**Jack Dongarra, Mathieu Faverge, Hatem Ltaief, Piotr Luszczek**

# High Performance Matrix Inversion Based on LU Factorization for Multicore Architectures

**presented by**

**Piotr Luszczek**

ICL UT

# Preliminaries

# Problem Statement

$A \in \boldsymbol{R}^{n \times n}$

$\quad \Rightarrow \quad PA = LU$

$\quad \Rightarrow \quad U \rightarrow U^{-1}$

$\quad \Rightarrow \quad L \rightarrow L^{-1}$

$\quad \Rightarrow \quad A^{-1} \in \boldsymbol{R}^{n \times n}$
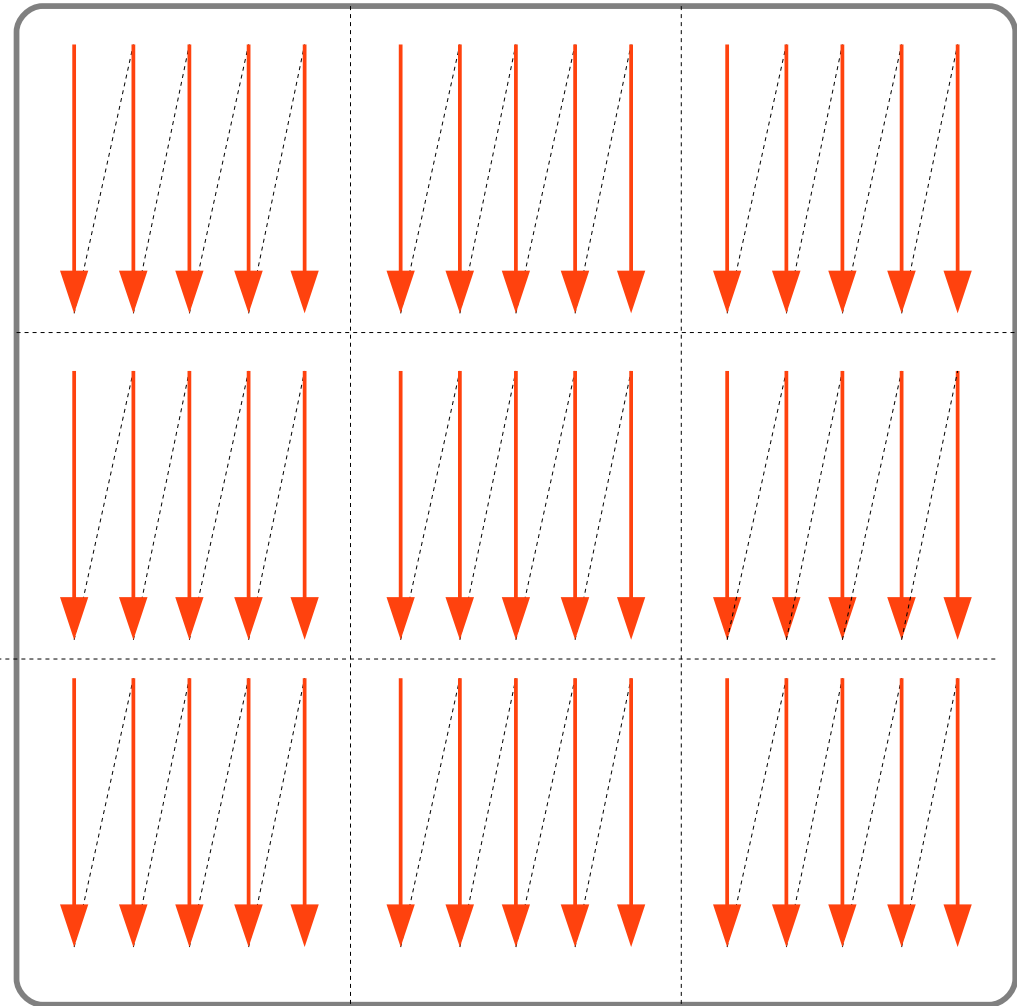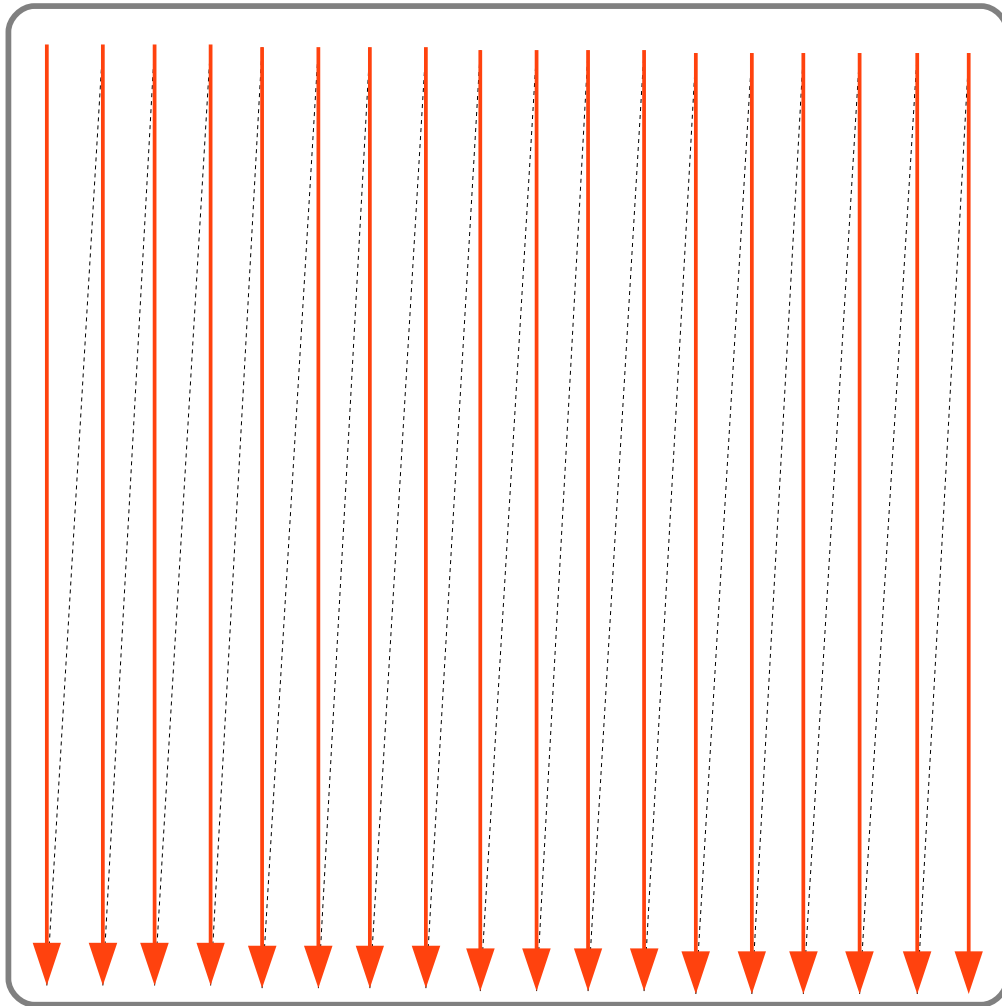
# To Keep in Mind...

