSIZE+1)*BLKSIZE+1
272             panelbuff_1(lb_p : lb_p+ub-lb)[partner] = panelbuff_1(lb : ub)
273         else
274             mloc_p = panelinfo_2(5)[partner]
275             lb_p = (mloc_p+BLKSIZE+1)*BLKSIZE+1
276             panelbuff_2(lb_p : lb_p+ub-lb)[partner] = panelbuff_2(lb : ub)
277         end if
278     end if
279 end if
280 deallocate(wptr, w)
281
282 call dtrsm('L', 'L', 'N', 'U', BLKSIZE, nn, 1.0d0, c, BLKSIZE, u, BLKSIZE)
283 if (mloc .gt. 0) call module_mmult(l(1), u(1), mloc, nn, iloc, jloc, mloc, coldiff, nn, m, rproc)
284
285 end with team
286
287 end subroutine update

```

D.8 module_solve.caf

```

1     ! solve the upper triangle matrix
2     subroutine backsolve(m, n)
3     use HPLmod
4     use support
5     integer :: i, j, m, n, numblks, pcoln, prown, pcolb, bnxt, bblk, m1, n1, blk, l, pcolp, prrowp
6     double precision, pointer, DIMENSION(:, :) :: a
7     double precision, pointer, DIMENSION(:) :: b
8     double precision, allocatable :: work(:)[*], buf(:)[*]
9     a => ab(1:m, 1:n)
10    b => ab(1:m, n+1)
11    numblks = (PROBLEMSIZE - 1) / BLKSIZE
12    pcoln = mod(numblks, NPCOL)           ! processors own the last column block of A
13    prown = mod(numblks, NPROW)         ! processors own the last row block of A
14    pcolp = pcoln
15    prrowp = prown
16    pcolb = mod(PROBLEMSIZE / BLKSIZE, NPCOL)
17    blk = PROBLEMSIZE - numblks * BLKSIZE
18    l = PROBLEMSIZE - blk
19
20    with team rteam
21        allocate(buf(m)[])
22        if (pcoln .ne. pcolb .and. mycol .eq. pcolb) buf(1:m)[pcoln] = ab(1:m, n+1)
23
24        call barrier()
25        if (pcoln .ne. pcolb .and. mycol .eq. pcoln) b = buf
26
27        if (mycol .ne. pcoln) b = 0

```

```

28
29 m1 = m
30 n1 = n
31 if (myrow .eq. prown .and. mycol .eq. pcoln) then
32   do i = m, m-blk+1, -1
33     b(i) = b(i) / a(i,i-m+n)
34     b(m-blk+1:i-1) = b(m-blk+1:i-1) - b(i) * a(m-blk+1:i-1,i-m+n)
35   end do
36   x(n-blk+1:n) = b(m-blk+1:m)
37 end if
38
39 if (myrow .eq. prown) m1 = m - blk
40 if (mycol .eq. pcoln) n1 = n - blk
41 allocate(work(min((NPCOL-1)*BLKSIZE,m))[])
42
43 do j = 1, numblks
44   pcolp = pcoln - 1
45   prowp = prown - 1
46   if (pcolp .lt. 0) pcolp = pcolp + NPCOL
47   if (prowp .lt. 0) prowp = prowp + NPROW
48   if (NPCOL .ne. 1) bnxt = min(1, (NPCOL-1)*BLKSIZE)
49   if (NPROW .eq. 1) bnxt = min(1, (NPROW-1)*BLKSIZE)
50   bblk = localsize(bnxt, max(1-bnxt,0), BLKSIZE, NPROW, myrow)
51
52   if (mycol .eq. pcoln) then
53     call team_broadcast(x(n1+1:n1+blk), prown, cteam)
54     if (bblk .gt. 0) then
55       b(m1-bblk+1:m1) = b(m1-bblk+1:m1) - matmul(a(m1-bblk+1:m1, n1+1:n1+blk), x(n1+1:n1+blk))
56       work(1:bblk)[pcolp] = b(m1-bblk+1:m1)
57     end if
58   end if
59
60   call barrier()
61
62   if (mycol .eq. pcolp .and. bblk .gt. 0) then
63     b(m1-bblk+1:m1) = b(m1-bblk+1:m1) + work(1:bblk)
64   end if
65
66   if (myrow .eq. prowp .and. mycol .eq. pcolp) then
67     do i = m1, m1-BLKSIZE+1, -1
68       b(i) = b(i) / a(i,i-m1+n1)
69       b(m1-BLKSIZE+1:i-1) = b(m1-BLKSIZE+1:i-1) - b(i) * a(m1-BLKSIZE+1:i-1,i-m1+n1)
70     end do
71     x(n1-BLKSIZE+1:n1) = b(m1-BLKSIZE+1:m1)
72   end if
73
74   if (mycol .eq. pcoln .and. m1-bblk .ge. 1) then
75     b(1:m1-bblk) = b(1:m1-bblk) - matmul(a(1:m1-bblk, n1+1:n1+blk), x(n1+1:n1+blk))
76   end if
77   l = 1 - BLKSIZE
78   pcoln = pcolp
79   prown = prowp
80   blk = BLKSIZE
81   if (myrow .eq. prown) m1 = m1 - blk
82   if (mycol .eq. pcoln) n1 = n1 - blk
83 end do
84
85 if (mycol .eq. pcolp) call team_broadcast(x(1:BLKSIZE), prowp, cteam)
86 deallocate(work, buf)
87 end with team
88
89 end subroutine backsolve

```

D.9 module_mmult.caf

