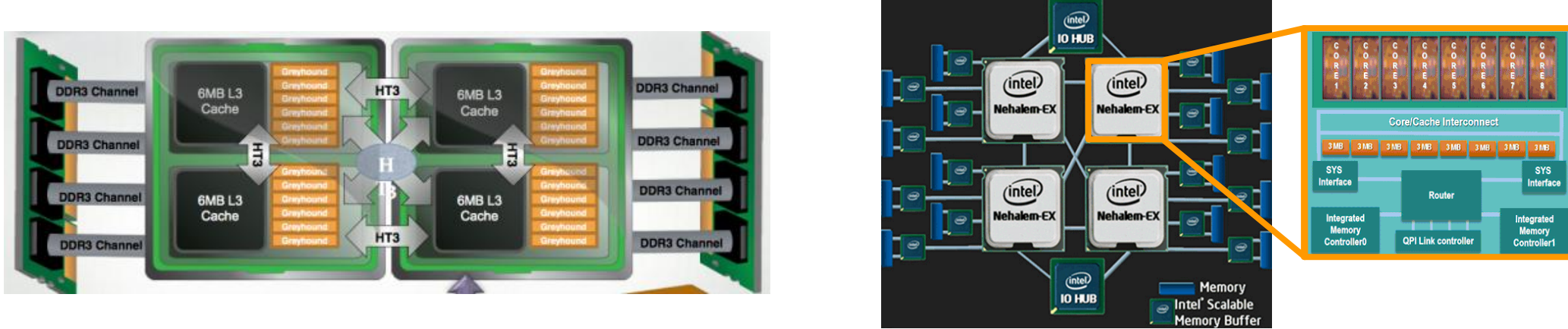


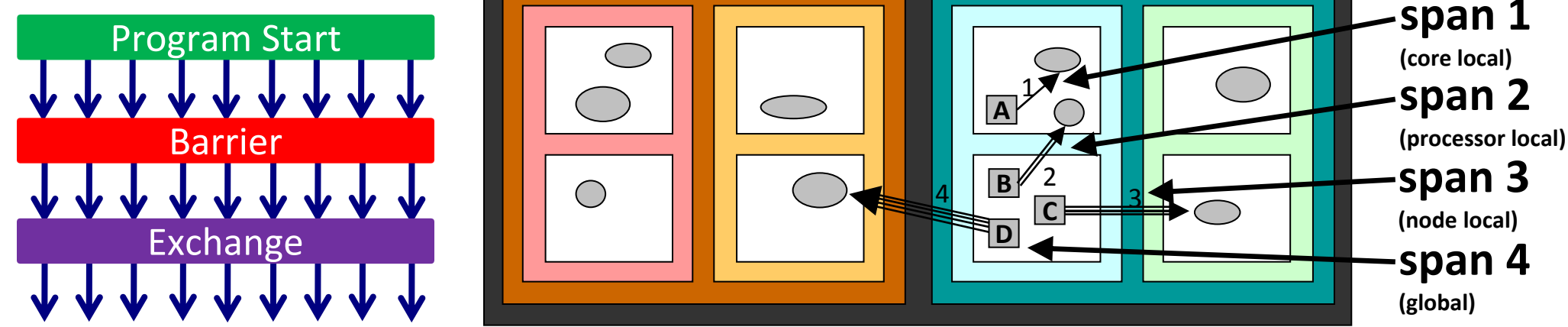


Hierarchical Machines and Algorithms

- Parallel machines are hierarchical



- Algorithms also may be hierarchical
 - Examples: merge sort, Barnes-Hut particle simulation
- Efficiency layer language must expose locality to users while also allowing general parallelism
- Our approach: start with parallelism model that has locality, add more general parallelism constructs
 - Single-Program, Multiple-Data parallelism model
 - Fixed number of threads, collective operations for synchronization, communication
 - Partitioned Global Address Space memory model
 - Threads can access data anywhere, locality encoded in pointer type



