

Modeling the LU Factorization for SMP Clusters

Jack Dongarra **Emmanuel Jeannot** Julien Langou

University of Tennessee
LORIA-INRIA Lorraine
University of Colorado

Outline of the talk

- 1 Introduction
- 2 Building the model
- 3 Computing the model parameters
- 4 experimental results
- 5 Conclusion

Introduction: LU Factorization

- A : a (squared, no-singular) matrix,
- LU factorization computes $A = P \cdot L \cdot U$ where:
 - L is a lower triangular unit matrix,
 - U is a upper triangular matrix,
 - P is a permutation matrix.
- Used to solve the problem $Ax = b$ where A and b are given (by performing two successive triangular solves on U then L).

- A : dense matrix of size $N \times N$,
- parallel factorization,
- 2-D **block-cyclic** distribution of the LU factorization as it is implemented in ScaLapack,
- $P \times Q = NP$ processor grid,
- block are squared of size NB .

Introduction: previous work

Scalapack Users Guide:

- crude model,
- shows the scalability of ScaLapack,
- Do not use the processor grid shape.

Lapack working note 43 (Dongarra, van de Geijn, Walker):

- simple model with 3 parameters:
 - the bandwidth of the network,
 - the latency of the network,
 - and the flop rate of the processor
- the value of these parameters does not depend on the subroutine of the algorithm.

Domas, Desprez, Tourancheau 96:

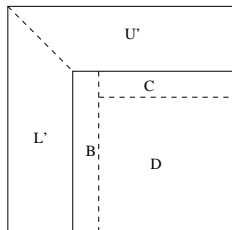
- Very fine model with more than 30 parameters.
- Do not show how to compute the parameters.
- Able compute the optimal grid size.

A new model:

- A detailed model :
 - take into account all the phenomena
 - derive a model for each step of the factorization and for the whole factorization.
- A fine parameterization: different parameters are derived for each subroutines used.
- An automatic parameter computing. We show how to simply compute these parameters.
- An Automatic grid shape computing.
- Enhancement to SMP clusters.

Partial factorization

After a partial factorization:



4 steps:

- 1 factorize the panel,
- 2 swapping row according to pivoting
- 3 compute block row of U
- 4 updating trailing submatrix

Modeling each step

Each step have different requirement in term of:

- computation,
- access pattern to the data,
- communication volume.

Therefore, for each substep xxx :

- γ_{xxx} : the time to perform one operation in this substep,
- β_{xxx} : the latency of the network for this substep
- α_{xxx} : the time to transfer one matrix element for this subroutine.

\Rightarrow 12 parameters.

Step 1: factorize the panel

For each column of the panel:

- 1 determine the pivot row.
- 2 swap the pivot row.
- 3 broadcast the pivot column.
- 4 update the local submatrix.

The cost for factorizing panel k is:

$$\begin{aligned} \beta_{\text{tf2}} \left(NB \left(2 \log_2(P) + \frac{P-1}{P} \right) + \log_2(Q) \right) &+ \alpha_{\text{tf2}} \left(\left(\log_2(P) + \frac{P-1}{P} \right) \frac{NB^2}{2} + NB \log_2(Q) \right) \\ &+ \frac{\gamma_{\text{tf2}}}{P} \left((N - (k-1)NB) NB^2 - \frac{1}{3} NB^3 \right) \end{aligned}$$

Factorizing all the panels ($k = [1, \dots, \lceil N/NB \rceil]$), costs:

$$\begin{aligned} \beta_{\text{tf2}} \left(N \left(2 \log_2(P) + \frac{P-1}{P} \right) + \frac{N}{NB} \log_2(Q) \right) &+ \alpha_{\text{tf2}} \left(\left(\log_2(P) + \frac{P-1}{P} \right) \frac{NNB}{2} + N \log_2(Q) \right) \\ &+ \frac{\gamma_{\text{tf2}}}{P} \left(\frac{1}{2} N(N + NB)NB - \frac{1}{3} NNB^2 \right) \end{aligned}$$

Step 2: swapping row according to pivoting

There is $\frac{P-1}{P}NB$ rows to swap between different processors on the average.