```

1  subroutine module_mmult(l,u,s,t,iloc,jloc,mloc,coldiff,nn,m,rproc)
2  use HPLmod
3  double precision :: l(s,BLKSIZE), u(BLKSIZE,t)
4  double precision, pointer, dimension(:,:) :: a
5  integer s, t, iloc, jloc, mloc, nn, coldiff, rproc, blk
6
7  a => ab(iloc:iloc+s-1, jloc+coldiff:jloc+coldiff+nn-1)
8  if (myrow .ne. rproc) then
9    call dgemm('N','N',s,t,BLKSIZE,-1.0d0,1(1,1),s,u(1,1),BLKSIZE,1.0d0,a(1,1),m)
10 else
11   blk = min(s, BLKSIZE)

```

```

12     if (s .le. blk) then
13         a(1:blk,:) = u(1:blk,:)
14     else
15         call dgemm('N','N',s-blok,t,BLKSIZE,-1.0d0,1(blk+1,1),s,u(1,1),BLKSIZE,1.0d0,a(blk+1,1),m)
16         a(1:blk,:) = u(1:blk,:)
17     end if
18 end if
19
20 end subroutine module_mmult

```

D.10 module_support.caf

```

1  module support
2  contains
3  subroutine reduce_max(in1, in2, out, nb)
4  integer :: nb
5  double precision :: in1, in2, out
6  out = max(in1, in2)
7  end subroutine reduce_max
8
9  subroutine reduce_sum(in1, in2, out, nb)
10 integer :: nb
11 double precision :: in1(nb/8), in2(nb/8), out(nb/8)
12 out = in1 + in2
13 end subroutine reduce_sum
14
15 ! use user-defined type
16 subroutine select_max(in1, in2, out, nb)
17 use HPLmod
18 integer :: nb, B
19 double precision :: in1(nb/8), in2(nb/8), out(nb/8)
20 B = BLKSIZE
21 if ((nb .eq. 8 .and. abs(in1(1)) .lt. abs(in2(1))) .or. &
22     (nb .eq. 8*(B+4) .and. (abs(in1(B+1)) .lt. abs(in2(B+1)) .or. &
23         (abs(in1(B+1)) .eq. abs(in2(B+1)) .and. in1(B+3) .lt. in2(B+3)))))) then
24     out = in2
25 else
26     out = in1
27 end if
28 end subroutine select_max
29
30 ! compute a local portion of the array from global index i (0 based)
31 function localsize(n, i, blk, np, myproc) result(mypart)
32 integer :: n, i, blk, np, myproc, numblks, locblks, mydist, inb, srcproc, inbs, mypart
33 inb = i / blk
34 srcproc = inb - (inb / np) * np
35 inbs = blk - (i - inb * blk)
36 numblks = (n - inbs) / blk + 1
37 locblks = numblks / np
38
39 if (myproc .eq. srcproc) then
40     mypart = n
41     if (n .gt. inbs) then
42         if (numblks .ne. locblks * np) mypart = inbs + locblks * blk
43         if (numblks .eq. locblks * np) mypart = n + (locblks - numblks) * blk
44     end if
45 else
46     mypart = 0
47     if (n .gt. inbs) then
48         mydist = myproc - srcproc
49         if (mydist .lt. 0) mydist = mydist + np
50         if (numblks .lt. np) then
51             if (mydist .lt. numblks) mypart = blk
52             if (mydist .eq. numblks) mypart = n - inbs + blk * (1 - numblks)
53         else
54             mydist = mydist - numblks + locblks * np
55             mypart = locblks * blk
56             if (mydist .lt. 0) mypart = mypart + blk
57             if (mydist .eq. 0) mypart = mypart + n - inbs - numblks * blk + blk
58         end if
59     end if
60 end if
61 end function localsize
62
63 end module support

```

E UTS

E.1 uts.caf

