

Scientific Libraries: MPI + Manycore Issues and Plans

Dense Numerical Linear Algebra Library at Scale

Jack Dongarra

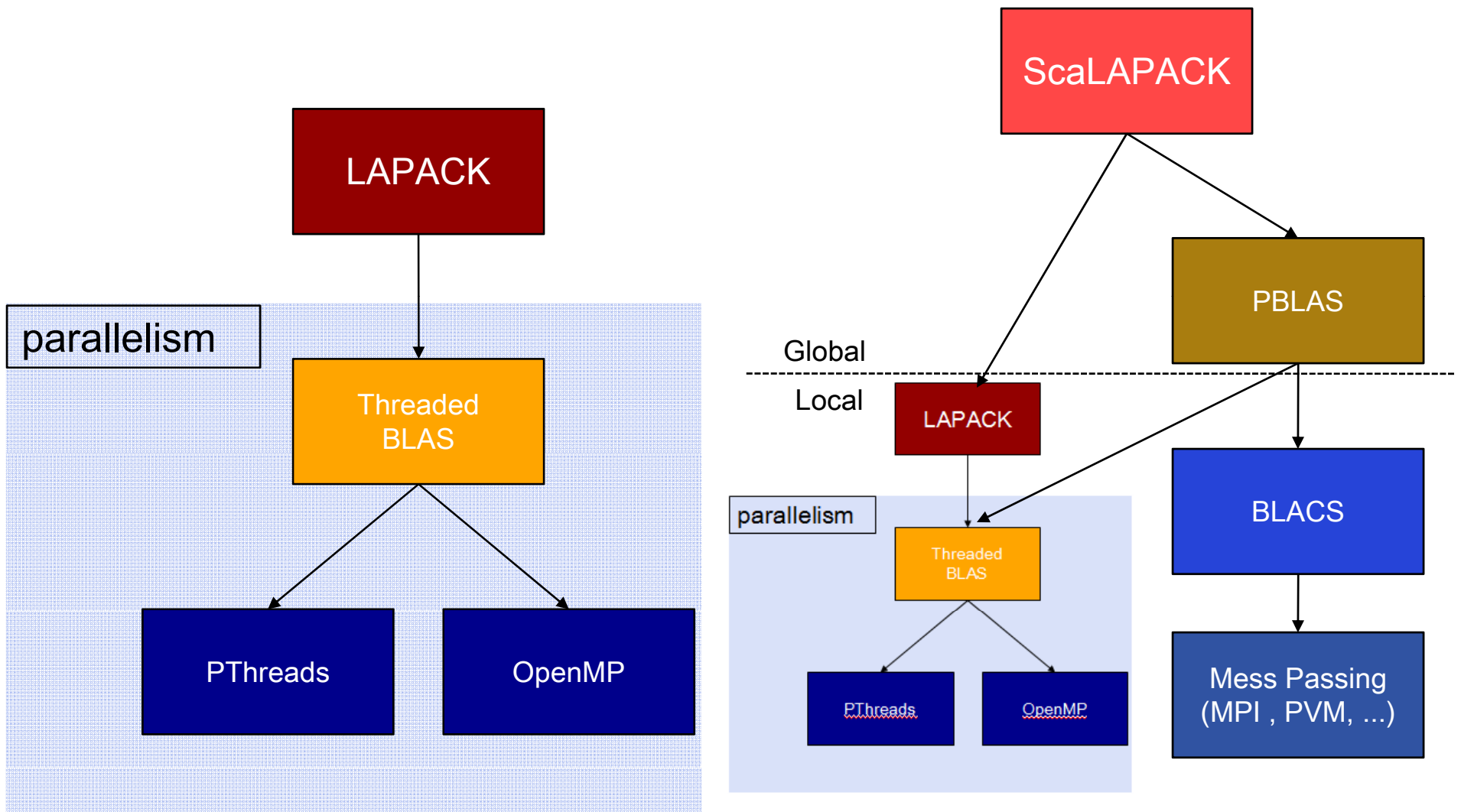
University of Tennessee
Oak Ridge National Laboratory



Major Changes to Software

- **Must rethink the design of our software**
 - **Another disruptive technology**
 - Similar to what happened with cluster computing and message passing
 - **Rethink and rewrite the applications, algorithms, and software**
- **Numerical libraries for example will change**
 - **For example, both LAPACK and ScaLAPACK will undergo major changes to accommodate this**

LAPACK and ScaLAPACK



About 1 million lines of code



Coding for an Abstract Multicore

Parallel software for multicores should have two characteristics:

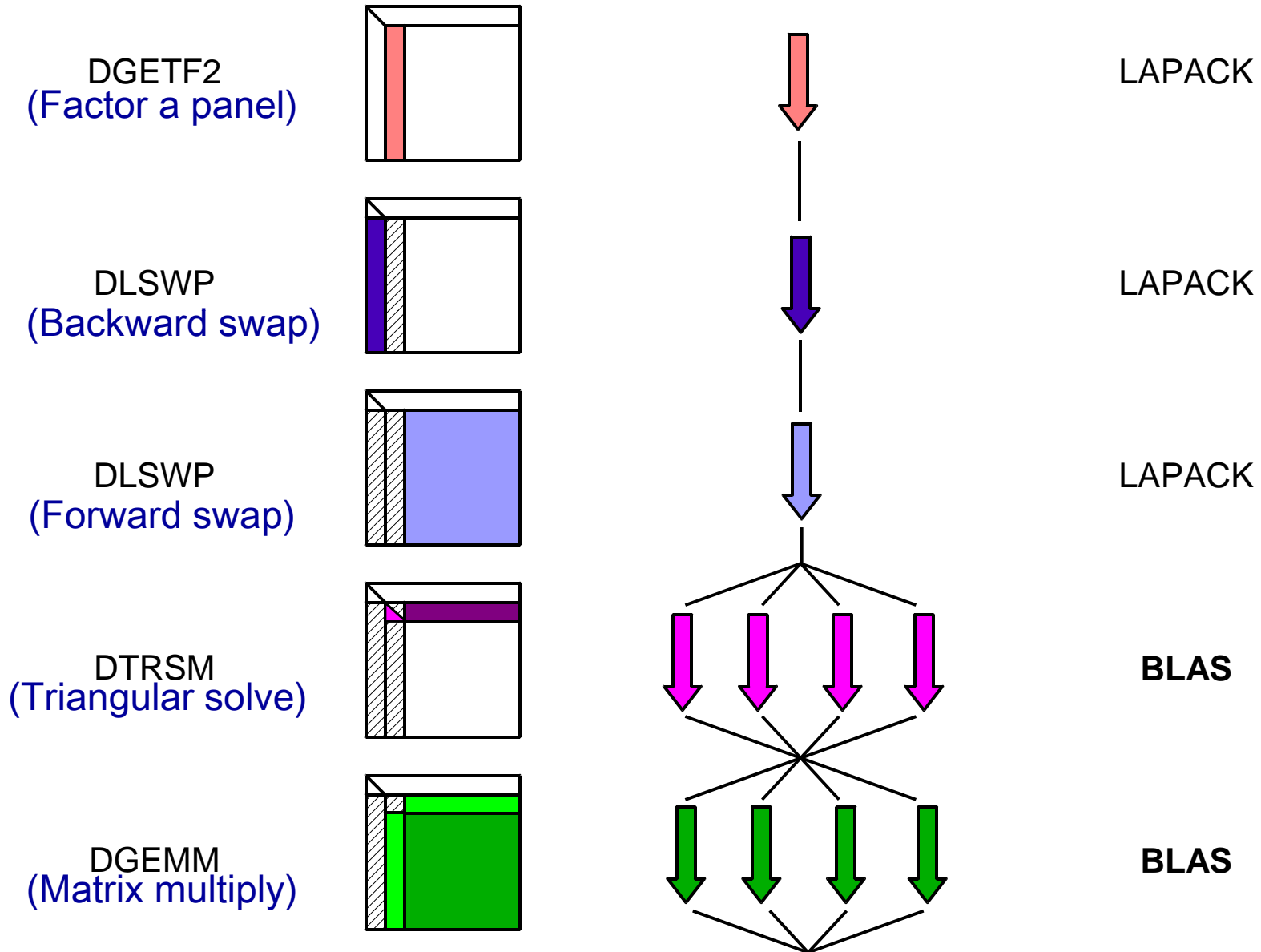
- **Fine granularity:**
 - High level of parallelism is needed
 - Cores will probably be associated with relatively small local memories. This requires splitting an operation into tasks that operate on small portions of data in order to reduce bus traffic and improve data locality.
- **Asynchronicity:**
 - As the degree of thread level parallelism grows and granularity of the operations becomes smaller, the presence of synchronization points in a parallel execution seriously affects the efficiency of an algorithm.