The cost for panel k is:

$$\beta_{\text{swp}} \frac{P-1}{P} NB + \alpha_{\text{swp}} \lceil (N - NB)/Q \rceil \frac{P-1}{P} NB + \gamma_{\text{swp}} \frac{NB}{P} \lceil (N - NB)/Q \rceil$$

The swapping cost for all the matrix is then:

$$\beta_{\text{swp}} N \frac{P-1}{P} + \alpha_{\text{swp}} N \lceil (N - NB)/Q \rceil \frac{P-1}{P} + \gamma_{\text{swp}} \frac{N}{P} \lceil (N - NB)/Q \rceil$$

Step 3: compute block row of U

This substep is done by :

- 1 distributing the factored upper block of the panel
- 2 performing a triangular solve

For panel k the cost is:

$$\log_2(Q)(\beta_{\text{trsm}} + \alpha_{\text{trsm}}NB^2) + \gamma_{\text{trsm}}NB^2[(N - kNB)/Q]$$

For all the whole factorization the total costs of this substep is:

$$\beta_{\text{trsm}} \log_2(Q) \frac{N}{NB} + \alpha_{\text{trsm}} \log_2(Q) NNB + \gamma_{\text{trsm}} \frac{N^2 NB}{2Q}$$

Step 4: updating trailing submatrix

This step consists in a parallel matrix multiplication:

- broadcasting of the rows,
- broadcasting of the columns,
- local multiplication of the blocks.

For panel k the cost is:

$$\begin{aligned} & \log_2(P)(\beta_{mm} + \alpha_{mm}NB \lceil (N - (k - 1)NB)/Q \rceil) \\ + & \log_2(Q)(\beta_{mm} + \alpha_{mm}NB \lceil (N - kNB)/P \rceil) \\ + & 2\gamma_{mm}NB \lceil (N - kNB)/P \rceil \lceil (N - kNB)/Q \rceil \end{aligned}$$

Hence, for the whole factorization the approximate cost is:

$$\begin{aligned} & \beta_{mm} \frac{N}{NB} (\log_2(Q) + \log_2(P)) \\ + & \alpha_{mm} \left(\log_2(Q) \frac{N^2}{2P} \log_2(P) \left(\frac{N^2}{2Q} + \frac{NNB}{Q} \right) \right) + \gamma_{mm} \frac{2N^3}{3NP} \end{aligned}$$

The p processors of a given SMP node can be arranged:

- 1 in a row-wise manner : $P \rightarrow P/p$, for counting the number of messages.
- 2 in a column-wise manner $Q \rightarrow Q/p$, for counting the number of messages.

Parameters evaluation: time spent in each step

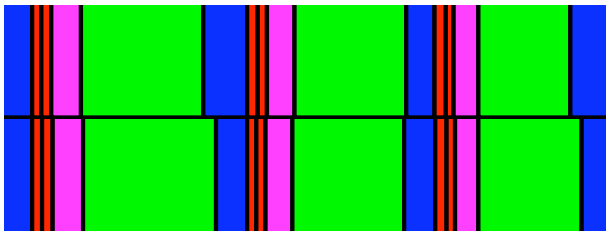
Strategy :

- ① Put timers for each routine implementing a step
- ② Run the LU factorization
- ③ Store the timing of each step for each panel and for each processor.
- ④ Sum the timings for computing the amount of time is spent in each step.
- ⑤ Fit the parameters with the timings: least square method

Summing the timings

The computation is not synchronized. Summing the timing is tricky.

Example: $P = 1$, $Q = 2$, $N = 2000$ and $NB = 64$.



Fitting the parameters

For each run:

- Build a matrix where is stored the value of the coefficient of the parameters
- Build a vector where you store the timing

Example : $N = 10000$, $NB = 32$, $P = 2$, $Q = 8$, we have: $T_{mm} = 24.8s$.

In A_{mm} we store :

- $\frac{N}{NB}(\log_2(Q) + \log_2(P)) = 1250$ in the first column,
- $\left(\log_2(Q)\frac{N^2}{2P} + \log_2(P)\left(\frac{N^2}{2Q} + \frac{NNB}{Q}\right)\right) = 6.25 \cdot 10^6$, in the second column,
- $\frac{2N^3}{3NP} = 4.1667 \cdot 10^{10}$ in the third column and,
- 24.8 in the same row of vector x .

With enough runs, a least square method on A_{mm} and x gives a value of β_{mm} , α_{mm} and γ_{mm}

In theory 3 runs are sufficient (one for each parameters).