```
1
2  module uts
3      use array_queue
4      implicit none
5
6      integer m !holds the number of children
7      integer*8 total_node_count !total number of nodes processed by each process
8
9      !lifeline threshold holds the minimum work a process has in queue
10     ! before pushing work to others.
11
12     !max_neighbor index controls the number of lifelines; as in
13     !
14     !maxlifeline = 2^(max_neighbor_index -1)
15
16     !gen_mx is the depth
17
18     integer world_size, my_rank, lifeline_threshold, gen_mx, max_neighbor_index
19
20     !steal_trunk_size/push_trunk_size ,generally set to 8 elements, hold the maximum
21     ! number of work items that can be obtained in a steal/push; this is because of the
22     ! AMmediumpacket size limit
23     integer steal_trunk_size, push_trunk_size, max_allowed_unsuccessful_steals, steal_threshold
24
25     character (kind=1, len=20) :: root_descriptor
26     integer yield_count, b_0, seed, log_world_size, last_node_index_work_pushed
27
28     !holds the lifelines; a process toggles a bit on the incoming_lifeline
29     ! to indicate the setup of a lifeline; the bit position indicates the relative position
30     ! (process_rank_on_which_lifeline - 2^bit position) of the process
31     integer(4) :: incoming_lifeline
32
33     integer :: activate_running
34     integer(8) tstart, tend, rate, alarm_count
35     real tsec
36
37     EXTERNAL          rng_spawn
38     EXTERNAL          rng_rand
39     INTEGER           rng_rand
40     EXTERNAL          rng_init
41
42 contains
43
44     !a standard initialization of the random generator
45     !in fortran. attribution:world wide web.
46
47     SUBROUTINE init_random_seed()
48         INTEGER :: i, n, clock
49         INTEGER, DIMENSION(:), ALLOCATABLE :: seed
50         CALL RANDOM_SEED(size = n)
51         ALLOCATE(seed(n))
52         CALL SYSTEM_CLOCK(COUNT=clock)
53         seed = clock + 37 * (my_rank + 1) * (/ (i - 1, i = 1, n) /)
54         CALL RANDOM_SEED(PUT = seed)
55         DEALLOCATE(seed)
56     END SUBROUTINE
57
58     !the kernel of the computation; the queue consists
59     ! of the nodes of the UTS tree whose subtree has to be built.
60     ! in this subroutine, a node is deleted from teh queue,
61     ! the number of children and the children's descriptors
62     ! are determined and added to the queue.
63     ! the children's descriptors are found using the rng_spawn
64     ! from the brg_sha1 library implementation.
65
66     subroutine process_work_item(descriptor, depth)
67         implicit none
68         integer*4, intent(in)::depth
69         integer*4 child_depth
70         character (*) descriptor
71         character (kind=1, len=20) :: child_descriptor
72         integer num_child, child_index
73         logical return_val
```



