

Programming Models for Scientific Computing on Leadership Computing Platforms: The Evolution of Coarray Fortran

John Mellor-Crummey
Department of Computer Science
Rice University

Acknowledgments

- CAF
 - Yuri Dotsenko
 - Cristian Coarfa
 - Bill Scherer
 - Laksono Adhianto
 - Guohua Jin
 - Fengmei Zhao

- HPF
 - Guohua Jin
 - Daniel Chavarria-Miranda
 - Vikram Adve

Outline

- Challenges for parallel languages
- A bit of parallel Fortran history
- Coarray Fortran, circa 1998 (CAF98)
- Assessment of CAF98
- A look at the emerging Fortran 2008 standard
- Recommendations for moving forward
- Open issues

Challenges for Parallel Languages

To succeed, a parallel programming language must ...

- be ubiquitous
 - multicore laptops
 - clusters on site
 - leadership-class machines at national centers
- be expressive
 - arbitrary algorithms
 - complex data structures
 - sophisticated parallelizations
- be productive
 - easy to write
 - easy to read and maintain
 - easy to reuse
- leverage legacy code: must be interoperable
- have a promise of future availability and longevity
- be supported by tools
- be efficient

On Being Efficient

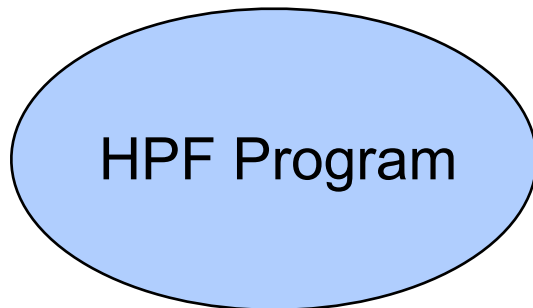
- Achieving high efficiency on multi- and many-core nodes
 - partitioning applications onto heterogeneous cores
 - match characteristics of application to architecture capabilities
 - achieving parallelism on node
 - multiple levels of parallelism: MIMD between cores, SIMD instructions
 - managing memory hierarchy at all levels: spatial and temporal reuse
 - GTC experience
 - organize the data structures for good spatial locality, vectorization
 - organize the loop nests to exploit temporal locality
 - dynamically reorganize particles
 - translate physical proximity of the particles into memory locality
- Scaling to large systems
 - partition well for good parallelism (avoid serialization, load imbalance)
 - avoid excessive communication frequency and volume
 - overlap communication latency with computation

1990s Vision: The Compiler was King

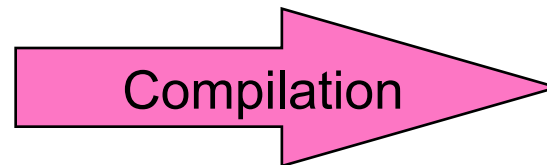
High Performance Fortran

Partitioning of data drives partitioning of computation, communication, and synchronization

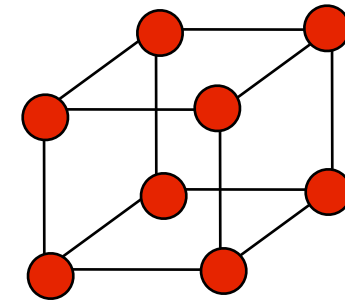
*Fortran program
+ data partitioning*



*Partition computation
Insert communication
Manage storage*



*Same answers as
sequential program*



Parallel Machine

Compiling HPF

- Partition data
 - follow user directives
- Select mapping of computation to processors
 - co-locate computation with data
- Analyze communication requirements
 - identify references that access off-processor data
- Partition computation by reducing loop bounds
 - schedule each processor to compute on its own data
- Insert communication
 - exchange values as needed by the computation
- Manage storage for non-local data

Some Lessons from HPF

Everything matters for good performance!

- Good data partitionings are essential for good parallelizations
 - e.g. BLOCK partitionings inferior to multipartitionings for line-sweeps
- Excess communication undermines scalability
 - both frequency and volume must be right
 - coalesce communication sets for multiple references
 - 41% lower message volume, 35% faster: NAS SP @ 64 procs
 - partially replicate computation to reduce communication
 - 66% lower message volume, 38% faster: NAS BT @ 64 procs
- User guidance is an invaluable supplement to analysis
 - e.g. HPF/JA directives to control communication; parallel loops
- Complex things are possible; realizing them can be challenging
 - RandomAccess, FFT
- If the compiler can't deliver, you're out of luck!