ManyCore - Parallelism for the Masses

- We are looking at the following concepts in designing the next numerical library implementation
 - Dynamic Data Driven Execution
 - Self Adapting / Auto Tuning
 - Block Data Layout
 - Mixed Precision in the Algorithm
 - Exploit Hybrid Architectures
 - Fault Tolerant Methods
 - Communication Avoiding Algorithms

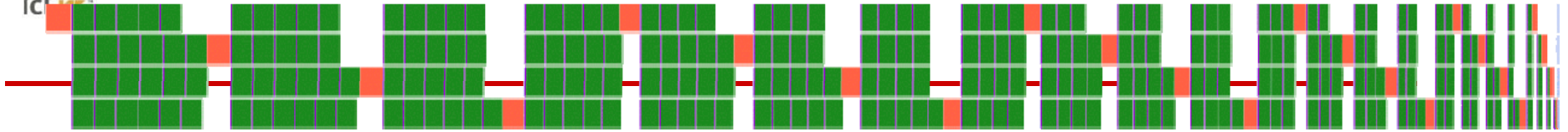
Steps in the LAPACK LU



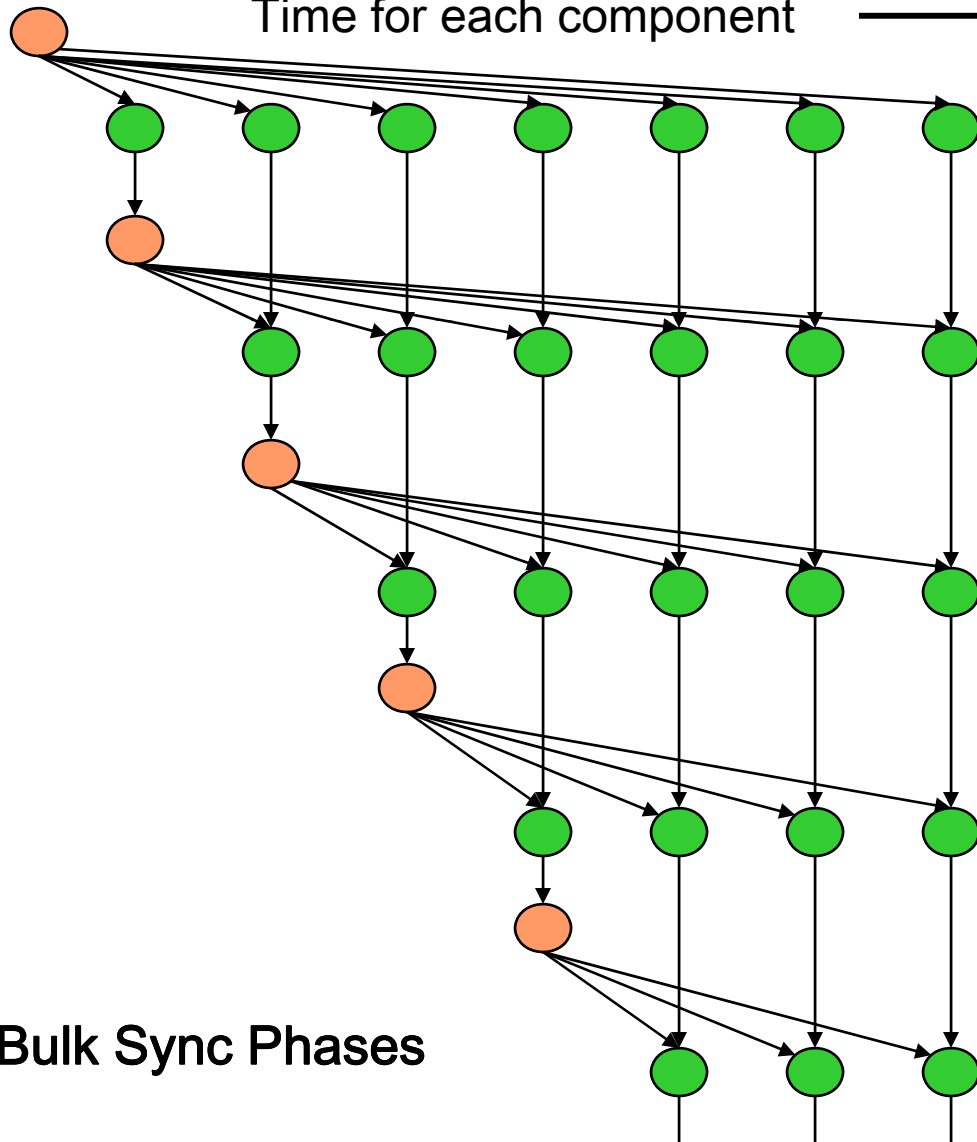


LU Timing Profile (4 core system)

Threads – no lookahead



Time for each component →



Bulk Sync Phases

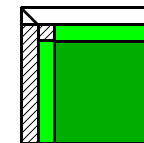
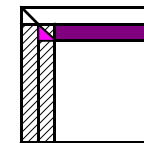
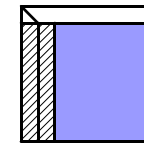
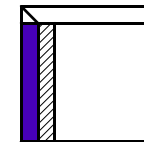
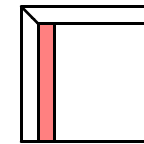
DGETF2