```

73
74         num_child = get_number_children(descriptor, depth)
75         total_node_count = total_node_count+1
76         child_depth = depth+1
77         do child_index=0,(num_child-1)
78             call rng_spawn(descriptor, child_descriptor, child_index)
79             if(insert_queue(child_descriptor, child_depth) .eq. .false.) then
80                 call exit()
81             endif
82         enddo
83     end subroutine
84
85     function get_lifeline(i) result(lifeline)
86         implicit none
87         integer i, lifeline
88         lifeline = mod(my_rank - (2**(i-1)) + world_size, world_size)
89     end function
90
91     function get_num_set_bits(inlifeline) result(setbitcnt)
92         implicit none
93         integer*4 inlifeline
94         integer setbitcnt, a
95         setbitcnt = 0
96         a = inlifeline
97         do while (a .ne. 0)
98             a = IAND(a, a-1)
99             setbitcnt = setbitcnt + 1
100        enddo
101    end function
102
103     function get_num_items_to_push(inlifeline) result(numitems)
104         implicit none
105         integer numitems
106         integer*4 inlifeline
107         numitems = CEILING((queue_count * 1.0) / (get_num_set_bits(inlifeline)+1.0))
108    end function
109
110     !push_work determines the processes registered on the incoming_lifeline
111     ! and pushed work to them. It also remembers the process for which
112     ! the least amount of work was pushed.
113
114     subroutine push_work()
115         implicit none
116         integer :: i, num_bits_to_beshifted, num_items_to_push, push_flag
117         integer*4 :: copy_inlifeline, copy_inlifeline2
118         logical :: do_not_set_anymore
119         integer :: new_last_node_index_work_pushed
120
121         do_not_set_anymore = .false.
122         i = last_node_index_work_pushed
123         new_last_node_index_work_pushed = last_node_index_work_pushed
124         num_bits_to_beshifted = (i-1)
125         copy_inlifeline2 = incoming_lifeline
126         incoming_lifeline = 0
127         copy_inlifeline = ISHFTC(copy_inlifeline2, -num_bits_to_beshifted, log_world_size)
128         num_items_to_push = get_num_items_to_push(copy_inlifeline)
129
130         do while (copy_inlifeline .ne. 0)
131             push_flag = IAND(1, copy_inlifeline)
132             if (push_flag .eq. 1) call push_work_while_sharing(i, num_items_to_push)
133             copy_inlifeline = ISHFT(copy_inlifeline, -1)
134             i = i + 1
135             if(i .eq. (log_world_size+1)) i=1
136             if ((queue_count) .lt. ((num_items_to_push)/4)) then
137                 if((copy_inlifeline .ne. 0) .and. (do_not_set_anymore .eq. .false.)) then
138                     last_node_index_work_pushed = i
139                     do_not_set_anymore = .true.
140                 endif
141             endif
142         end do
143
144         if((last_node_index_work_pushed .ne. 1) .and. do_not_set_anymore) then
145             copy_inlifeline = 2**(log_world_size) - 1
146             copy_inlifeline = ISHFT(copy_inlifeline, last_node_index_work_pushed-1)
147             copy_inlifeline2 = IAND(copy_inlifeline2, copy_inlifeline)
148             incoming_lifeline = IOR(incoming_lifeline, copy_inlifeline2)
149         endif

```