In practice with less than 6 runs, we obtain good results

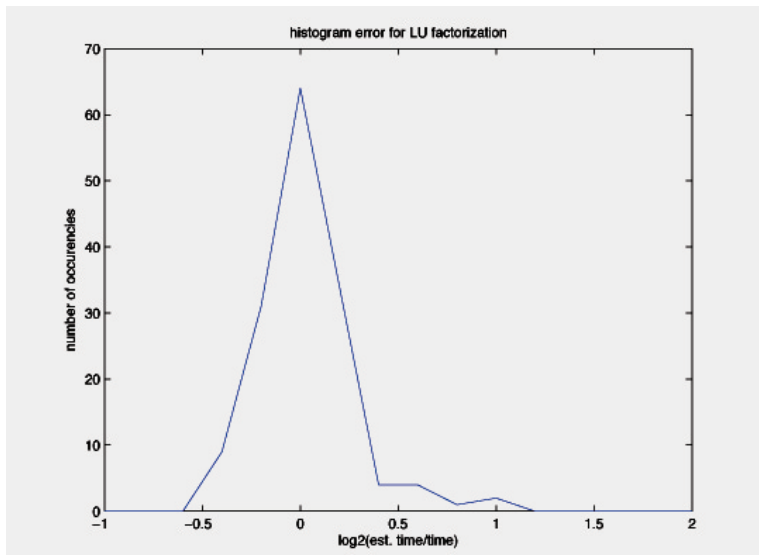
⇒ Less than one hour to compute the parameter of the model of a given machine

⇒ Easy to build automatically the model at installation time

- 48 dual-processors, AMD Opteron 64 bits, Giga ethernet interconnect
- ScaLapack.
- Precision of the model
- Computing the grid-shape

Histogram of error of our model

$NB = 32$, $N \in \{10000, 20000, 30000\}$ and $P \in [1, 64]$.



Predicting the grid shape: $N = 30000$ and $P \leq 64$

| NP | bestP | bestQ | time | estP | estQ | time | diff | error |
|----|-------|-------|---------|------|------|---------|--------|--------|
| 12 | 2 | 6 | 609.956 | 1 | 12 | 615.227 | -5.27 | -0.86% |
| 14 | 2 | 7 | 536.245 | 1 | 14 | 545.136 | -8.89 | -1.63% |
| 16 | 2 | 8 | 472.505 | 1 | 16 | 487.092 | -14.59 | -2.99% |
| 21 | 2 | 10 | 388.700 | 1 | 21 | 397.544 | -8.84 | -2.22% |
| 23 | 2 | 11 | 362.586 | 1 | 23 | 370.557 | -7.97 | -2.15% |
| 39 | 3 | 13 | 235.064 | 2 | 19 | 236.332 | -1.27 | -0.54% |
| 48 | 3 | 16 | 195.734 | 2 | 24 | 196.408 | -0.67 | -0.34% |
| 49 | 3 | 16 | 195.734 | 2 | 24 | 196.408 | -0.67 | -0.34% |
| 50 | 3 | 16 | 195.734 | 2 | 25 | 196.212 | -0.48 | -0.24% |
| 51 | 3 | 17 | 188.982 | 2 | 25 | 196.212 | -7.23 | -3.68% |
| 54 | 3 | 18 | 182.037 | 2 | 27 | 184.544 | -2.51 | -1.36% |
| 55 | 3 | 18 | 182.037 | 2 | 27 | 184.544 | -2.51 | -1.36% |
| 57 | 3 | 19 | 172.508 | 2 | 28 | 177.105 | -4.60 | -2.60% |
| 58 | 3 | 19 | 172.508 | 2 | 29 | 177.021 | -4.51 | -2.55% |
| 59 | 3 | 19 | 172.508 | 2 | 29 | 177.021 | -4.51 | -2.55% |
| 60 | 3 | 20 | 167.131 | 2 | 30 | 169.264 | -2.13 | -1.26% |
| 61 | 3 | 20 | 167.131 | 2 | 30 | 169.264 | -2.13 | -1.26% |
| 62 | 3 | 20 | 167.131 | 2 | 31 | 170.350 | -3.22 | -1.89% |
| 63 | 3 | 21 | 159.817 | 2 | 31 | 170.350 | -10.53 | -6.18% |
| 64 | 3 | 21 | 159.817 | 2 | 32 | 162.626 | -2.81 | -1.73% |

Find optimal grid-size in 68% of the cases.

Worst case: 6.18%..

- Provided a model for LU factorization
- Presented how to compute the parameters
- Results show that the model is precise and is guess very good grid-shape