In the vast majority
of
practical computational problems,
it is unnecessary
and
inadvisable
to
actually compute $A^{-1}$.

*Forsythe, Malcolm, and Moler*

# Data Layouts for Matrix Elements

**Column-major (LAPACK and derivatives)**

**Tile (PLASMA)**

# Tasks and DAGs

# Block LU Inversion

- For each panel *LU factorization*
  DGETF2( )
  DLASWP( )
  DLASWP( )
  DTRSM( )
  DGEMM( )

- For each panel *Invert U*
  DTRMM( )
  DTRSM( )
  DTRTI2( )

- For each panel *Invert L*
  DLACPY( )
  DLASET( )
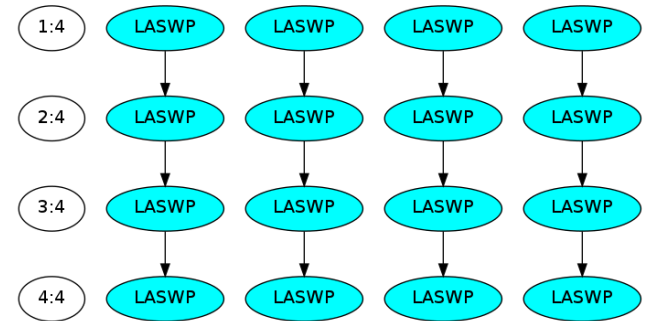  DGEMM( )
  DTRSM( )

- DLASWP( ) *column interchanges*

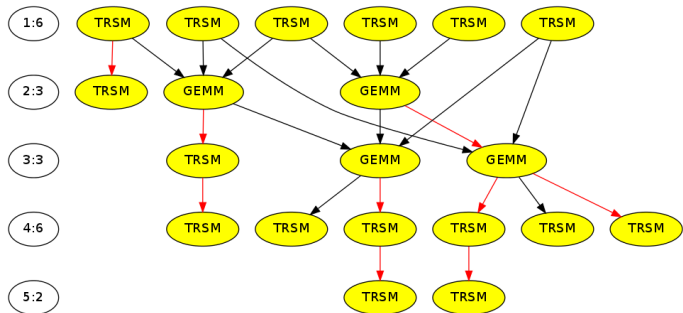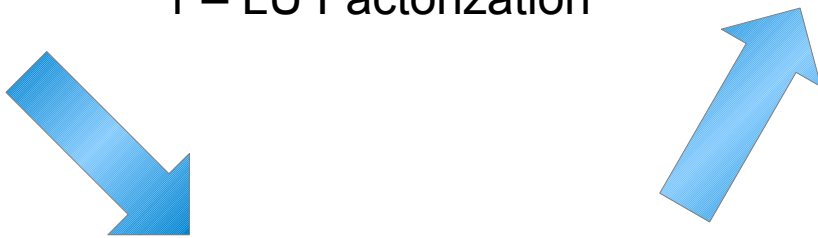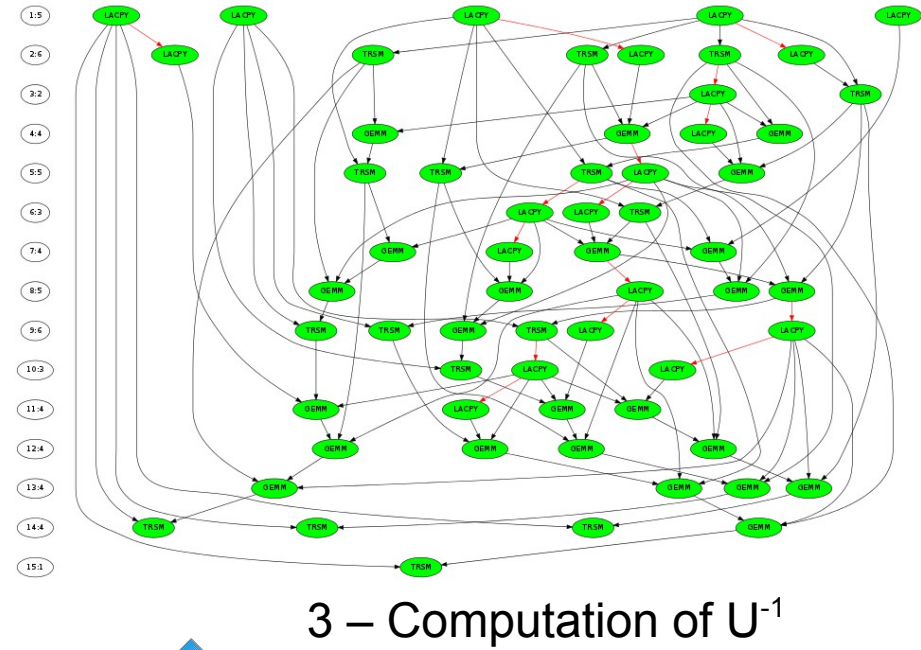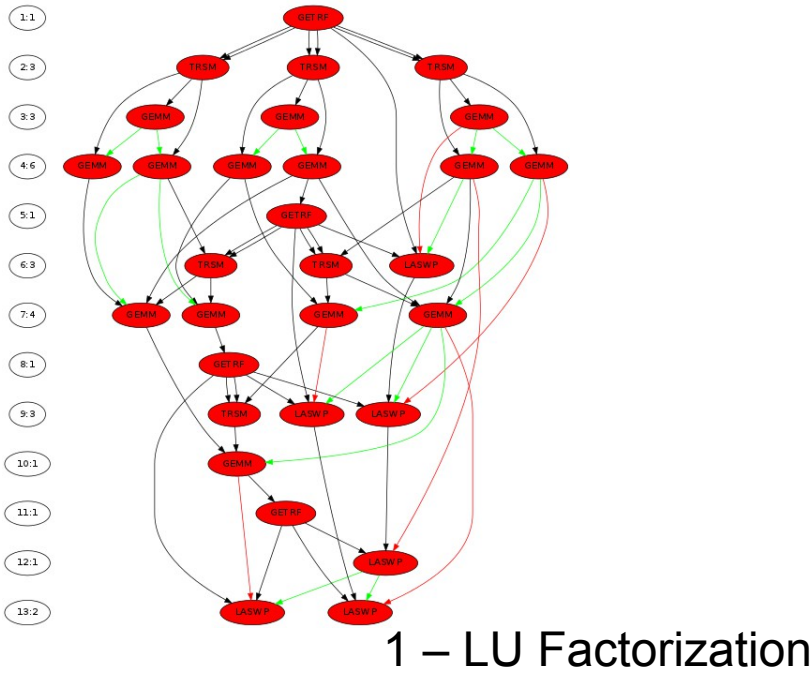# Tile LU Inversion