- Implementation in Titanium language, using GASNet runtime

Hierarchical Thread Teams

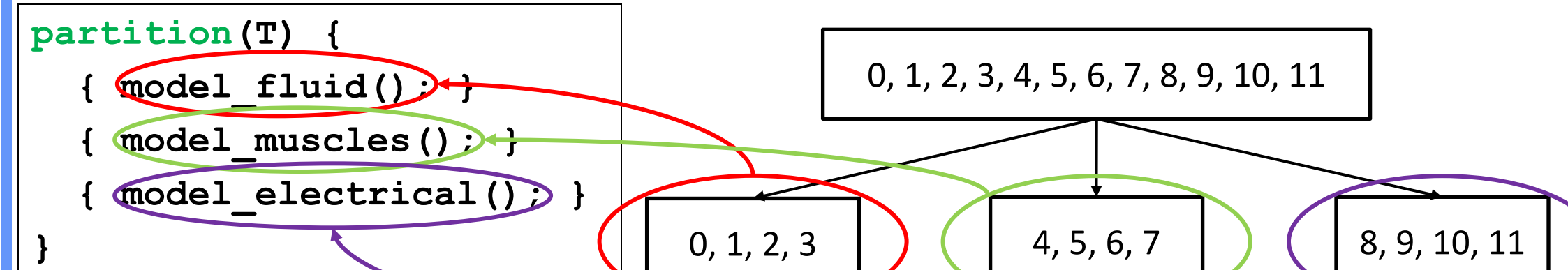
- Thread *teams* are basic units of cooperation
 - Groups of threads that cooperatively execute code
 - Collective operations over teams
- Teams should be hierarchical
 - Match hierarchical nature of machines, algorithms
- Team data structure
 - ```

graph TD
 Root[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11] --> T1[0, 1, 2, 3]
 Root --> T2[4, 5, 6, 7]
 Root --> T3[8, 9, 10, 11]
 T1 --> T1_1[1, 3, 2]
 T1 --> T1_2[0]
 T3 --> T3_1[9, 8]
 T3 --> T3_2[10, 11]

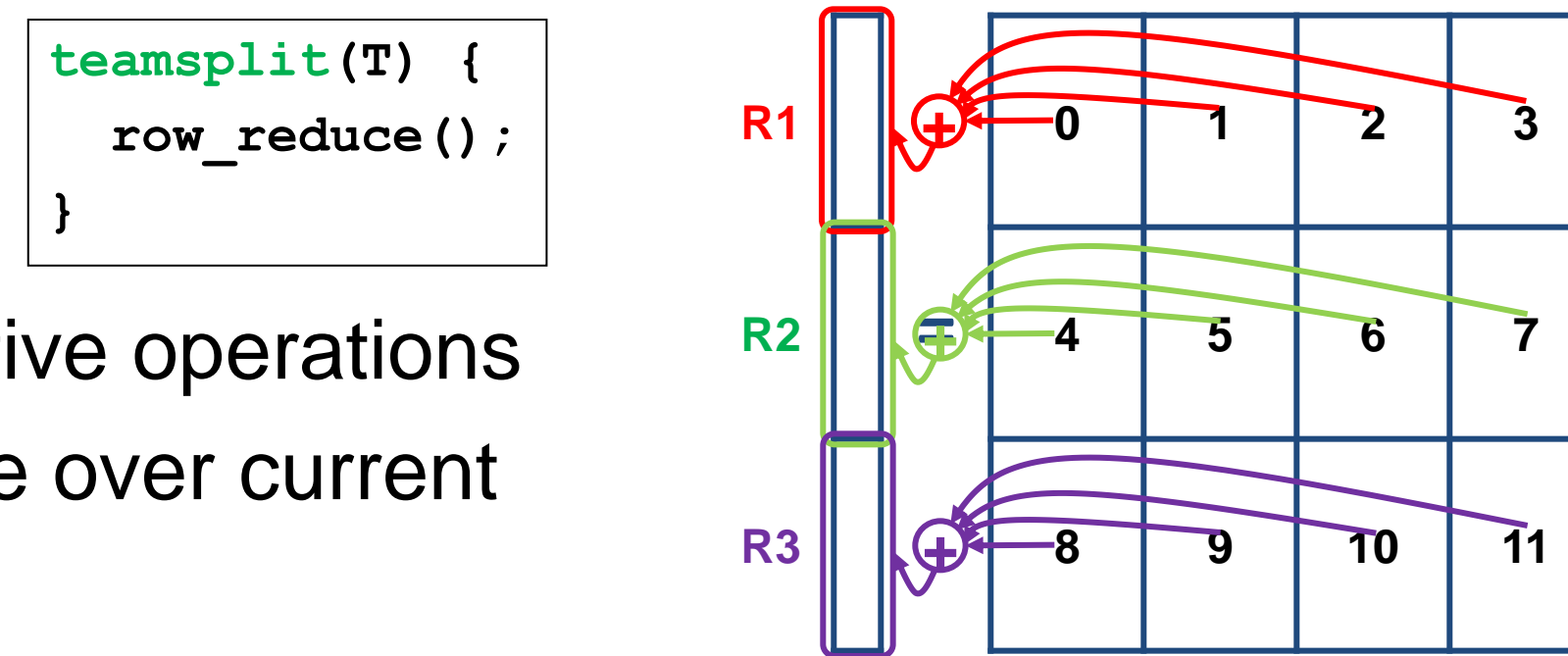
```
- Library functions provided to facilitate team creation
- Need to provide ability to query machine structure
  - At the moment, we provide a function to create a team hierarchy that loosely matches the machine