DLSWP

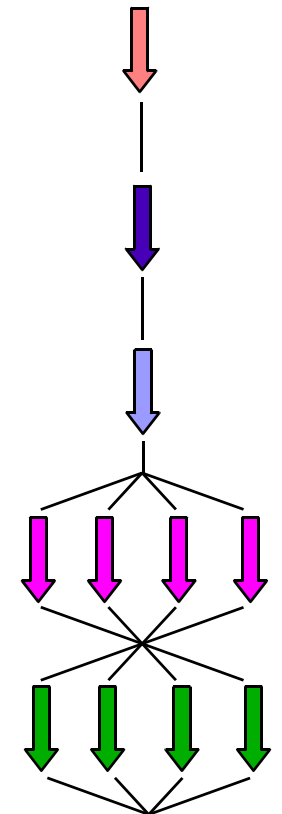
DLSWP

DTRSM

DGEMM

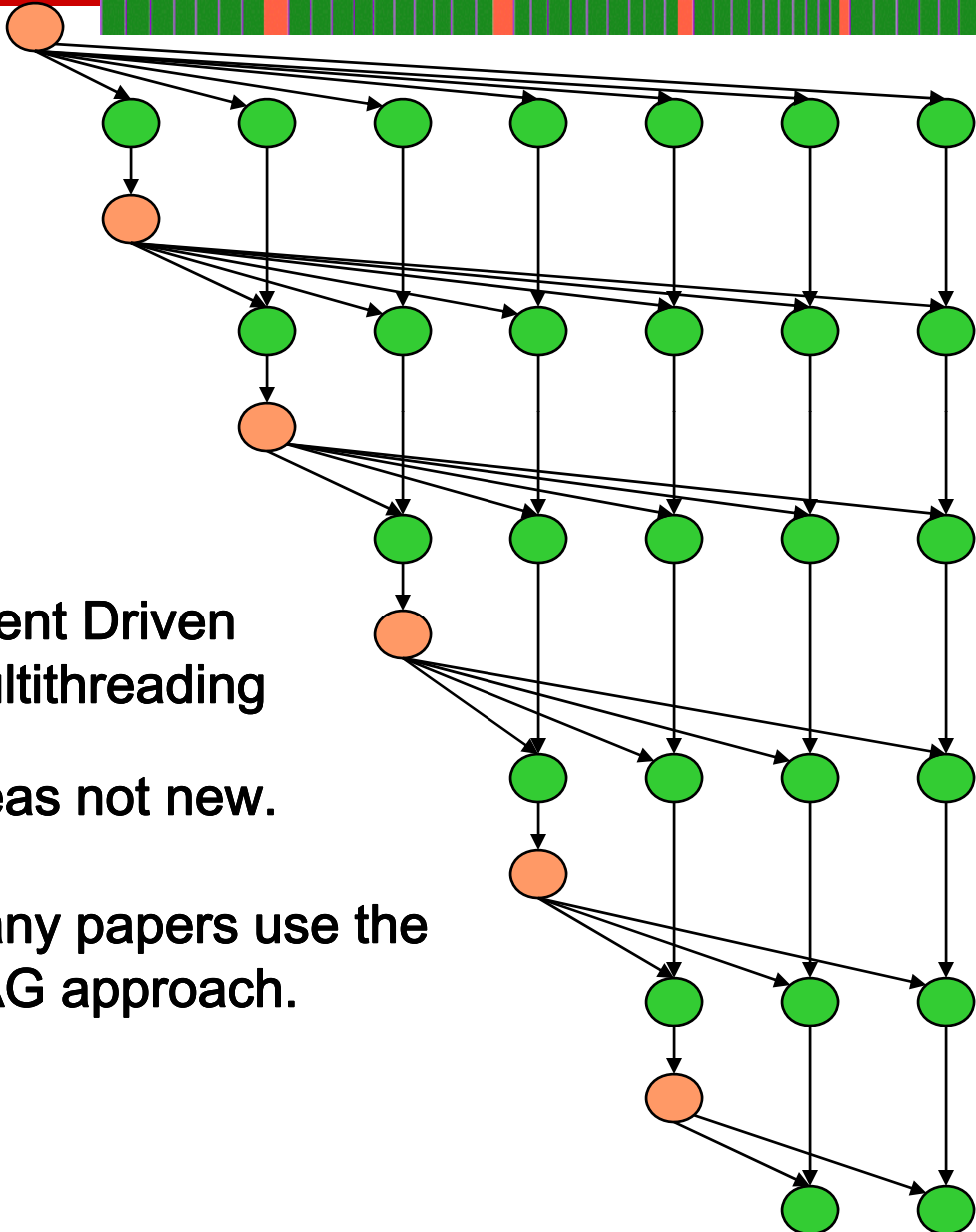
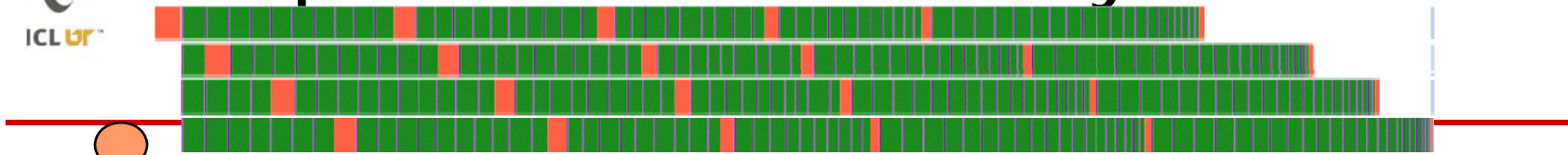


-  DGETF2
-  DLASWP(L)
-  DLASWP(R)
-  DTRSM
-  DGEMM





Adaptive Lookahead - Dynamic



Event Driven
Multithreading

Ideas not new.

Many papers use the
DAG approach.

```
while(1)
  fetch_task();
  switch(task.type) {
    case PANEL:
      dgetf2();
      update_progress();
    case COLUMN:
      dlaswp();
      dtrsm();
      dgemm();
      update_progress();
    case END:
      for()
        dlaswp();
      return;
  }
}
```

**Reorganizing
algorithms to use
this approach**

Redesign

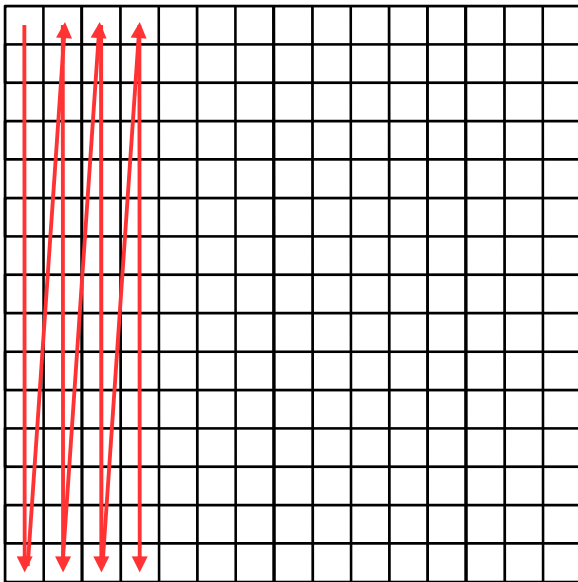
- **Asynchronicity**
 - Avoid fork-join (Bulk sync design)
- **Dynamic Scheduling**
 - Out of order execution
- **Fine Granularity**
 - Independent block operations
- **Locality of Reference**
 - Data storage - Block Data Layout



Achieving Fine Granularity

Fine granularity may require novel data formats to overcome the limitations of BLAS on small chunks of data.