Partitioned Global Address Space Languages

- Global address space
 - one-sided communication (GET/PUT) simpler than msg passing
- Programmer has control over performance-critical factors
 - data distribution and locality control lacking in OpenMP
 - computation partitioning } HPF & OpenMP compilers
 - communication placement } must get this right
- Data movement and synchronization as language primitives
 - amenable to compiler-based communication optimization

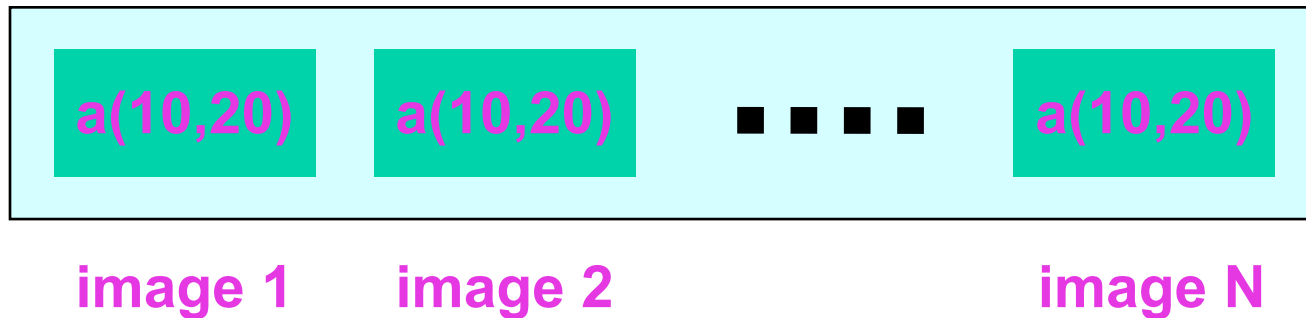
Coarray Fortran in 1998 (CAF98)

Numrich and Reid [ACM Fortran Forum, 17:2, 1998]

- SPMD process images
 - fixed number of images during execution
 - images operate asynchronously
- Two-level memory model for managing locality (local vs. remote)
 - real $y(20, 20)$ a private 20x20 array in each image
 - real $y(20, 20)$ [*] a shared 20x20 array in each image
- Simple one-sided shared-memory communication
 - $x(:,r) = y(:,l)$ [right] copy column from processor right into local column
 - $b(:,:) = a[p,q]\%data(:,:)$ fill a local array from remote processor [p,q]
- Synchronization
 - critical section mutual exclusion
 - `synch_all()` barrier
 - `sync_team(notify, wait)`
 - `synch_memory` memory fence
- Pointers and dynamic allocation of pointer components

One-sided Communication with Coarrays

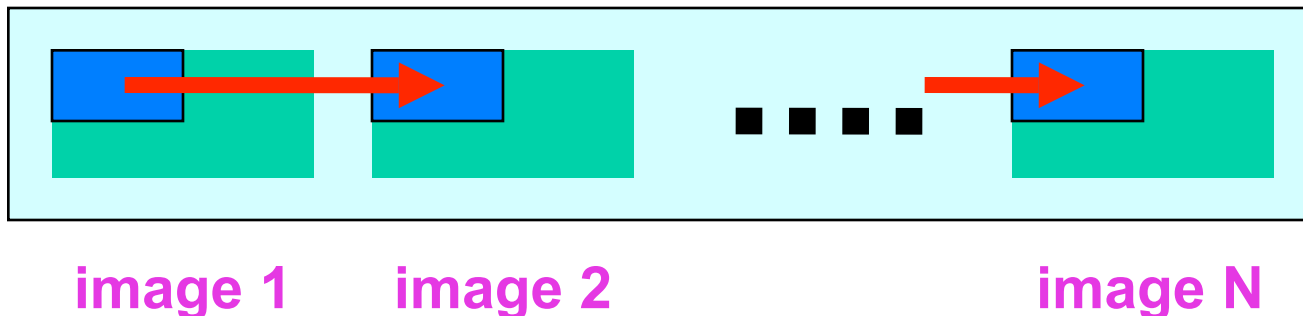
```
integer a(10,20) [*]
```



```
if (this_image() > 1) then
```

```
  a(1:5,1:10) = a(1:5,1:10) [this_image() - 1]
```

```
endif
```