## Team Usage

- Thread teams may execute different code



- Threads may execute same code as part of different teams of cooperating threads



- Collective operations operate over current team

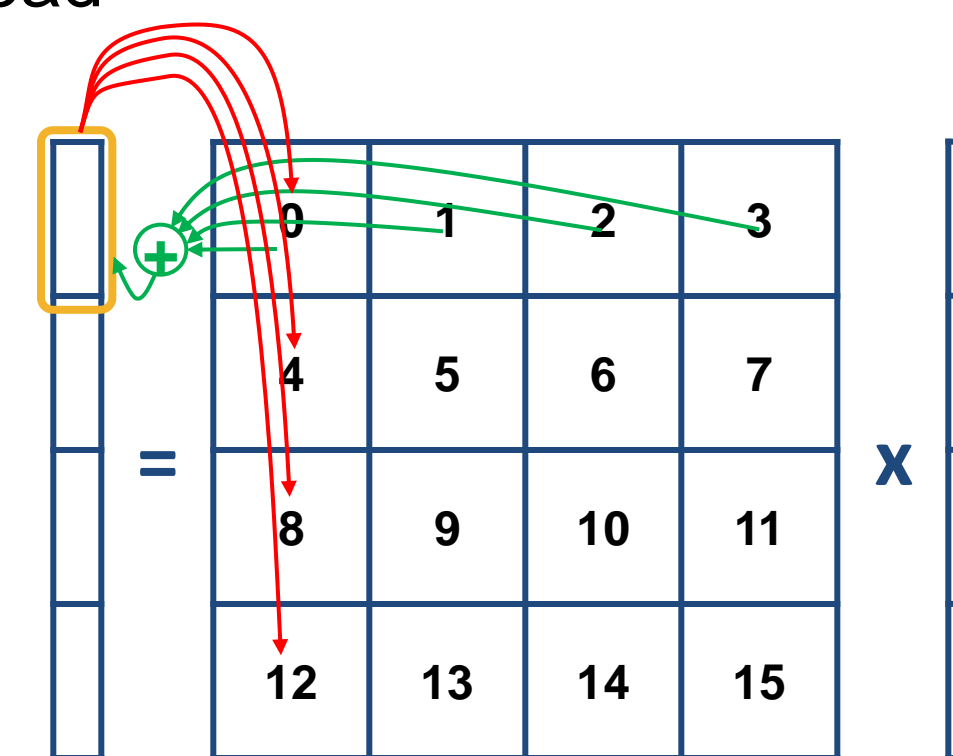
## Conjugate Gradient Library

- NAS conjugate gradient application originally written by Kaushik Datta
- Iterative solver with sparse matrix-vector multiply in each iteration
  - Matrix divided among threads by row and column
  - Result of one iteration is input of next
  - Reduction across rows, broadcast along columns
  - Original code had hand-written collectives; we use built-in team collectives instead