Column-Major

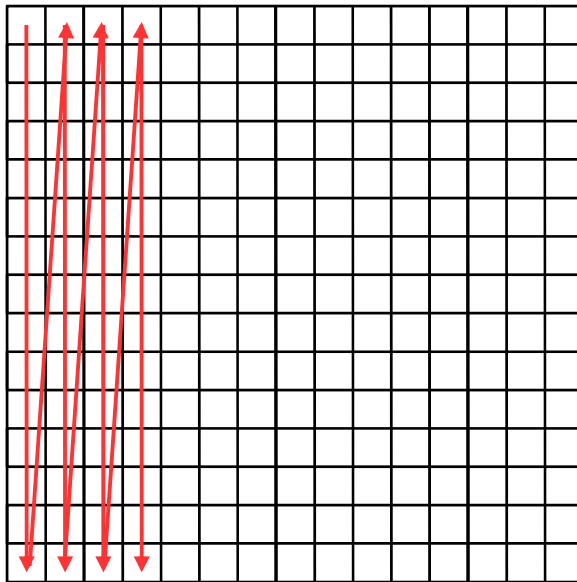




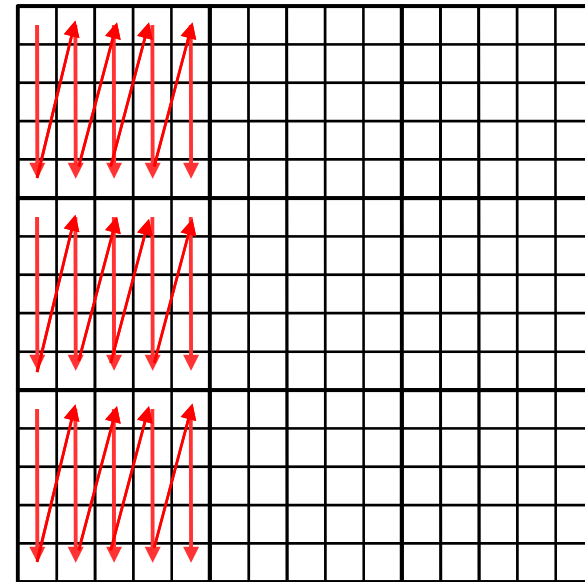
Achieving Fine Granularity

Fine granularity may require novel data formats to overcome the limitations of BLAS on small chunks of data.

Column-Major



Blocked

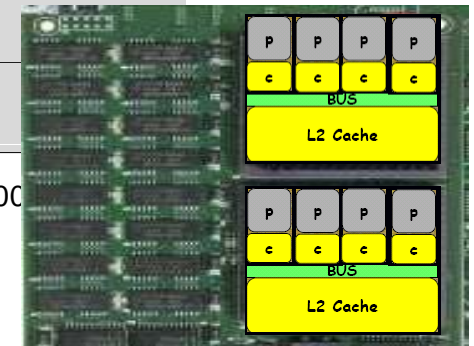
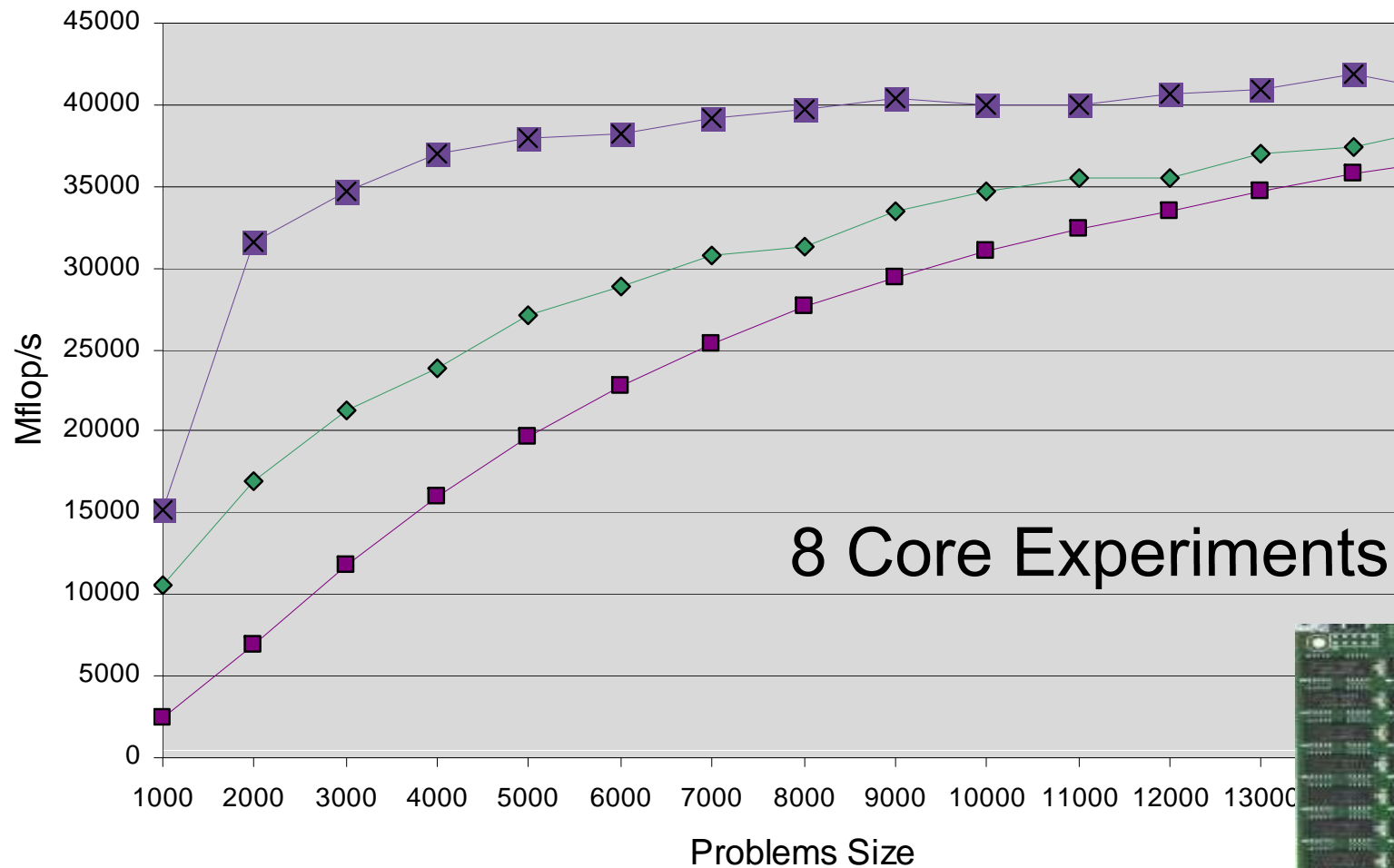