The Power of CAF98

A Finite Element Example (Numrich, Reid; 1998)

```
subroutine assemble(start, prin, ghost, neib, x)
  integer :: start(:),prin(:),ghost(:),neib(:),k1, k2, p
  real :: x(:) [*]
  call sync_all(neib)
  do p = 1, size(neib) ! Add contribs. from ghost regions
    k1 = start(p); k2 = start(p+1)-1
    x(prin(k1:k2))=x(prin(k1:k2))+x(ghost(k1:k2))[neib(p)]
  enddo
  call sync_all(neib)
  do p = 1, size(neib) ! Update the ghosts
    k1 = start(p); k2 = start(p+1)-1
    x(ghost(k1:k2))[neib(p)] = x(prin(k1:k2))
  enddo
  call sync_all
end subroutine assemble
```

High-level Assessment of CAF

- Advantages
 - admits sophisticated parallelizations with compact expression
 - doesn't require as much from compilers
 - yields scalable high performance today with careful programming
 - if you put in the effort, you can get the results
- Disadvantages
 - users code data movement and synchronization
 - tradeoff between the abstraction of HPF vs. control of CAF
 - optimizing performance can require intricate programming
 - buffer management is fully exposed!
 - expressiveness is a concern for CAF
 - insufficient primitives to express a wide range of programs

A Closer Look at CAF98 Details

- Strengths
 - one-sided data access can simplify some programs
 - vectorized access to remote data can be efficient
 - amortizes communication startup costs
 - data streaming can hide communication latency (e.g. on the Cray X1)
- Weaknesses
 - synchronization can be expensive
 - single critical section was very limiting
 - `synch_all`, `synch_team` were not sufficient
 - `synch_all`: barriers are a heavy-handed mechanism
 - `synch_team` semantics required $O(n^2)$ pairwise communication
 - rolling your own collectives doesn't lead to portable high performance
 - latency hiding is impossible in important cases
 - procedure calls had implicit barriers to guarantee data consistency
 - communication couldn't be overlapped with a procedure call

Emerging Fortran 2008 Standard, Feb. 2008

Coarray features being considered for inclusion

- Single and multidimensional coarrays
- Collective allocation of coarrays to support a symmetric heap
- Named critical sections for structured mutual exclusion
 - improves over single critical section in CAF98
- Synch_all and synch_team
 - pre-arrange teams with form_team for efficient synchronization
- Pairwise non-blocking synchronization with notify/query
 - support communication/computation overlap
- Collective communication intrinsics for convenience & portability

Are F2008 Coarrays Ready for Prime Time?

Questions worth considering

1. What types of parallel systems are viewed as the important targets for Fortran 2008?
2. Does Fortran 2008 provide the set of features necessary to support parallel scientific libraries that will help catalyze development of parallel software using the language?
3. What types of parallel applications is Fortran 2008 intended to support and is the collection of features proposed sufficiently expressive to meet those needs?
4. Will the collection of coarray features described provide Fortran 2008 facilitate writing portable parallel programs that deliver high performance on systems with a range of characteristics?

1. Target Architectures?

CAF support must be ubiquitous or (almost) no one will use it

- Important targets
 - clusters and leadership class machines
 - multicore processors and SMPs
- Difficulties
 - F2008 CAF lacks flexibility, which makes it a poor choice for multicore
 - features are designed for regular, SPMD
 - multicore will need better one-sided support
 - current scalable parallel systems lack h/w shared memory
 - e.g. clusters, Blue Gene, Cray XT, SiCortex
 - big hurdle for third-party compiler vendors to target scalable systems