- For each diagonal tile
  -DGETRFR()*parallel recursive LU*
  for each tail tile panel
    -DLASWP( )
  for each tail tile
    -DGEMM( )
  for each left tile panel
    -DLASWP( )

- For each diagonal tile *Invert U*
  for each tile in panel
    -DTRSM( )
  for each tail tile
    -DGEMM( )
  for each left panel tile
    -DTRSM( )
  -DTRTRI( )

- For each left tile *Invert L*
  -DLACPY( )
  -DLASET( )
  ...

ICL UT

# Queuing Functions with QUARK

```
QUARK_Insert_Task(

    panel_LU_task,

    M, matrix_1 , INPUT,

    N, matrix_2 , INOUT,

    1,  result  , OUTPUT,

    K, buffer  , SCRATCH,

        0);
```
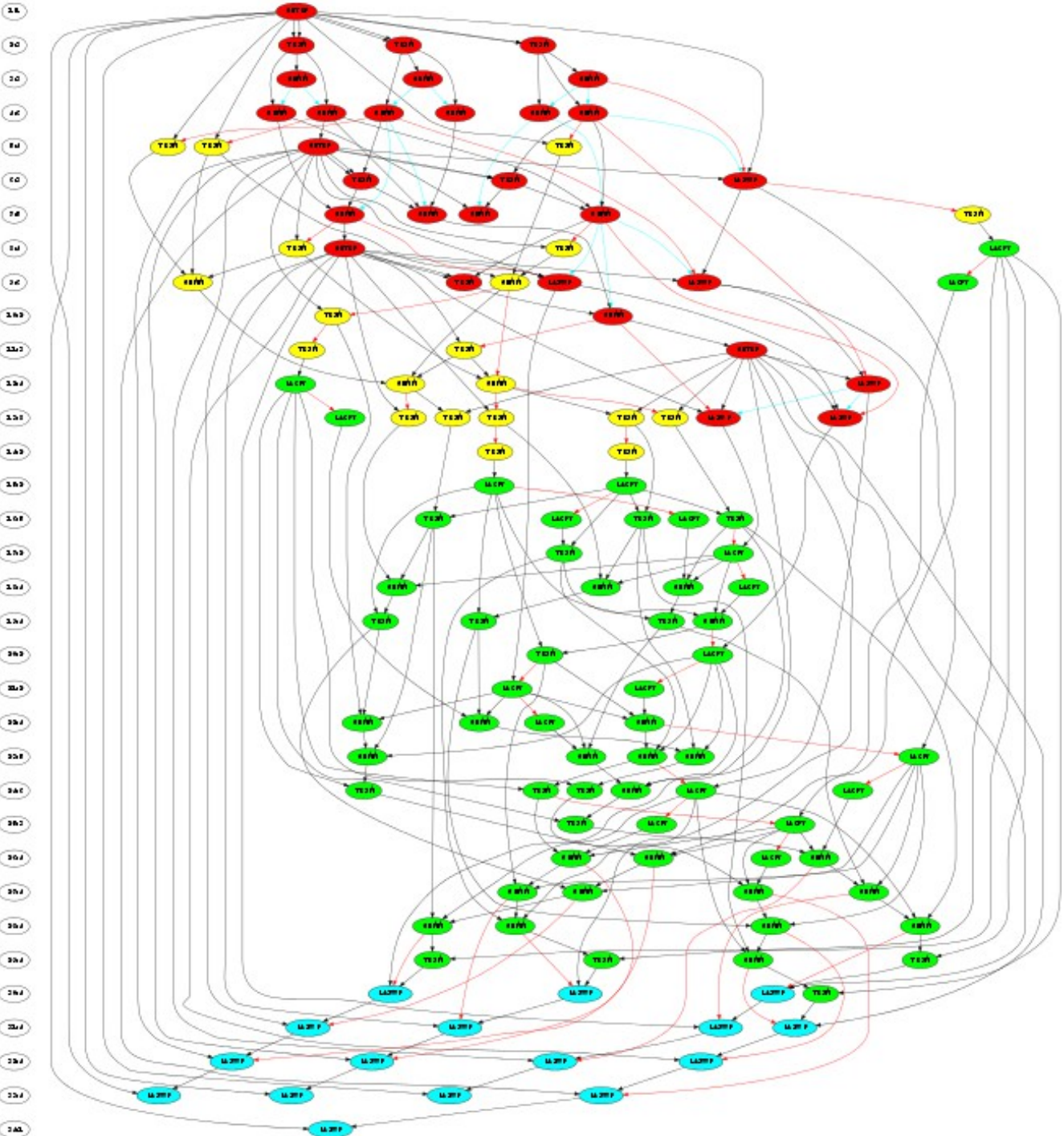
# DAGs of Tasks, Each State Separately
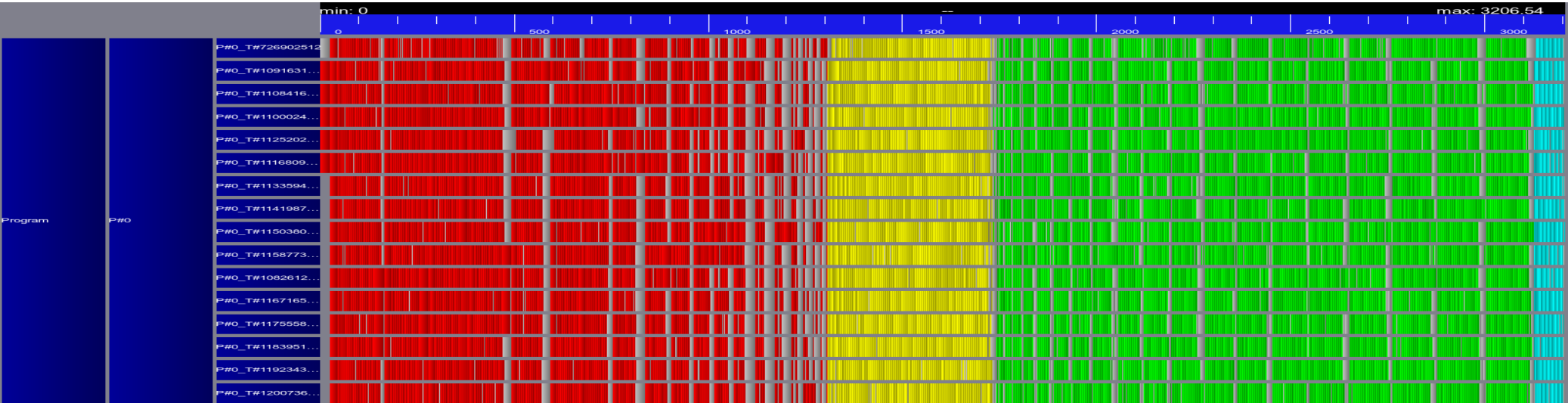