```

150         if(last_node_index_work_pushed .eq. new_last_node_index_work_pushed) then
151             last_node_index_work_pushed = 1
152         endif
153     end subroutine
154
155     !activate; each process after finishing its quota of work performs
156     !a random steal indicated by the line 188 followed by setting up of
157     !lifelines.
158     !yield_count controls when a process should look at pushing work/
159     !allow stealing of work.
160
161     recursive subroutine activate()
162         character (kind=1, len=20) :: descriptor
163         integer steal_from_img, processed_count, num_items_to_push
164         integer i, j, lifeline, push_flag, neighbor_index, next_neighbor
165         logical return_val
166         integer*4 copy_inlifeline, depth
167
168         activate_running = 1
169         processed_count = 0
170         do while(queue_count .gt. 0)
171             if(delete_queue_end(descriptor, depth) .eq. .false.) then
172                 call exit()
173             endif
174             processed_count = processed_count + 1
175             call process_work_item(descriptor, depth)
176             if (processed_count .eq. yield_count) then
177                 processed_count = 0
178                 call caf_async_advance
179                 if ((incoming_lifeline .ne. 0) .and. (queue_count .ge. lifeline_threshold)) then
180                     call push_work()
181                 endif
182             endif
183         enddo
184
185         activate_running = 0
186         if(world_size .gt. 1) then
187             steal_from_img = get_random_image(my_rank, my_rank)
188             spawn steal_work_spawn(my_rank, 0)[steal_from_img]
189
190             neighbor_index = 0
191             do while (neighbor_index .lt. max_neighbor_index)
192                 next_neighbor = mod(my_rank+(2**neighbor_index), world_size)
193                 spawn set_lifelines(my_rank, neighbor_index)[next_neighbor]
194                 neighbor_index = neighbor_index + 1
195             enddo
196         end if
197     end subroutine
198
199     !pack a set of work items to push to a process.
200
201     subroutine push_work_while_sharing(neighbor_index, num_items_to_push)
202     implicit none
203         integer size, j, neighbor_index, dest, num_items_to_push
204         logical return_val
205         type (queue_element) :: push_work_descriptor(10)
206         character (kind=1, len=20) :: descriptor
207         integer*4 :: lifeline_status, depth
208
209         size = min(queue_count-1, num_items_to_push)
210         size = min(size, push_trunk_size)
211         if(size .gt. 0) then
212             do j=1,size
213                 if(delete_queue_begin(descriptor, depth) .eq. .false.) then
214                     call exit()
215                 endif
216                 push_work_descriptor(j)%desc = descriptor
217                 push_work_descriptor(j)%depth = depth
218             enddo
219             dest = get_lifeline(neighbor_index)
220             spawn copy_item_and_activate(push_work_descriptor, size)[dest]
221         endif
222     end subroutine
223
224     !a successful steal/(work pushed from aprocess) arrive into a queue
225     !via this below function.
226

```

```

227     subroutine copy_item_and_activate(descriptor, num_items)
228         implicit none
229         type(queue_element) :: descriptor(10)
230         integer num_items, i
231         logical return_val
232
233         do i=1,num_items
234             if(insert_queue(descriptor(i)%desc, descriptor(i)%depth) .eq. .false.) then
235                 call exit()
236             endif
237         enddo
238
239         if (activate_running .eq. 0) then
240             activate_running = 1
241             call activate()
242         endif
243     end subroutine
244
245     !apart from setting of lifeline, a process also looks
246     ! at stealing from the node
247
248     subroutine set_lifelines(home, neighbor_index)
249         implicit none
250         integer :: home, neighbor_index, i, size, next_neighbor
251         logical :: return_val
252         integer*4 :: lifeline_status, depth
253         type(queue_element) :: steal_work_descriptor(10)
254         character (kind=1, len=20) :: descriptor
255
256         incoming_lifeline = IOR(incoming_lifeline, ISHFT(1, neighbor_index))
257         if (queue_count .gt. steal_threshold) then
258             size = min(steal_trunk_size, (queue_count/2))
259         do i=1,size
260             if(delete_queue_begin(descriptor, depth) .eq. .false.) then
261                 call exit()
262             endif
263             steal_work_descriptor(i)%desc = descriptor
264             steal_work_descriptor(i)%depth = depth
265         enddo
266         spawn copy_item_and_activate(steal_work_descriptor, size)[home]
267         end if
268     end subroutine
269
270     !this function is spawned on a process to steal work.
271
272     subroutine steal_work_spawn(home, num_unsuccessful_steal_attempt)
273         implicit none
274         integer next_neighbor, home, num_unsuccessful_steal_attempt, steal_from_img
275         character (kind=1, len=20) :: descriptor
276         integer*4 depth
277         type(queue_element) :: steal_work_descriptor(10)
278         integer i, lifeline_img, work_count, size, neighbor_index
279         logical return_val
280
281         if (queue_count .gt. steal_threshold) then
282             size = min(steal_trunk_size, (queue_count/2))
283         do i=1,size
284             if(delete_queue_begin(descriptor, depth) .eq. .false.) then
285                 call exit()
286             endif
287             steal_work_descriptor(i)%desc = descriptor
288             steal_work_descriptor(i)%depth = depth
289         enddo
290         spawn copy_item_and_activate(steal_work_descriptor, size)[home]
291         else
292             num_unsuccessful_steal_attempt = num_unsuccessful_steal_attempt + 1
293         endif
294     end subroutine
295
296     !determines the number of children; the algorithm is followed from
297     ! uts 1.1. library released from ohio state university.
298
299     function get_number_children(descriptor, depth) result(n)
300         implicit none
301         integer*4, intent(in)::depth
302         character (*) descriptor
303         integer v, n

```