2. Adequate Support for Libraries?

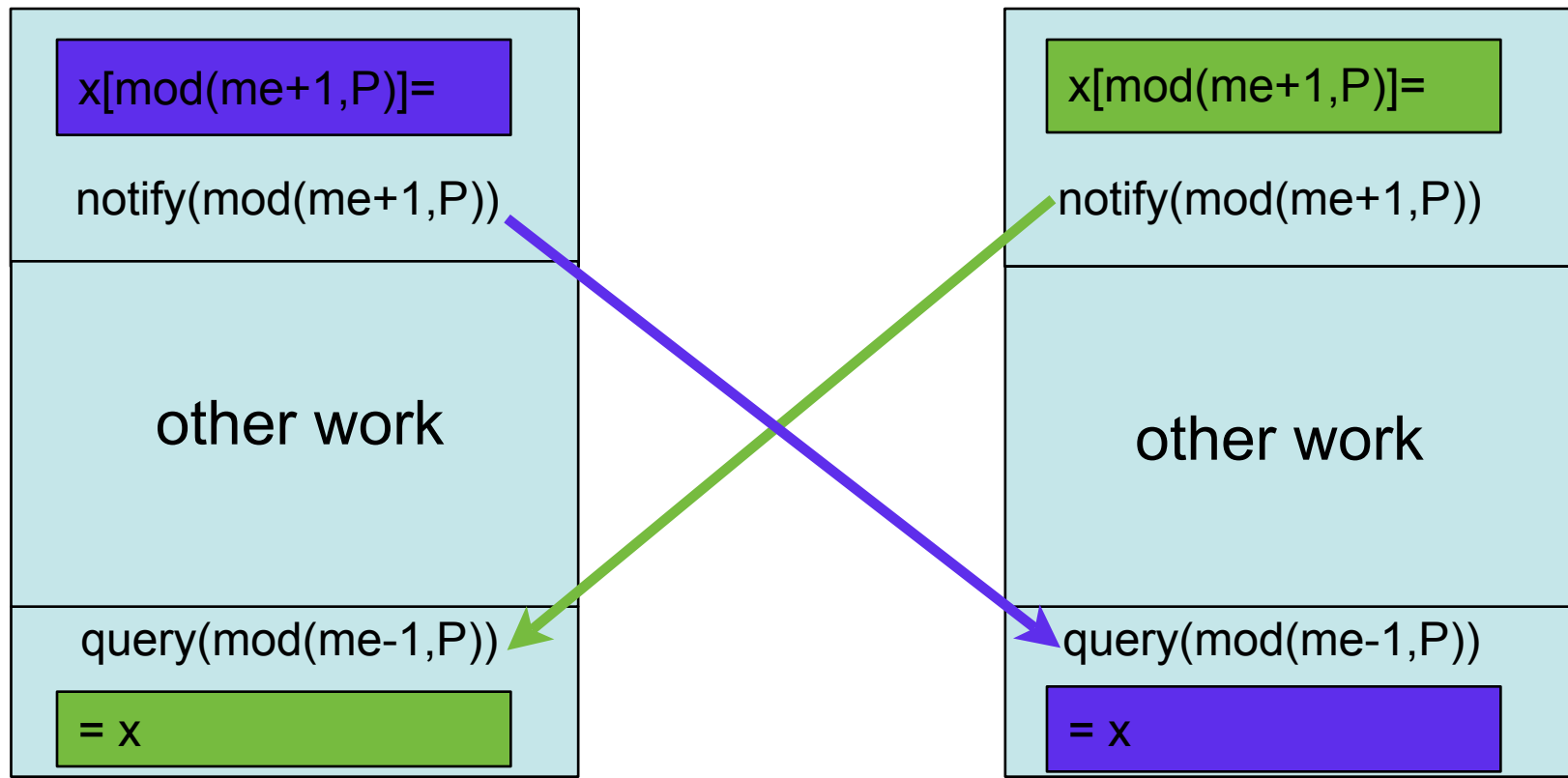
Lessons from MPI: Library needs [MPI 1.1 Standard]

- **Safe communication space:** libraries can communicate as they need to, without conflicting with communication outside the library
- **Group scope for collective operations:** allow libraries to avoid unnecessarily synchronizing uninvolved processes
- **Abstract process naming:** allow libraries to describe their communication to suit their own data structures and algorithms
- **Provide a means to extend the message-passing notation:** user-defined attributes, e.g., extra collective operations