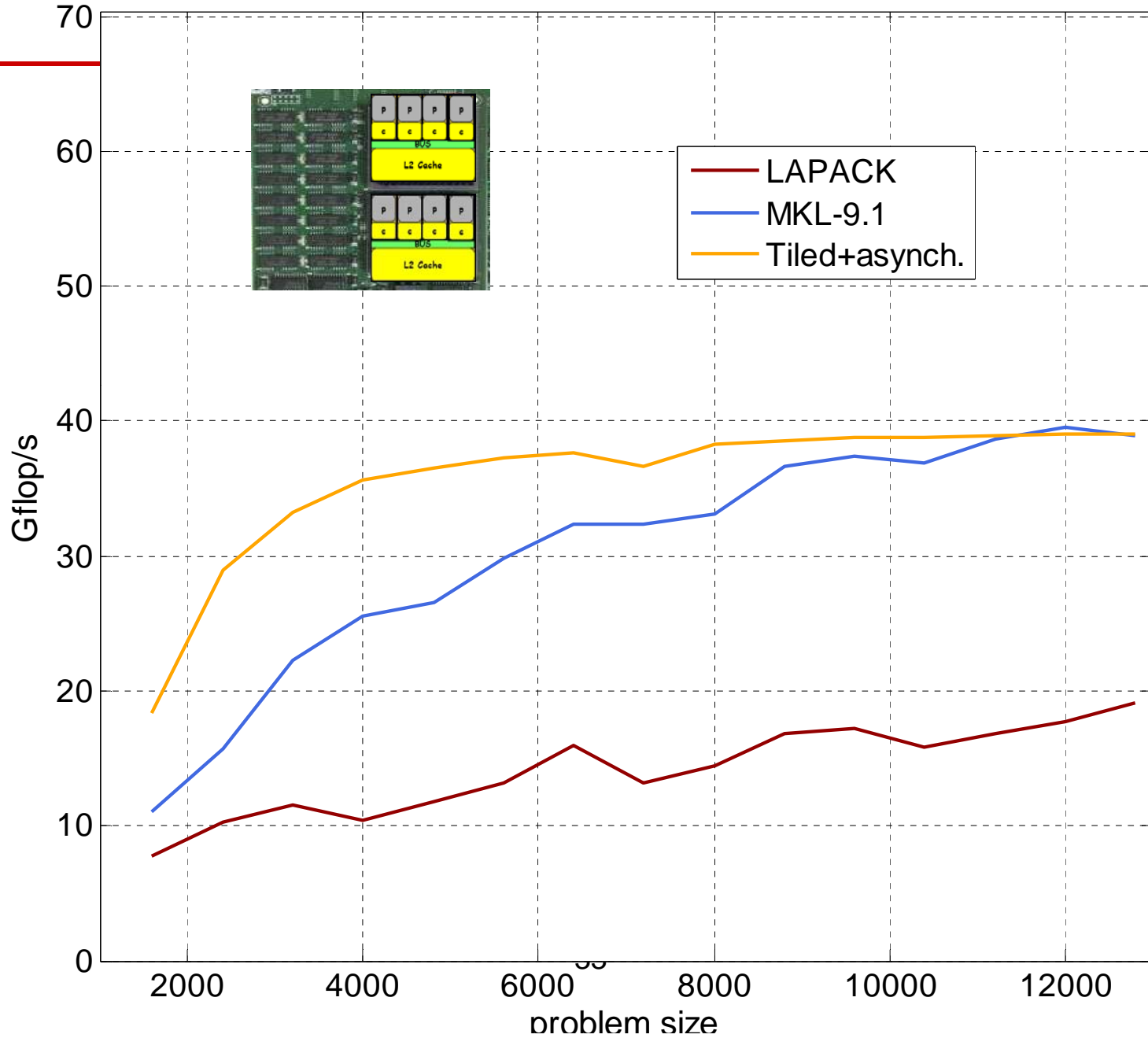
Intel's Clovertown Quad Core

3 Implementations of LU factorization
Quad core w/2 sockets per board, w/ 8 Treads

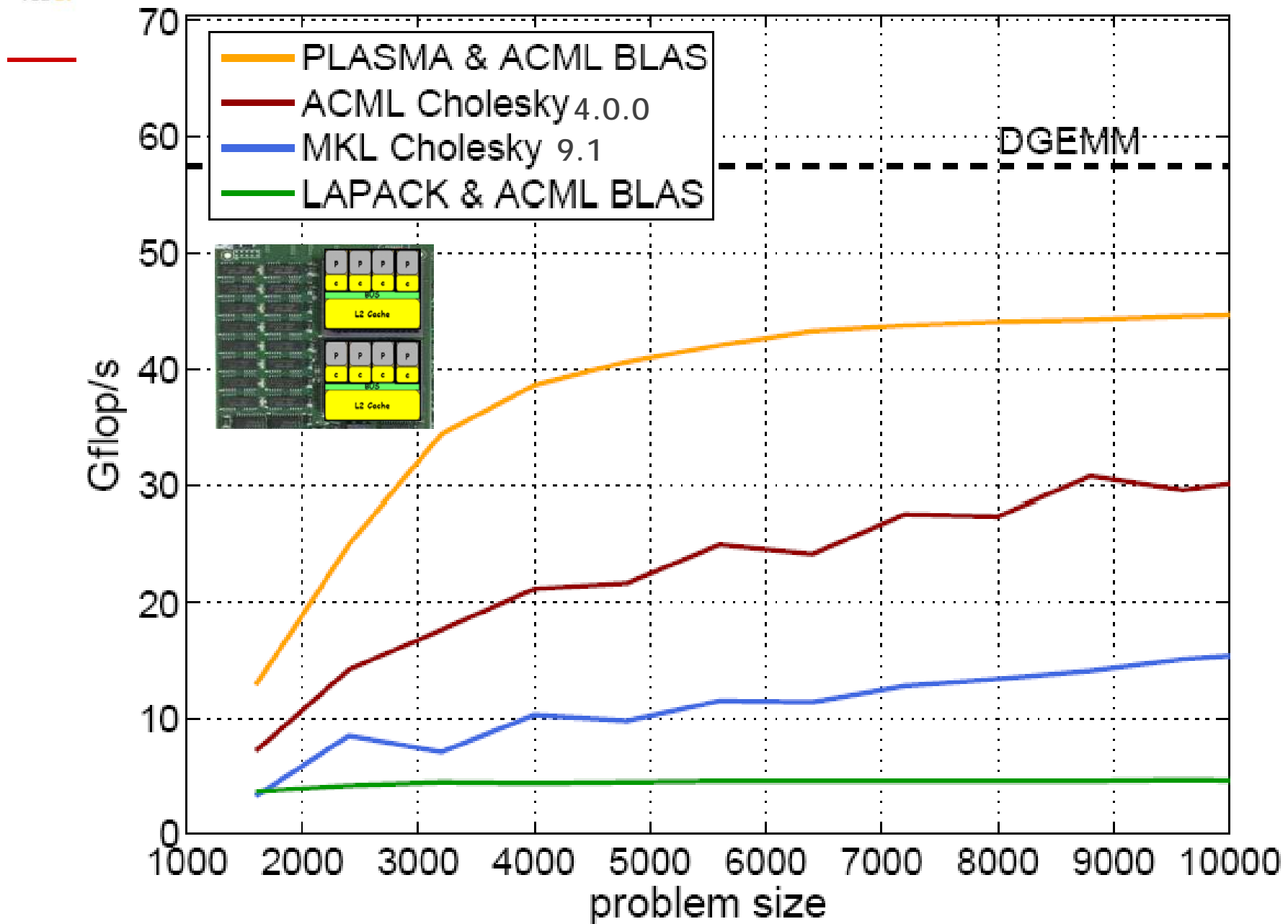
1. LAPACK (BLAS Fork-Join Parallelism)
2. ScaLAPACK (Mess Pass using mem copy)
3. DAG Based (Dynamic Scheduling)