1 – LU Factorization

2 – Computation of L$^{-1}$

3 – Computation of U$^{-1}$

4 – Column swapping

# DAGs of Tasks, All Stages Overlapped
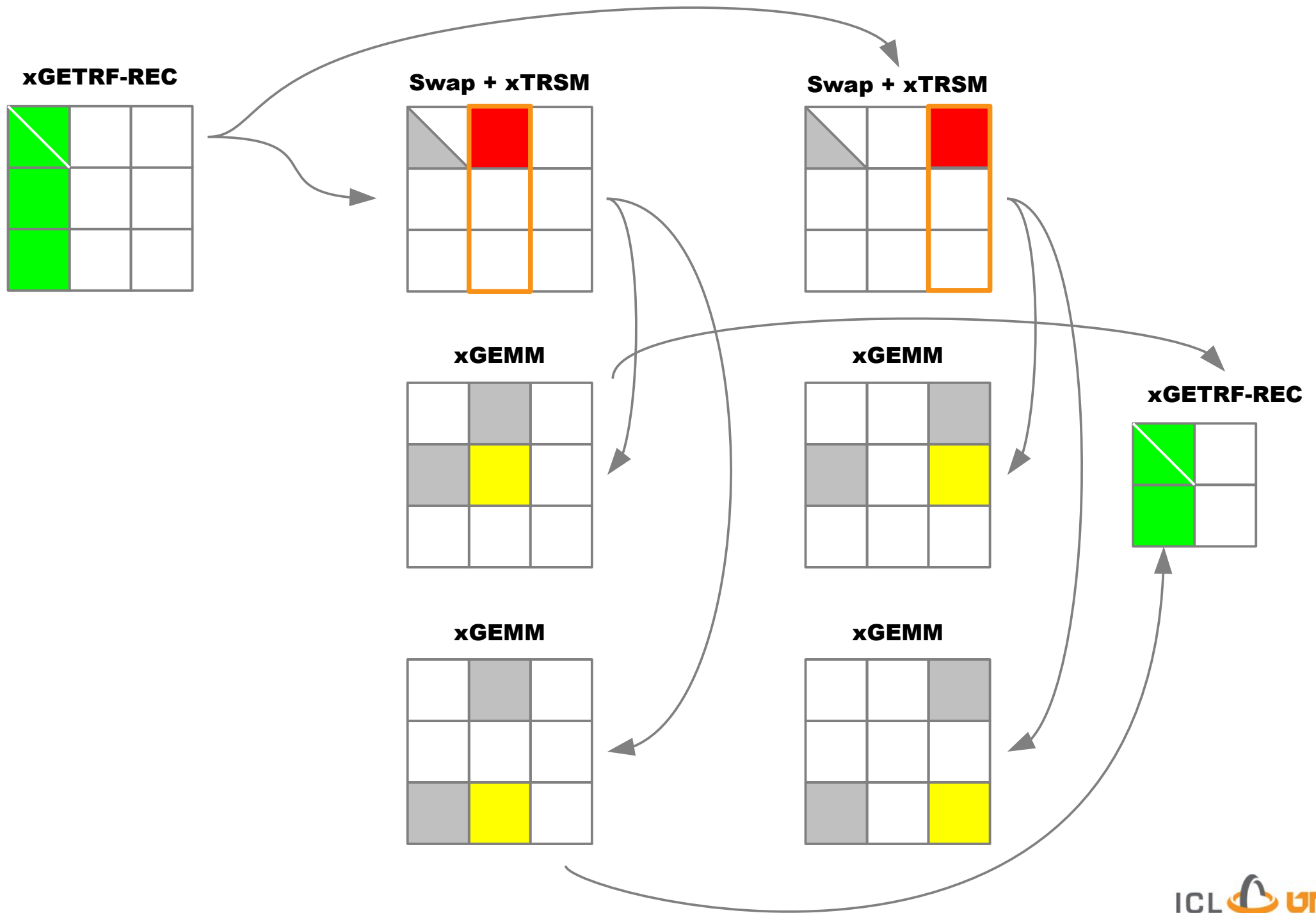
# Execution Traces
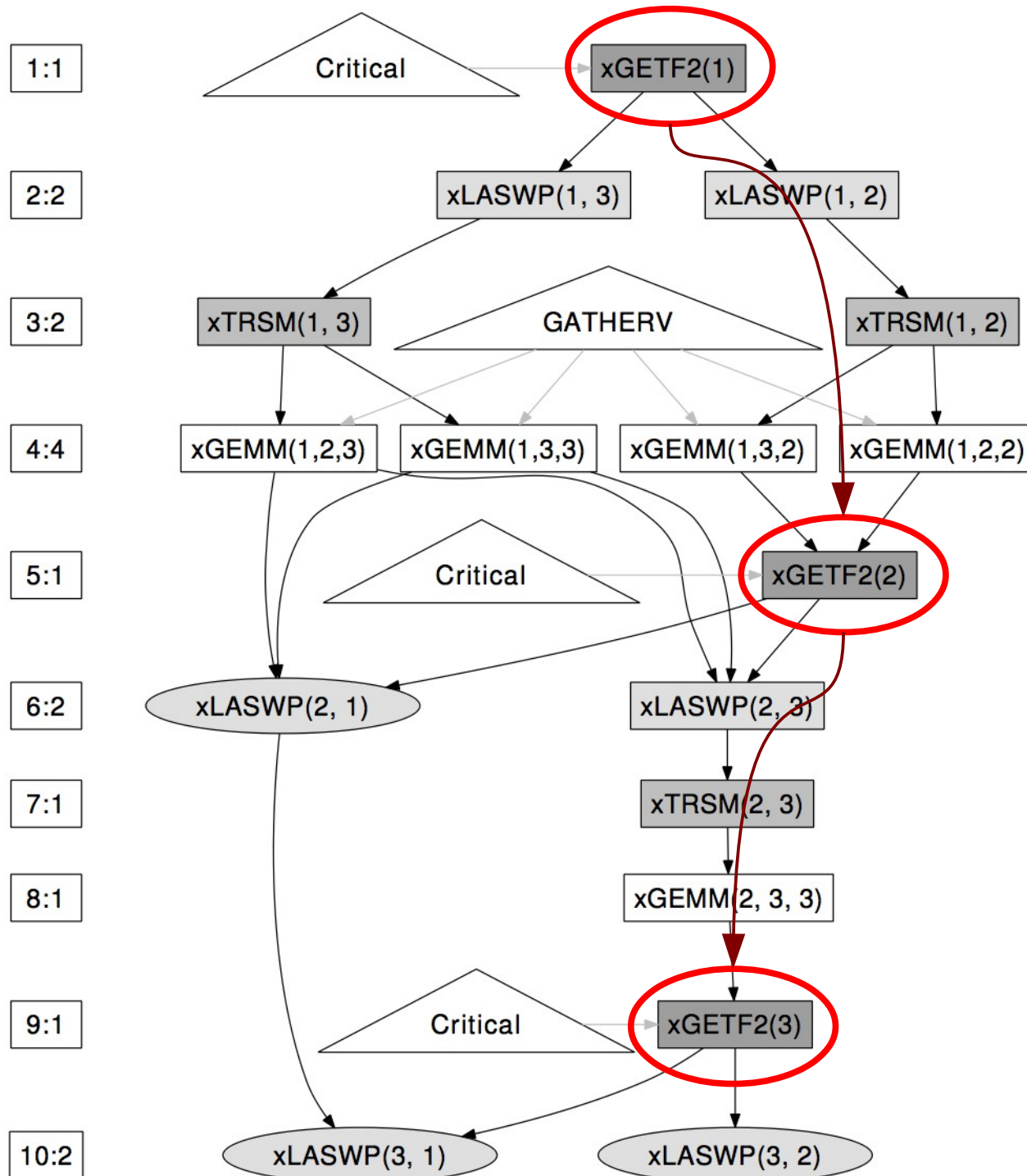
## No Overlap of Stages