Lack of a Safe Communication Space (Part 1)

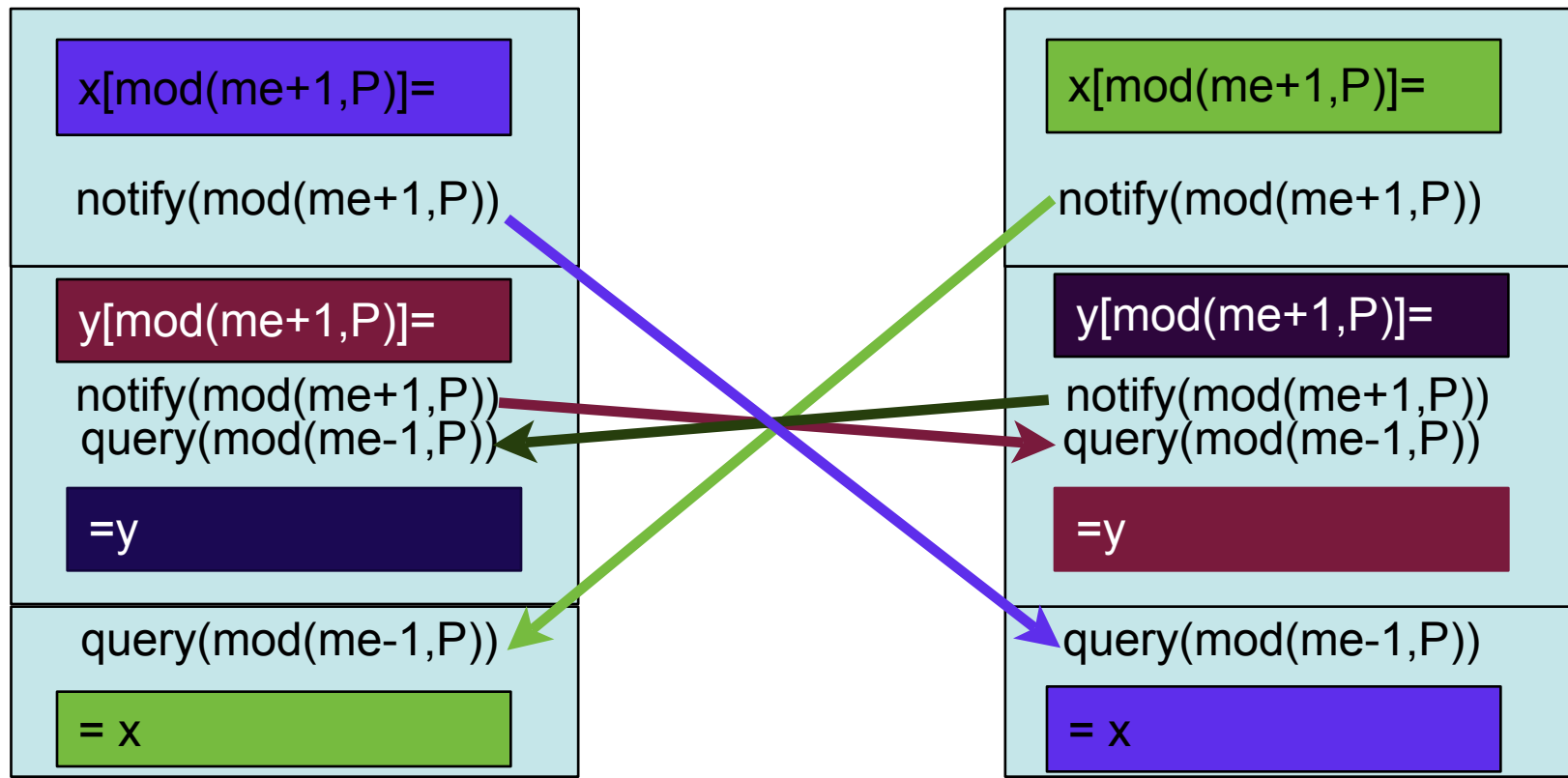
- F2008 has a single synchronization channel between process pairs
- Multiple channels are essential for library encapsulation