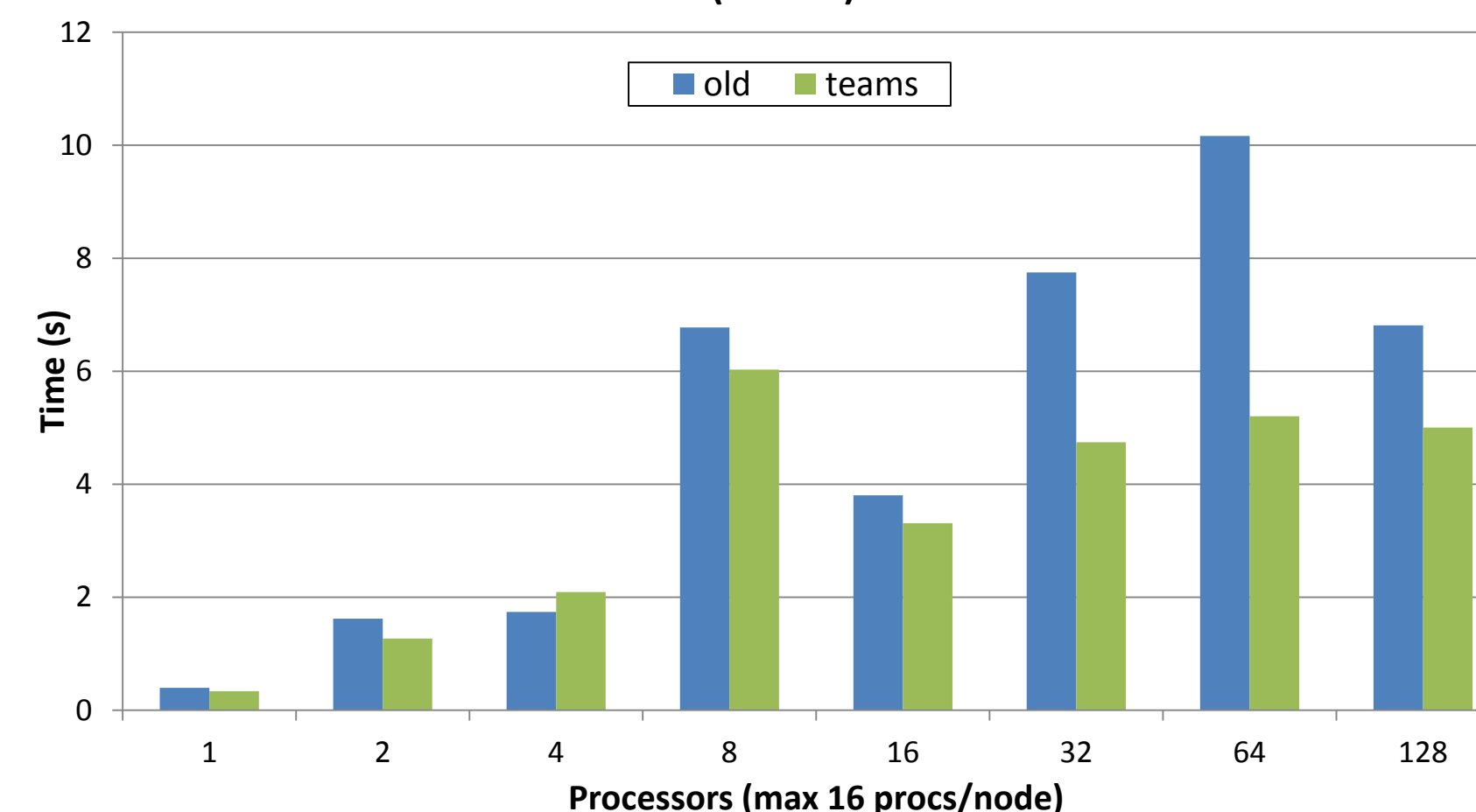
```

teamsplit(rowTeam) {
 Reduce.add(mtmp, myResults0, rpivot);
}
if (reduceCopy)
 myOut.copy(allResults[reduceSource]);
teamsplit(columnTeam) {
 myOut.vbroadcast(cpivot);
}