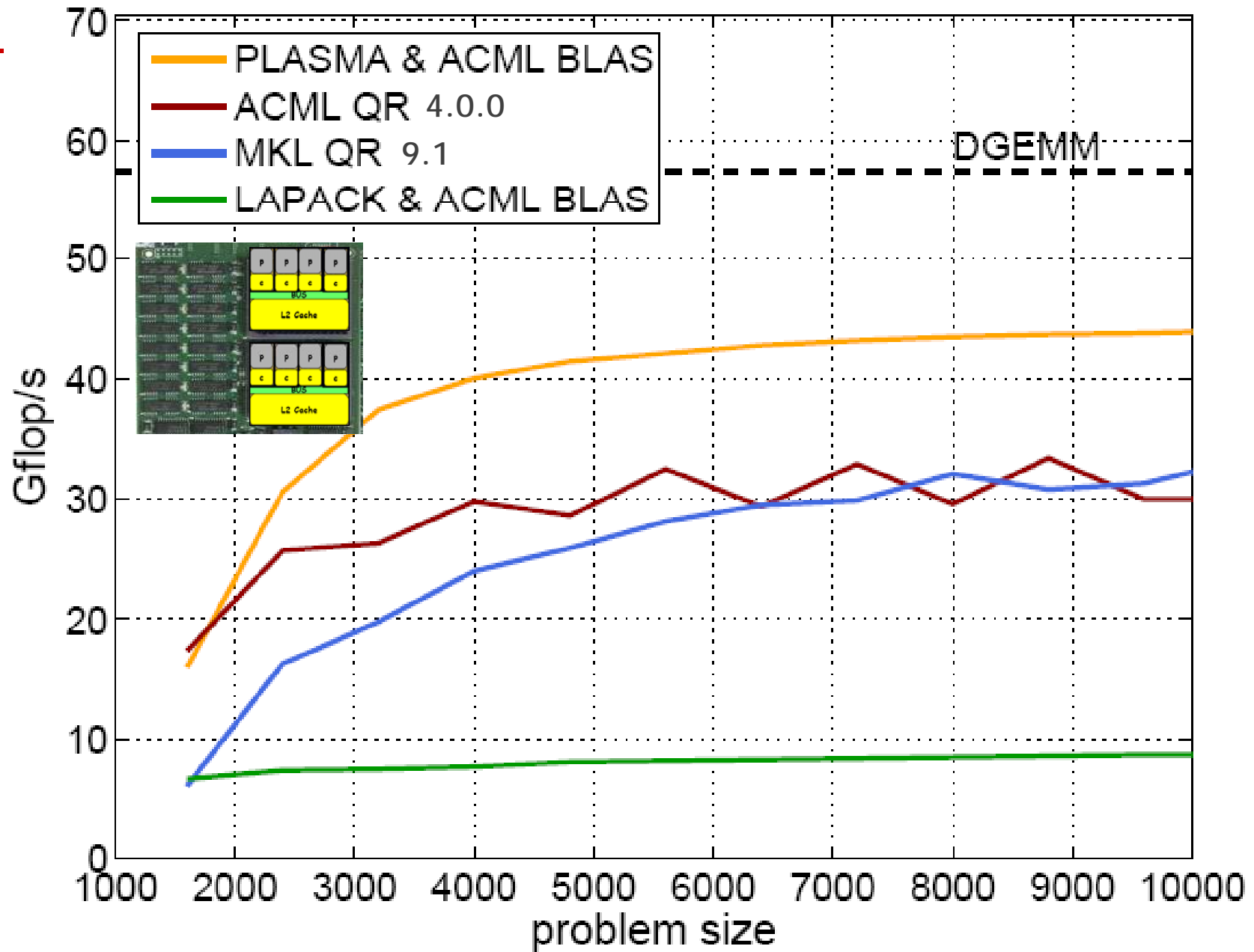
LU -- 8-way dual Opteron -- MKL-9.1



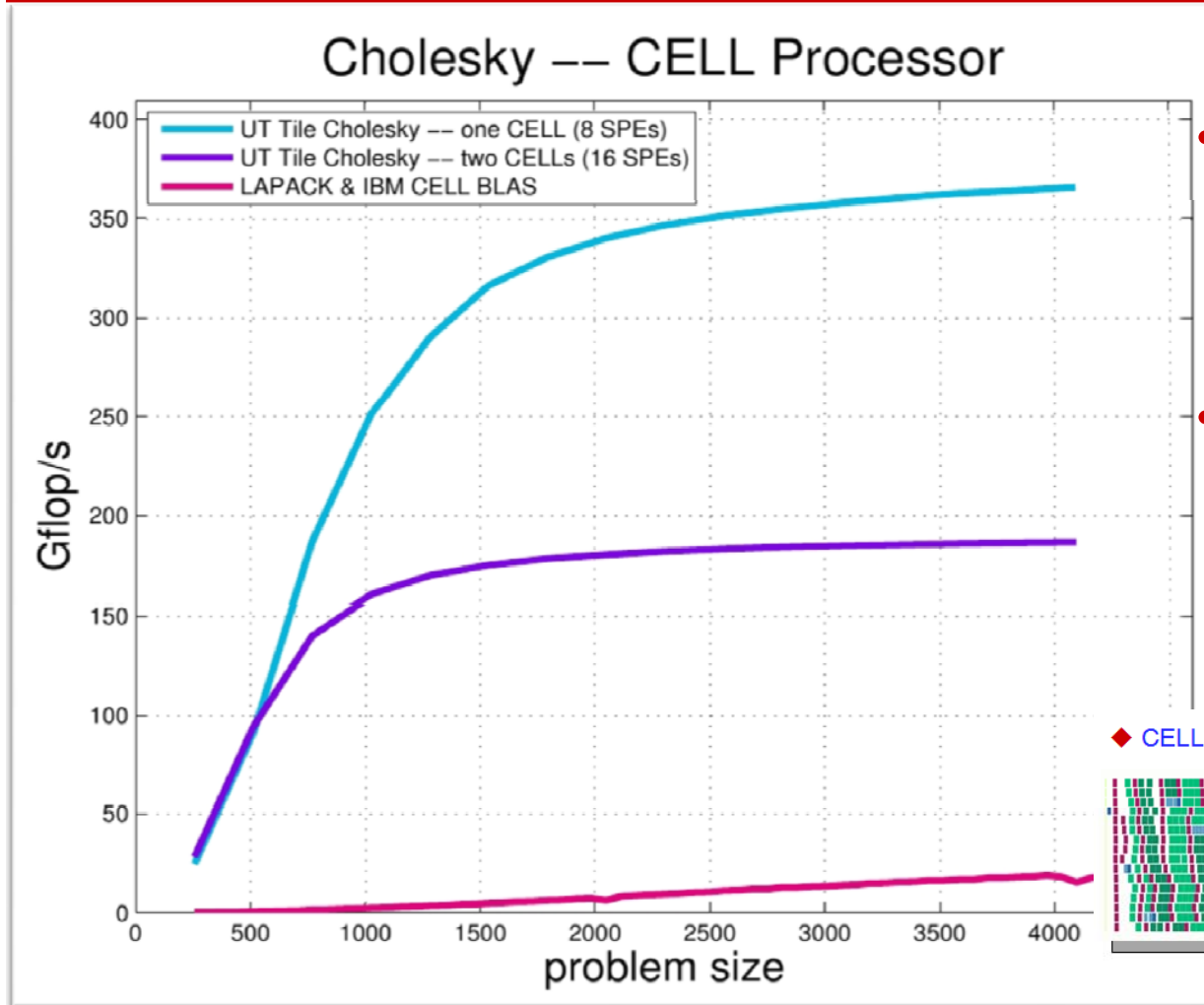
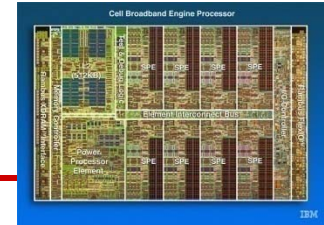
Cholesky -- quad-socket, dual-core Opteron



QR -- quad-socket, dual-core Opteron



Cholesky on the CELL

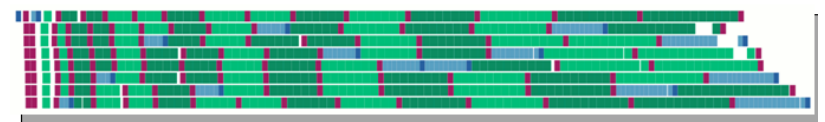


- **1 CELL (8 SPEs)**
 - 186 Gflop/s
 - 91 % peak
 - 97 % SGEMM peak
- **2 CELLS (16 SPEs)**
 - 365 Gflop/s
 - 89 % peak
 - 95 % SGEMM peak

◆ CELL Cholesky - 16 cores

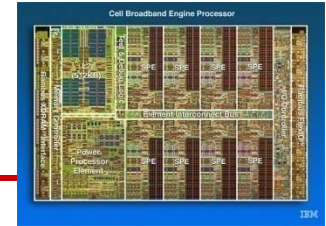


◆ CELL Cholesky - 8 cores



Single precision results on the Cell

Cholesky on the CELL



- 1 CELL (8 SPEs)
 - 186 Gflop/s
 - 91 % peak
 - 97 % SGEMM peak
- 2 CELLS (16 SPEs)
 - 365 Gflop/s
 - 89 % peak
 - 95 % SGEMM peak