Consider the following



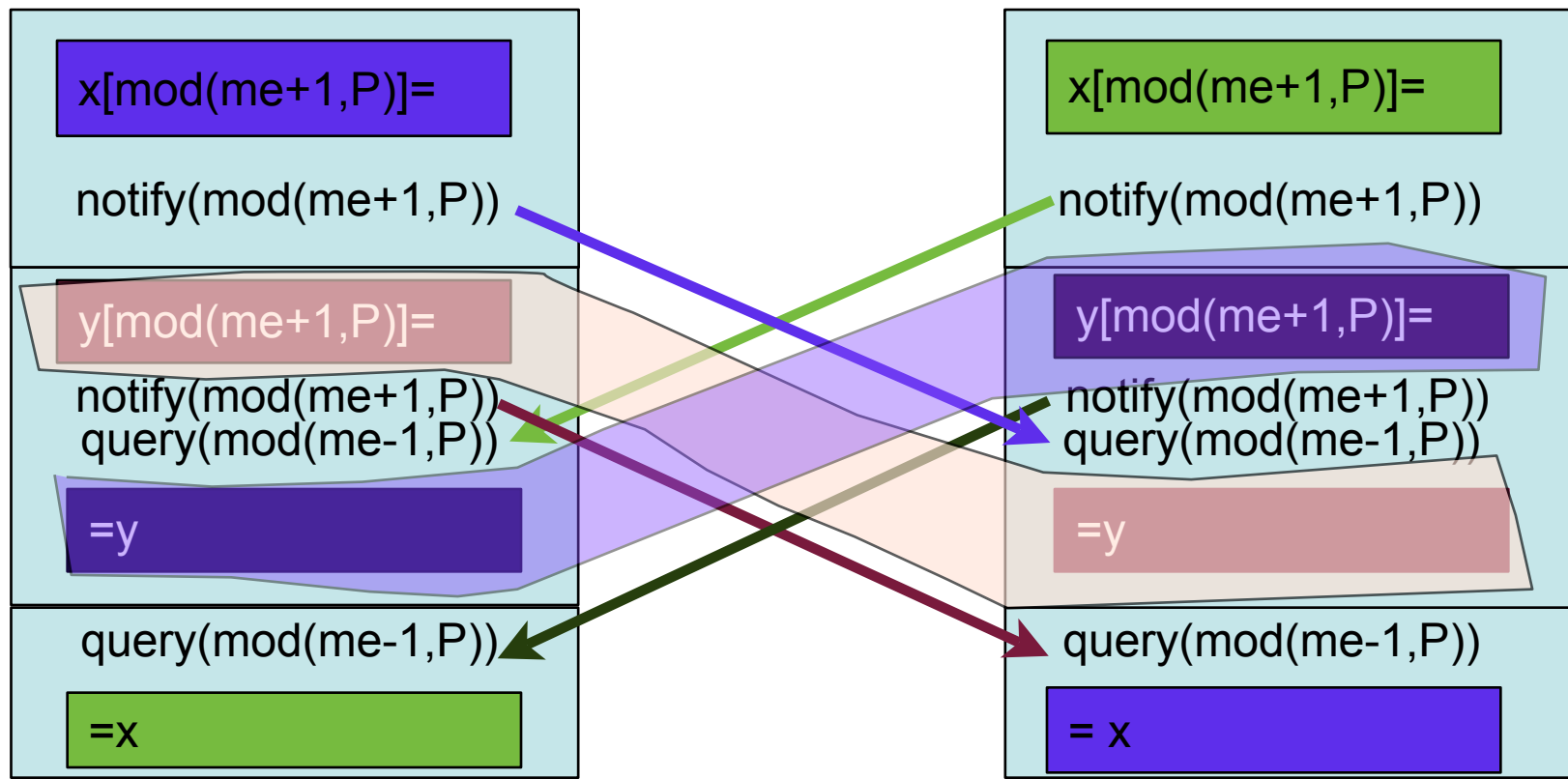
Lack of a Safe Communication Space (Part 2)

What if “other work” was synchronized?