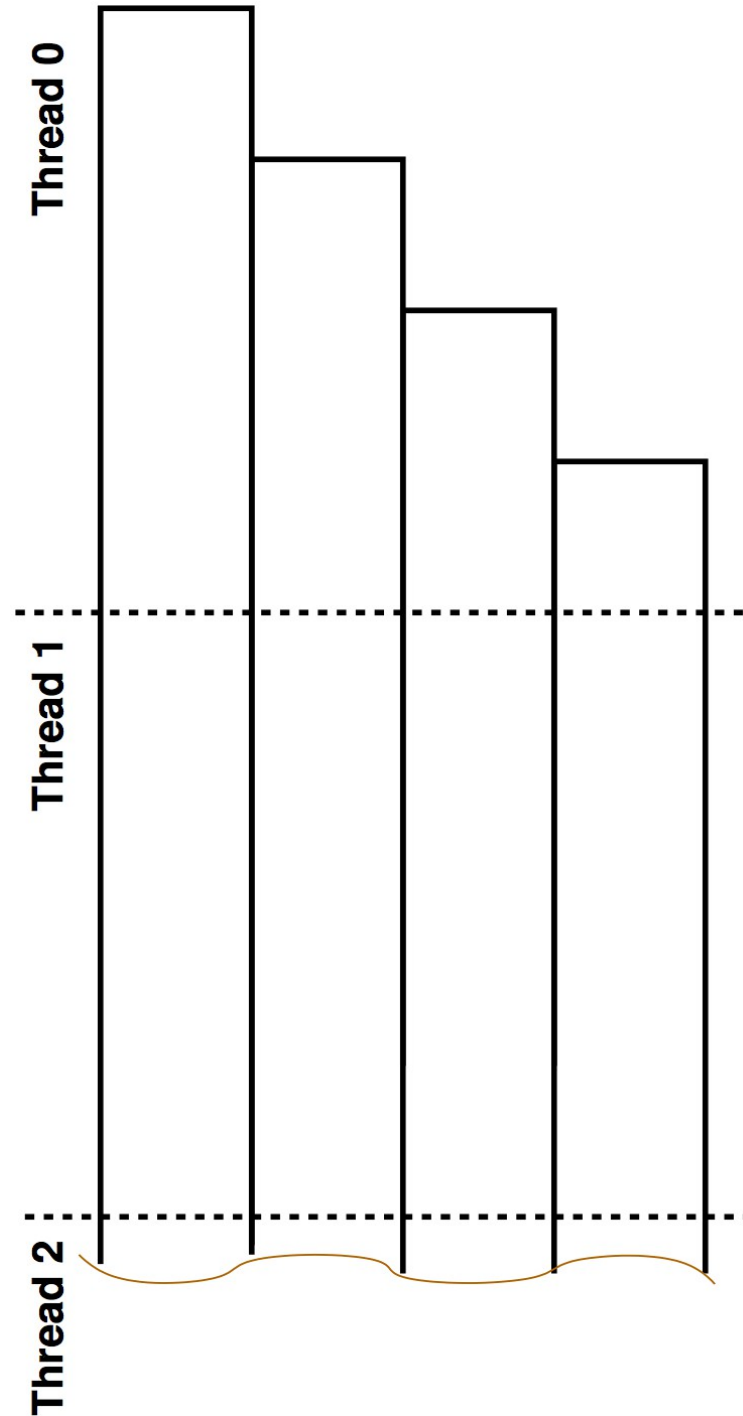## Overlap of Stages

# The Case for Nested Parallelism

# Panel Factorization as the Sequential Bottleneck

# Panel Factorization is On Critical Path of DAG

# Parallel Panel Factorization: Data Partitioning

# Parallel Panel Factorization: Algorithm

```
function xGETRFR(M, N, column) {
    if N == 1 {

        idx = split_IxAMAX(...)

        gidx = combine_IxAMAX(idx)

        split_xSCAL(...)
    } else {
        xGETRFR(M, N/2, column)

        xLASWP(...)

        split_xTRSM(...)

        split_xGEMM(...)

        xGETRFR(M, N-N/2, column+N/2)

        xLASWP(...)
    }
}
```

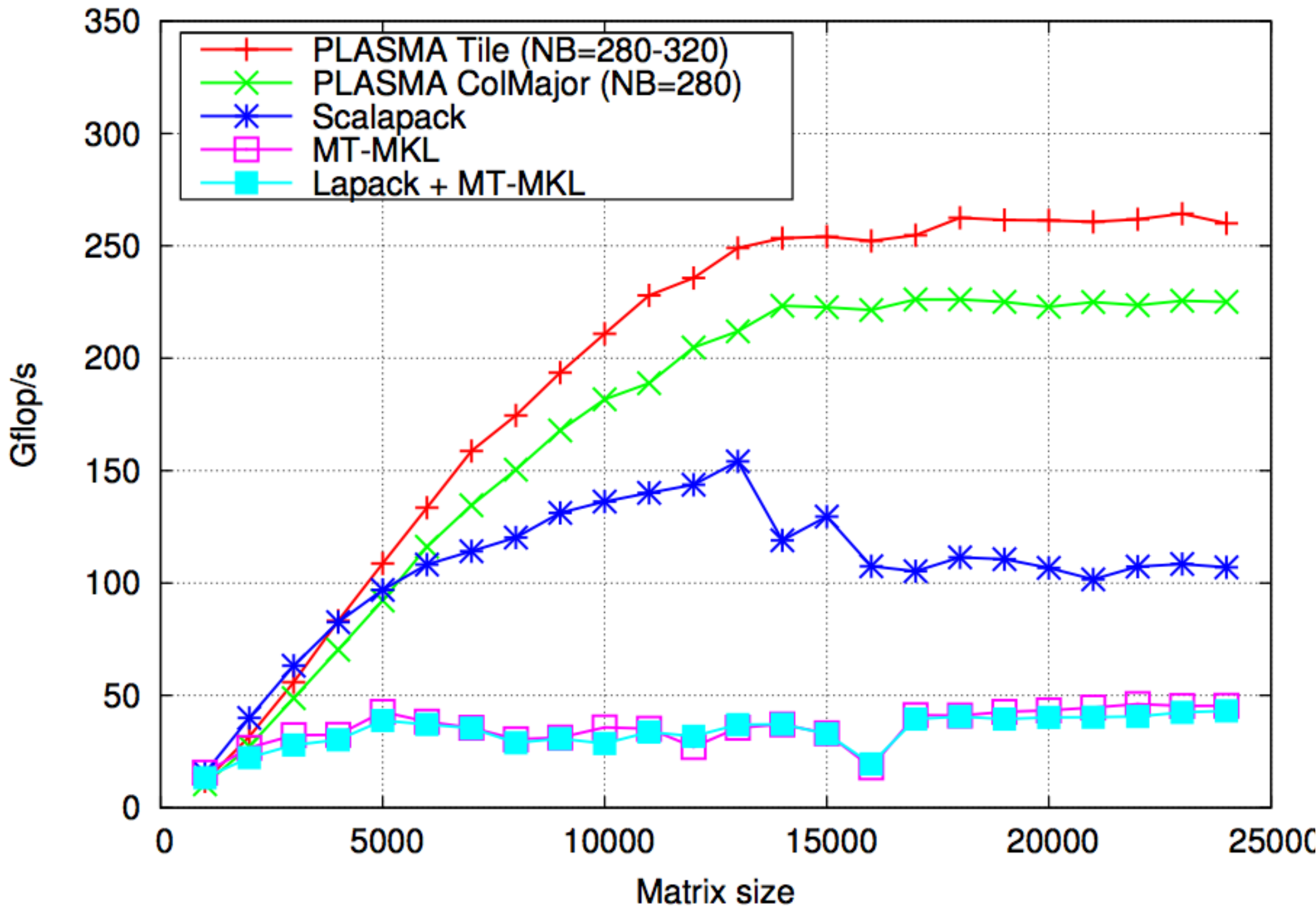| | |
|---|---|
| *single column, recursion stops* | |
| *compute local maximum of modulus* | |
| *combine local results* | |
| *scale local data* | |
| *recursive call to factor left half* | |
| *pivoting forward* | |
| *triangular solve* | |
| *Schur's complement* | |
| *recursive call to factor right half* | |
| *pivoting backward* | |