```

304         real p_0, d, b_i
305         v = rng_rand(descriptor)
306         if(v .lt. 0) then
307             v = 0
308         endif
309         d = v
310         d = (d/HUGE(v))
311
312         if (depth .lt. gen_mx) then
313             b_i = m
314         else
315             b_i = 0
316         endif
317         p_0 = (1.0 / (1.0 + b_i))
318
319         if( (d .eq. 1) .or. (p_0 .eq. 1)) then
320             n = 0
321         else
322             n = (floor(alog(1-d)/alog(1-p_0)))
323         endif
324     end function
325
326     !used for verification
327
328     subroutine send_my_count(other_cnt)
329         integer*8 other_cnt
330         total_node_count = total_node_count + other_cnt
331     end subroutine
332
333     function get_random_image(avoid_img1, avoid_img2) result(img)
334         integer avoid_img1, avoid_img2, img
335         character (kind=1, len=20) :: descriptor
336         real harvest
337         call RANDOM_NUMBER(harvest)
338         img = mod (INT(harvest * world_size), world_size)
339         if(img .eq. avoid_img1) img = mod(img+2, world_size)
340         if(img .eq. avoid_img2) img = mod(img+2, world_size)
341     end function
342
343 end module
344
345 program s
346     use uts
347     implicit none
348     character (kind=1, len=20):: child_descriptor
349     integer*4 depth, step
350     integer i,j, k, desc_count, img, write_to, work_items_per_image
351     logical return_val
352     type(queue_element) :: push_work_descriptor(10)
353     type(queue_element) :: push_desc(10)
354
355     gen_mx = 18
356     m=4
357     seed = 19
358     total_node_count=0
359     b_0 = m
360
361     my_rank = team_rank()
362     world_size = team_size()
363
364     call init_random_seed()
365     call init_queue(40000)
366
367     log_world_size = (alog(real(world_size))/alog(real(2)))
368     max_allowed_unsuccessful_steals = 1
369     incoming_lifeline = 0
370     last_node_index_work_pushed = 1
371     max_neighbor_index = log_world_size
372
373     yield_count = 64
374     lifeline_threshold = 1*(max_neighbor_index)
375     steal_threshold = 2
376     push_trunk_size = 8
377     steal_trunk_size = 8
378
379     step = 1
380     call team_barrier()

```

```

381     call system_clock(tstart, rate)
382     call rng_init(root_descriptor, seed)
383
384     !process 0 generates some nodes before distributing to others.
385
386     if(my_rank .eq. 0) then
387         do i=0,(b_0-1)
388             call rng_spawn(root_descriptor, child_descriptor, i)
389             return_val = insert_queue(child_descriptor, 1)
390         enddo
391         do while(queue_count .lt. 4*(world_size/step))
392             if(delete_queue_begin(child_descriptor, depth) .eq. .false.) then
393                 call exit()
394             endif
395             call process_work_item(child_descriptor, depth)
396         enddo
397     endif
398
399     work_items_per_image = 4
400     !activate_running = 1
401     call team_barrier()
402
403     !use function shipping to distribute work to others;
404     !after distribution the processes initially dormant start working
405     ! via the call activate present in copy_item_and_activate
406     finish
407         if(my_rank .eq. 0) then
408             do i=1,world_size-1, step
409                 j=1
410                 do while((j .le. work_items_per_image) .and. (queue_count .gt. 0))
411                     return_val = delete_queue_begin(push_desc(j)%desc, push_desc(j)%depth)
412                     j = j+1
413                 enddo
414                 if(j .gt. 1) then
415                     spawn copy_item_and_activate(push_desc, (j-1))[i]
416                 endif
417             enddo
418             call activate()
419         endif
420     end finish
421
422     if (0 .eq. my_rank) then
423         call system_clock(tend, rate)
424         tsec = real(tend-tstart) / rate
425         write (*,*) "Elapsed time: ", tsec
426     endif
427
428     call team_barrier()
429
430     !printing out the work done by each process
431     do i=0,world_size-1
432         if(my_rank .eq. i) then
433             write(*,*) my_rank, "s total node count is: ", total_node_count
434         endif
435         call team_barrier()
436     enddo
437
438     !accumulating the total node count
439     finish
440         if(my_rank .ne. 0) then
441             spawn send_my_count(total_node_count)[0]
442         endif
443     end finish
444
445     if(my_rank .eq. 0) then
446         write(*,*) my_rank, " has the total node count: ", total_node_count
447     endif
448     call team_barrier()
449 end program

```

E.2 queue.caf