Lack of a Safe Communication Space (Part 3)

Lack of encapsulation leads to data races



Lack of Support for Process Subsets

A library can't conveniently operate on a process subset

- Multidimensional coarrays
 - must be allocated across all process images
 - can't conveniently employ this abstraction for a process subset
- Image naming
 - all naming of process images is global
 - makes it harder to work within process subsets
 - must be cognizant of their embedding in the whole
- Allocation/deallocation
 - libraries shouldn't unnecessarily synchronize uninvolved processes
 - but ... coarrays in F2008 require
 - global collective allocation/deallocation
 - serious complication for coupled codes on process subsets
 - complete loss of encapsulation

Lack of a Means to Extend Collectives

- A variety of collective reduction subroutines
 - sum, product, maxloc, maxval, minloc, minval
 - any, all, count
- Shortcomings
 - no support for user-defined reduction operators
 - no support for scan reductions
 - missing feature: all-to-all

3. Target Application Domains?

- Can F2008 support applications that require one-sided update of mutable shared dynamic data structures?
- No. Two key problems
 - can't add a piece to a partner's data structure
 - F2008 doesn't support remote allocation
 - F2008 doesn't support pointers to remote data
 - F2008 doesn't support remote execution using "function shipping"
 - synchronization is inadequate
 - named critical sections overly limit concurrency
 - only one process active per static name
 - unreasonable to expect users to "roll their own"
- As defined, F2008 useful for halo exchanges on dense arrays

4. Adequate Support for Writing Fast Code?

- Lack of support for hiding synchronization latency
 - F2008 notify/query specification precludes overlap
 - no split-phase barrier
- Lack of support for exploiting locality in machine topology
- Lack of a precise memory model
 - developers must use loose orderings where possible
 - must be able to reason about what behaviors one should expect
 - programs must be resilient to reorderings

Lessons from MPI 1.1

What capabilities are needed for parallel libraries?

- Abstraction for a group of processes
 - functions for constructing and manipulating process groups
- Virtual communication topologies
 - e.g. cartesian, graph
 - neighbor operation for indexing
- Multiple communication contexts
 - e.g. parallel linear algebra uses multiple communicators
 - rows of blocks, columns of blocks, all blocks

Recommendations for Moving Forward (Part 1)

- Only one-dimensional co-arrays
 - no collective allocation/deallocation: require users to synchronize
- Team abstraction that represents explicitly ordered process groups
 - deftly supports coupled codes, linear algebra
 - enables renumbering to optimize embedding in physical topology
- Topology abstraction for groups: cartesian and graph topologies
 - cartesian is a better alternative to k-D coarrays
 - supports processor subsets, periodic boundary conditions as well
 - graph is a general abstraction for all purposes
- Multiple communication contexts
 - apply notify/query to semaphore-like variables for multiple contexts
- Add support for function shipping
 - spawn remote functions for latency avoidance
 - spawn local functions to exploit parallelism locally
 - lazy multithreading and work stealing within an image

Recommendations for Moving Forward (Part 2)

- Better mutual exclusion support for coordinating activities
 - short term: critical sections using lock variables
 - longer term: conditional ATOMIC operations based on transactional memory
- Enhanced support for collectives
 - add support for user-defined reduction operators
 - add support for scan reductions
 - add support for all-to-all operations
- Add multiversion variables
 - simplify producer/consumer interactions in a shared memory world
- Add global pointers to enable flat access to data

Open Questions (Part 1)

- Atomic operations: what is the best approach for scalable systems?
 - X10 approach: similar to k-way compare-and-swap?
 - transactional memory?
 - Bocchino, Adve, Chamberlain: Cluster STM [PPoPP '08]
 - weak atomicity; transactional and non-transactional data in a phase
 - flexible distributed atomic operations
 - function shipping for latency avoidance
 - shortcomings: rigid SPMD model (1 task per PE)

Open Questions (Part 2)

Synchronization with dynamic threading

- Barriers with dynamic threading: who participates?
- Alternatives
 - Cilk's "SYNC"
 - a SYNC in a procedure blocks until all its spawned work finishes
 - limited to rigid nested fork-join synchronization
 - X10's "FINISH" construct
 - all computation and threads inside a FINISH block must complete
 - more flexible than Cilk's model
 - an entire nested computation can complete to a single FINISH
 - FINISH is global
- Proposed approach
 - support FINISH on processor subsets (CAF teams)