ICL

# Quick Performance Experiment

# Results

# Performance on AMD MagnyCours, 4x12=48 cores



Legend:
- PLASMA Tile (NB=280-320)
- PLASMA ColMajor (NB=280)
- Scalapack
- MT-MKL
- Lapack + MT-MKL
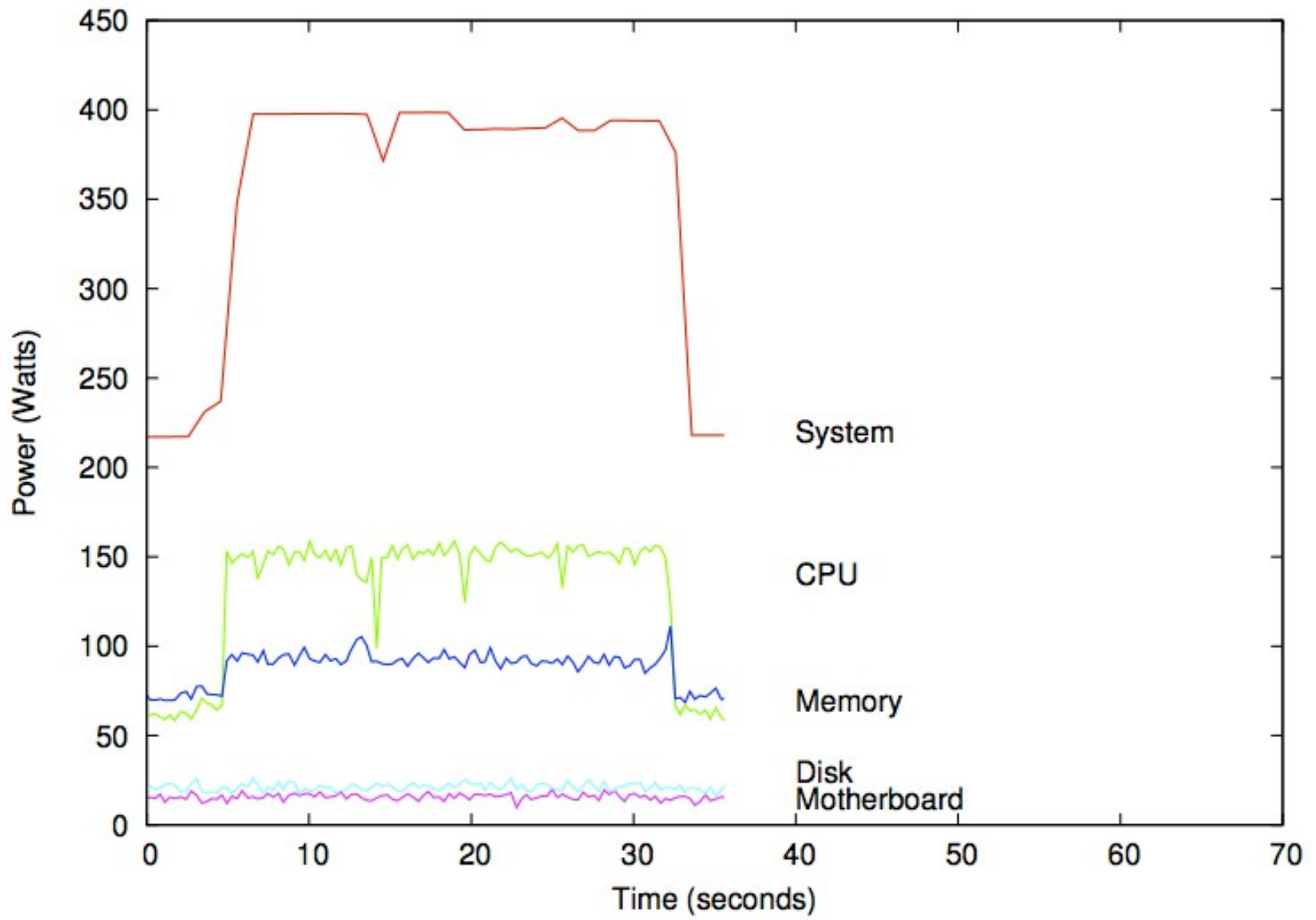
Y-axis: Gflop/s
X-axis: Matrix size