◆ CELL Cholesky – 16 cores

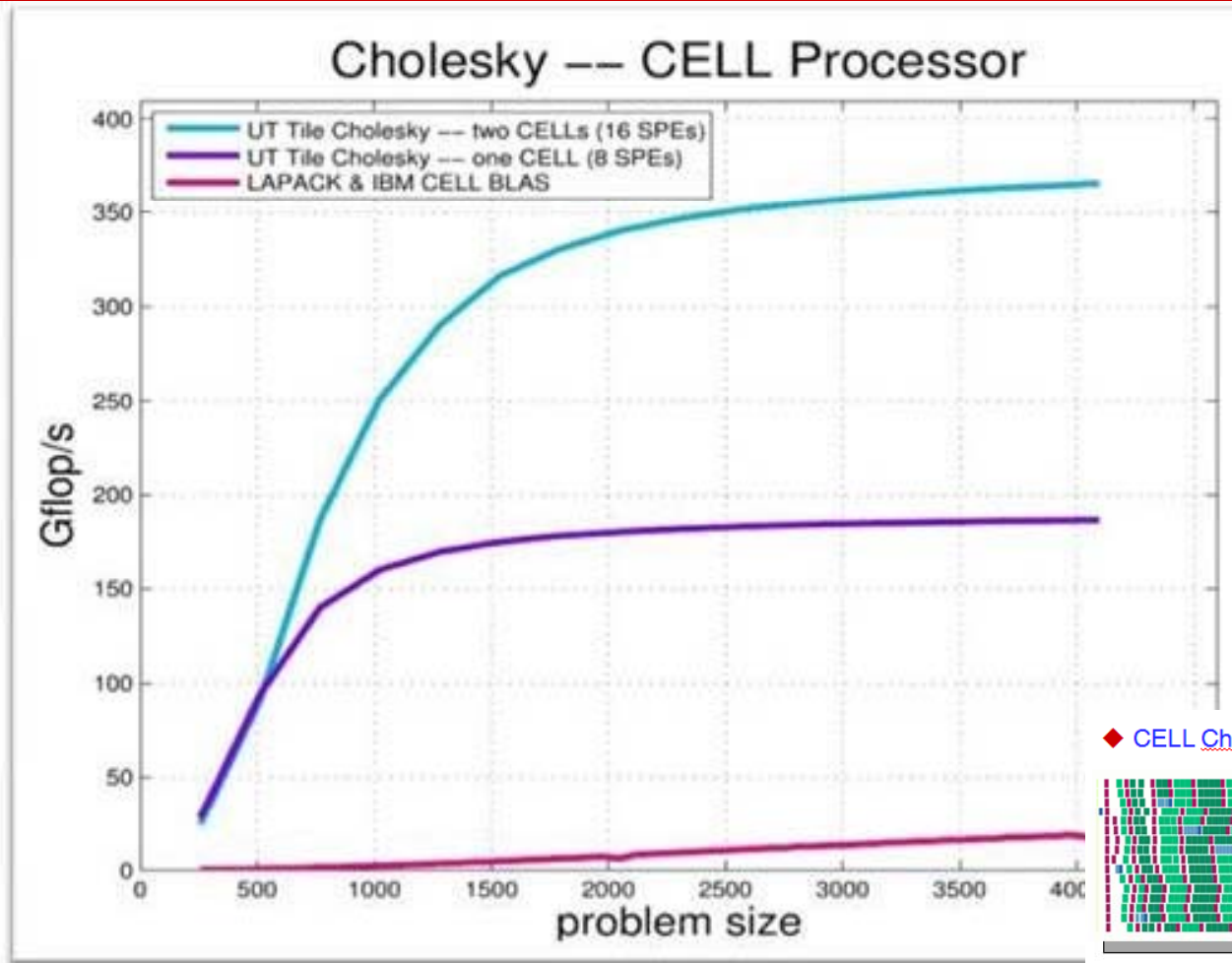
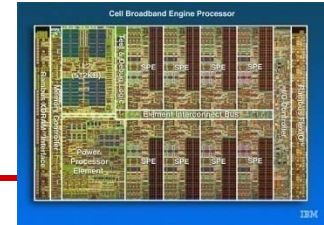


◆ CELL Cholesky – 8 cores



Single precision results on the Cell

Cholesky on the CELL



- **1 CELL (8 SPEs)**
 - 186 Gflop/s
 - 91 % peak
 - 97 % SGEMM peak
- **2 CELLS (16 SPEs)**
 - 365 Gflop/s
 - 89 % peak
 - 95 % SGEMM peak