```

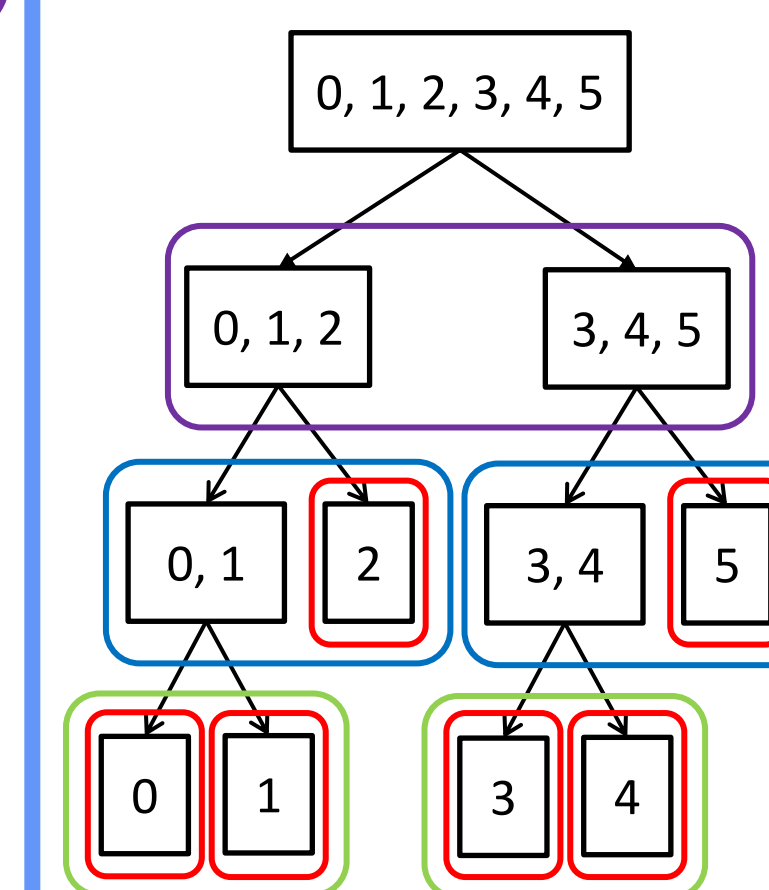


NAS CG SPMV Communication Time (Cray XE6)  
(Class B)



## Parallel Sorting Library

- Parallel sorting application using new hierarchical constructs
- SMP sort uses team constructs to merge results from different threads in parallel



- Phase 1: all threads sort
- Phase 2: t0 and t3 merge
- Phase 3: t0 and t3 merge
- Phase 4: t0 merges

```

void divideTeam(Team t) {
 if (t.size() > 1) {
 t.splitTeam(2);
 divideTeam(t.child(0));
 divideTeam(t.child(1));
 }
}

single void sortAndMerge(Team t) {
 if (Ti.numProcs() == 1) {
 sequentialSort(myData);
 } else {
 teamsplit(t) {
 sortAndMerge(Ti.currentTeam());
 }
 Ti.barrier();
 if (Ti.thisProc() == 0) {
 int otherProc = myProc + t.child(0).size();
 merge(data(myProc), data(otherProc));
 }
 }
}

```

- Initial distributed sort uses existing, unoptimized sample sort code to distribute elements among nodes, SMP sort to sort on each node
- 12 lines of code, 5 minutes to solution
- However, doesn't scale beyond 8 nodes

- New version uses all threads to distribute elements among nodes
- Still uses SMP sort to sort on each node, with no changes required to SMP code

Sample Sort v0.9a (Cray XT4)  
(10,000,000 elements/proc, 10,000 samples/proc)