```

1     !typical double ended queue implementation
2
3     module type_queue_element
4         type queue_element

```

```

5         sequence
6         !to hold the descriptor of the tree node
7         character (kind=1,len=20) :: desc
8         !to hold the depth of the tree node
9         integer*4 :: depth
10    end type queue_element
11 end module
12 module array_queue
13     use type_queue_element
14     implicit none
15 type (queue_element), allocatable, dimension (:) :: queue
16 integer :: queue_head, queue_count, queue_tail
17 integer :: max_elements
18 contains
19     subroutine init_queue(max_elem)
20     implicit none
21     integer, intent(in) :: max_elem
22     queue_head = -1
23     queue_tail = -1
24     queue_count = 0
25     max_elements = max_elem
26     allocate (queue(0:max_elem-1))
27 end subroutine
28 function insert_queue(descriptor, depth) result(inserted)
29 implicit none
30 logical :: inserted
31 character (kind=1,len=20), intent(in) :: descriptor
32 integer*4, intent(in) :: depth
33 if(queue_count .eq. 0) then
34     queue_head = 0
35     queue_tail = 0
36     queue(queue_tail)%desc = descriptor
37     queue(queue_tail)%depth = depth
38     queue_count = queue_count + 1
39     inserted = .true.
40 else if (queue_count .eq. max_elements) then
41     inserted = .false.
42 else
43     queue_tail = (mod(queue_tail+1,max_elements))
44     queue(queue_tail)%desc = descriptor
45     queue(queue_tail)%depth = depth
46     queue_count = queue_count + 1
47     inserted = .true.
48 endif
49 end function
50 function delete_queue_begin(descriptor, depth) result(success)
51 implicit none
52 logical :: success
53 character (kind=1,len=20) :: descriptor
54 integer*4,intent(out) :: depth
55 if(queue_count .eq. 0) then
56     success = .false.
57 else if(queue_count .eq. 1) then
58     descriptor = queue(queue_head)%desc
59     depth = queue(queue_head)%depth
60     queue_head = -1
61     queue_tail = -1
62     queue_count = 0
63     success = .true.
64 else
65     descriptor = queue(queue_head)%desc
66     depth = queue(queue_head)%depth
67     queue_head = mod(queue_head+1,max_elements)
68     queue_count = queue_count - 1
69     success = .true.
70 endif
71 end function
72 function delete_queue_end(descriptor, depth) result(success)
73 implicit none
74 logical :: success
75 character (kind=1,len=20) :: descriptor
76 integer*4,intent(out) :: depth
77 if(queue_count .eq. 0) then
78     success = .false.
79 else if(queue_count .eq. 1) then
80     descriptor = queue(queue_tail)%desc
81     depth = queue(queue_tail)%depth

```

```
82         queue_head = -1
83         queue_tail = -1
84         queue_count = 0
85         success = .true.
86     else
87         descriptor = queue(queue_tail)%desc
88         depth = queue(queue_tail)%depth
89         queue_tail = queue_tail - 1
90         if(queue_tail .lt. 0) queue_tail=max_elements-1
91         queue_count = queue_count - 1
92         success = .true.
93     endif
94     end function
95 end module
```