◆ CELL Cholesky - 16 cores



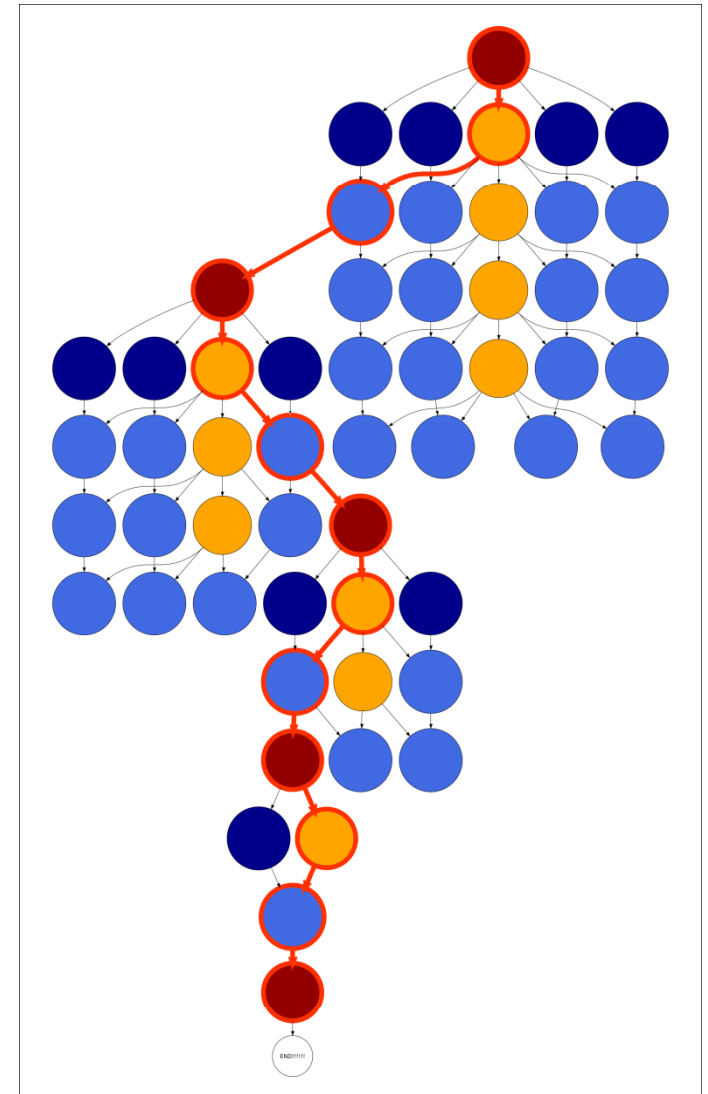
◆ CELL Cholesky - 8 cores



Single precision results on the Cell

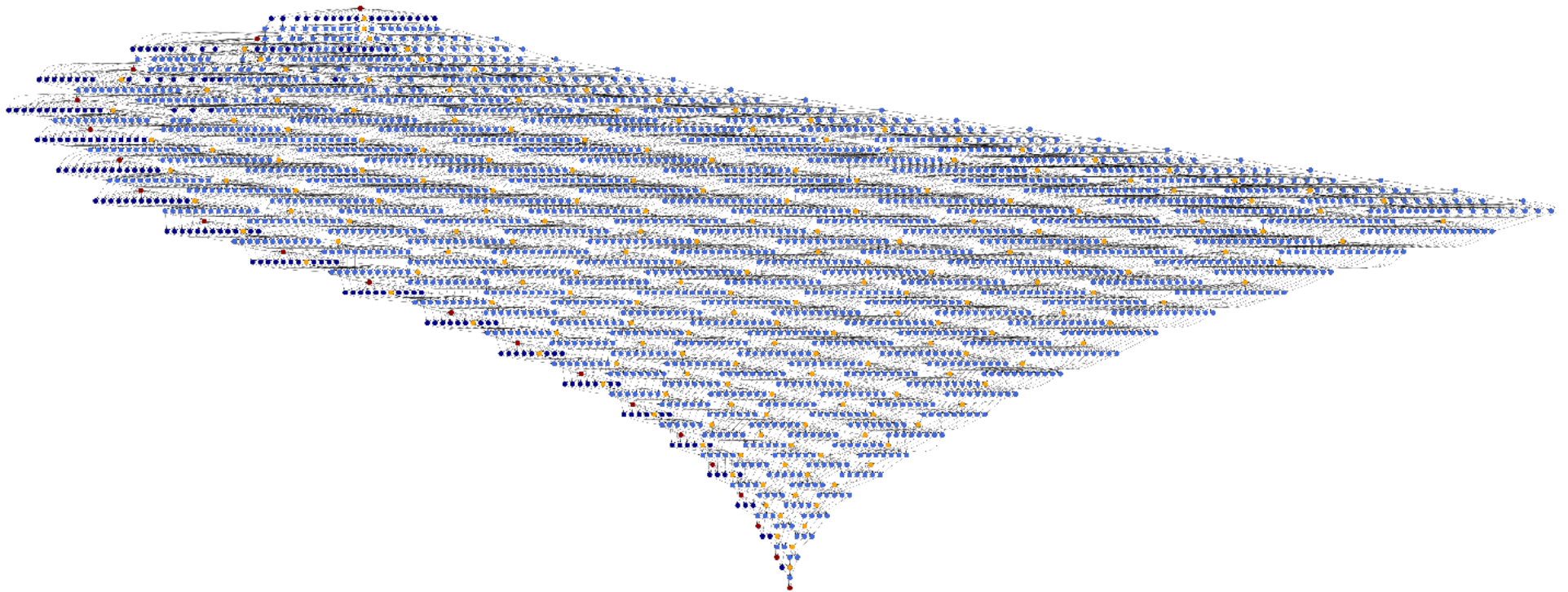
If We Had A Small Matrix Problem

- We would generate the DAG, find the critical path and execute it.
- DAG too large to generate ahead of time
 - Not explicitly generate
 - Dynamically generate the DAG as we go
- Machines will have large number of cores in a distributed fashion
 - Will have to engage in message passing
 - Distributed management
 - Locally have a run time system



The DAGs are Large

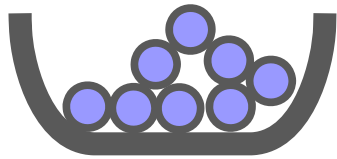
- Here is the DAG for a factorization on a 20 x 20 matrix



- For a large matrix say $O(10^6)$ the DAG is huge
- Many challenges for the software

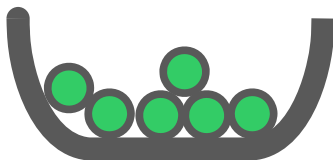


Each Node or Core Will Have A Run Time System



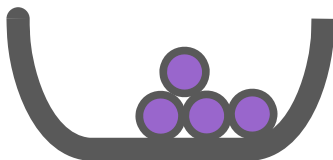
BIN 1

- ◆ some dependencies satisfied
- ◆ waiting for all dependencies



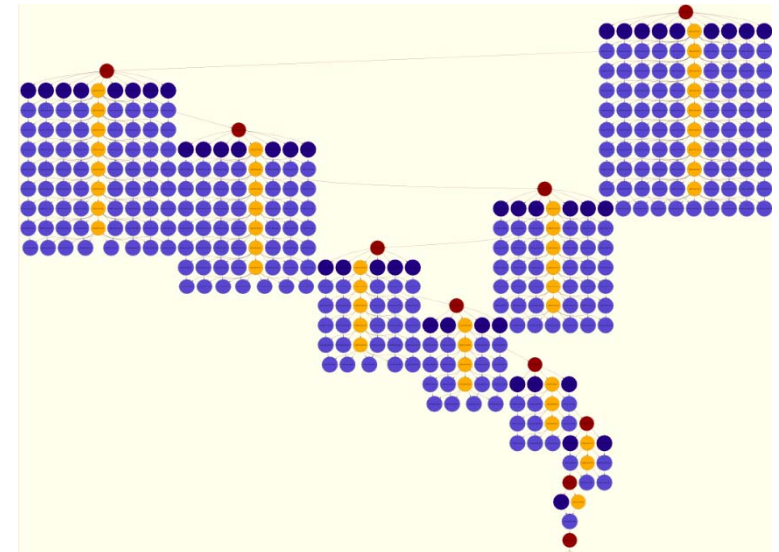
BIN 2

- ◆ all dependencies satisfied
- ◆ some data delivered
- ◆ waiting for all data



BIN 3

- ◆ all data delivered
- ◆ waiting for execution





DAG and Scheduling

- DAG is dynamically generated and implicit
- Everything designed for distributed memory systems
- Runtime system on each node or core
- Run time
- Bin 1
 - See if new data has arrived
- Bin 2
 - See if new dependences are satisfied
 - If so move task to Bin 3
- Bin 3
 - Exec a task that's ready
 - Notify children of completion
 - Send data to children
 - If no work do work stealing



Some Questions

- What's the best way to represent the DAG?
- What's the best approach to dynamically generating the DAG?
- What run time system should we use?
 - We will probably build something that we would target to the underlying system's RTS.
- What about work stealing?
 - Can we do better than nearest neighbor work stealing?
- What does the program look like?
 - Experimenting with Cilk, Charm++, UPC, Intel Threads
 - I would like to reuse as much of the existing software as possible



PLASMA Collaborators

- **U Tennessee, Knoxville**
 - Jack Dongarra, Julie Langou, Stan Tomov, Jakub Kurzak, Hatem Ltaief, Alfredo Buttari, Julien Langou, Piotr Luszczek, Marc Baboulin
- **UC Berkeley**
 - Jim Demmel, Ming Gu, W. Kahan, Beresford Parlett, Xiaoye Li, Osni Marques, Yozo Hida, Jason Riedy, Vasily Volkov, Christof Voemel, David Bindel
- **Other Academic Institutions**
 - UC Davis, CU Denver, Florida IT, Georgia Tech, U Maryland, North Carolina SU, UC Santa Barbara, UT Austin, LBNL
 - TU Berlin, ETH, U Electrocomm. (Japan), FU Hagen, U Carlos III Madrid, U Manchester, U Umeå, U Wuppertal, U Zagreb, UPC Barcelona, ENS Lyon, INRIA
- **Industrial Partners**
 - Cray, HP, Intel, Interactive Supercomputing, MathWorks, NAG, NVIDIA, Microsoft