# LU Inversion's Power Profile: LAPACK

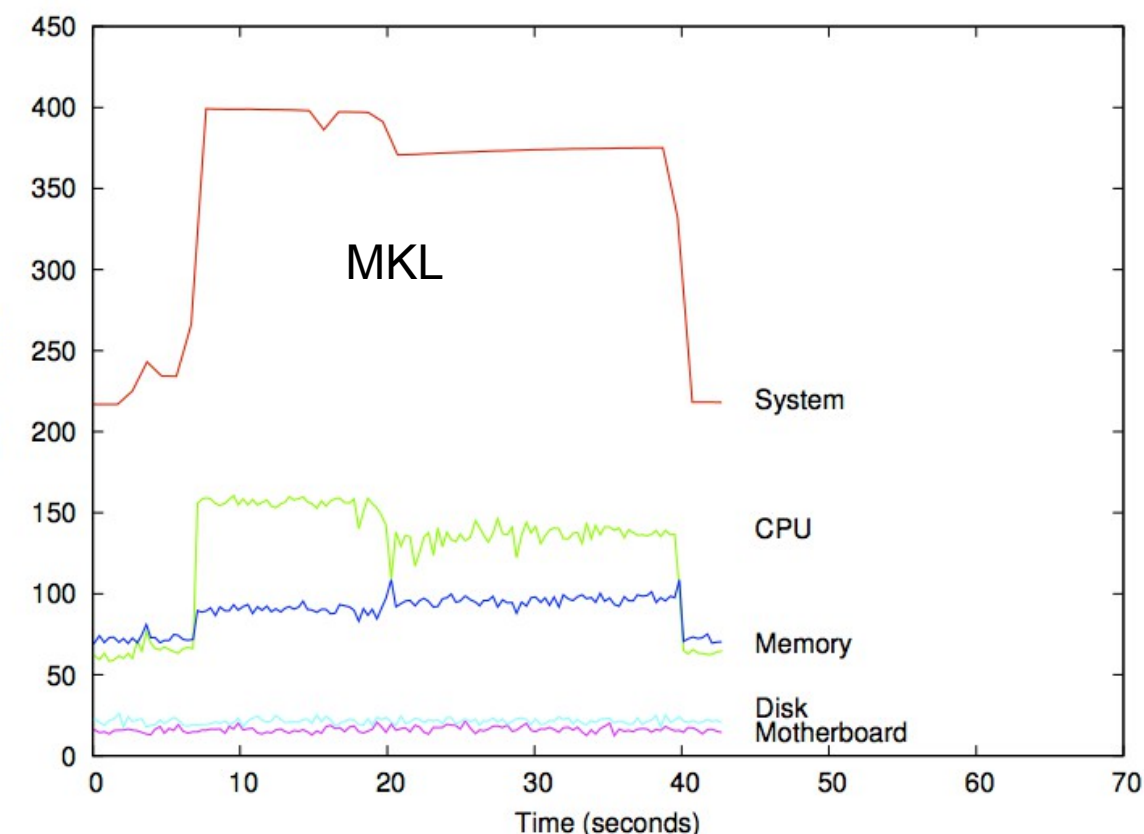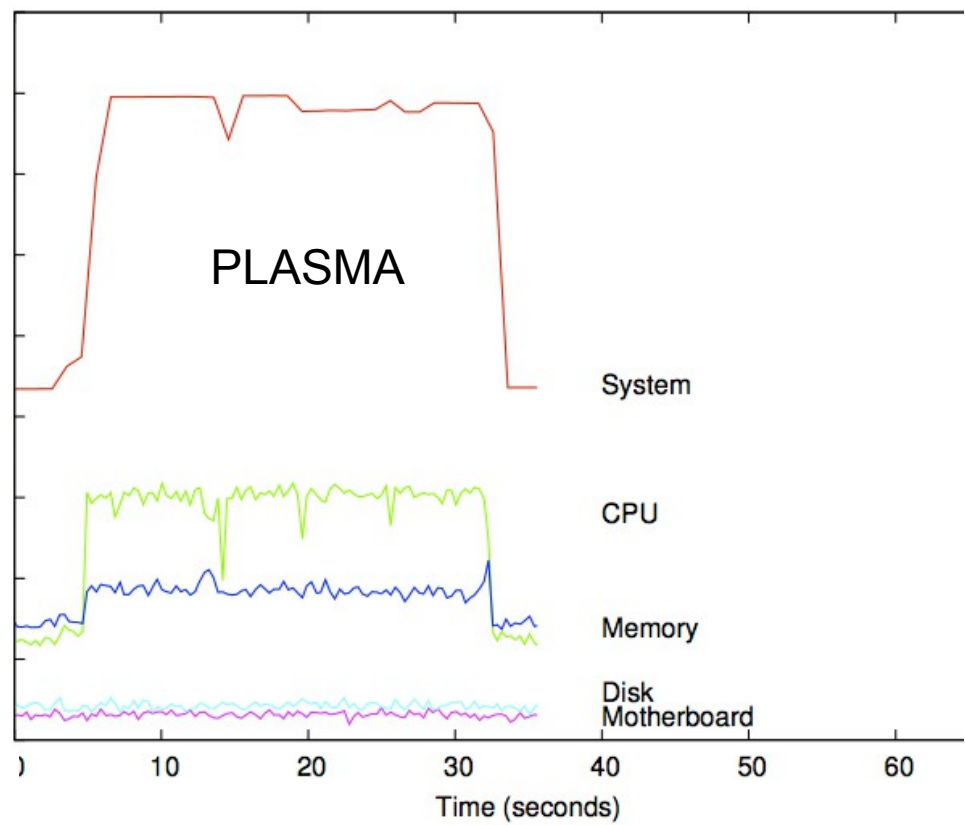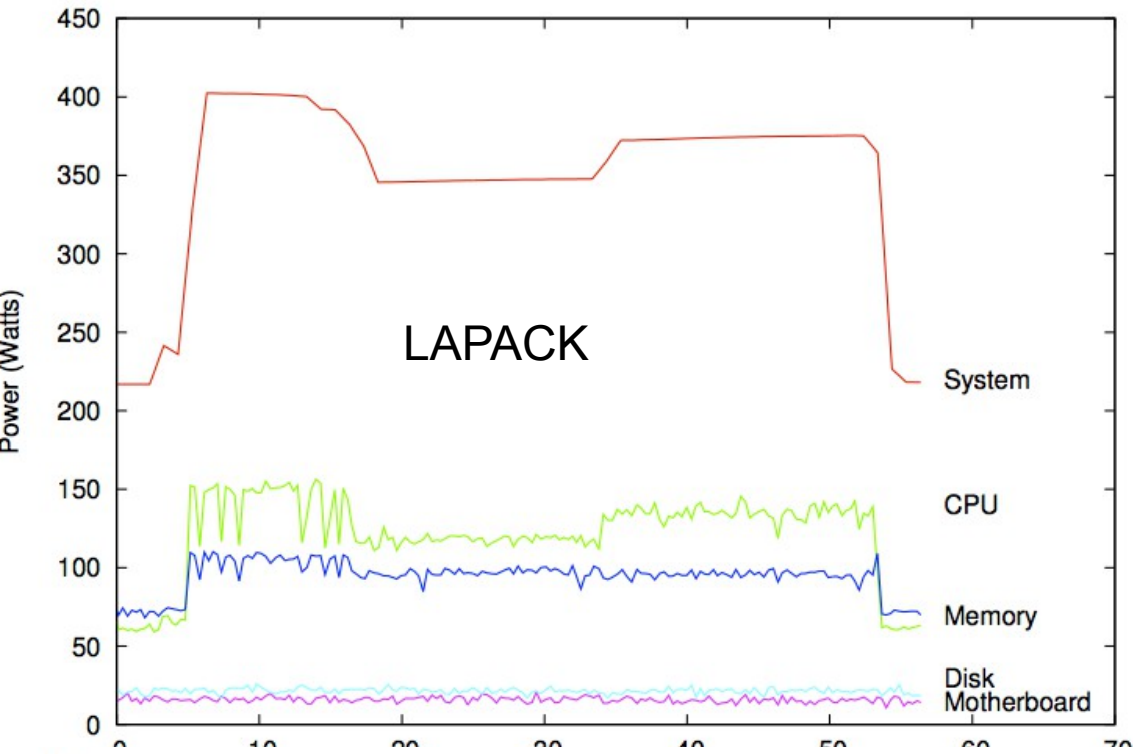# LU Inversion's Power Profile: MKL

# LU Inversion's Power Profile: PLASMA

This work was sponsored
by
NSF, DOE, and Microsoft