

Design of Scalable Dense Linear Algebra Libraries for Multicore Processors and Multi-GPU Platforms

Enrique S. Quintana-Ortí
quintana@icc.uji.es

High Performance Computing & Architectures Group
Universidad Jaime I de Castellón (Spain)

Braunschweig – July, 2008



Joint work with:

Sergio Barrachina

Maribel Castillo

Francisco D. Igual

Rafael Mayo

Gregorio Quintana-Ortí

Rafael Rubio

Ernie Chan

Robert van de Geijn

Field G. Van Zee

Universidad Jaime I (Spain)

The University of Texas at Austin

Supported by:

National Science Foundation (NSF)

Spanish Office of Science

National Instruments

NVIDIA

...



Who has a multicore processor on the desktop/laptop?

Who has a recent graphics card on the desktop/laptop?



Who has a multicore processor on the desktop/laptop?

Are you using more than 1 core?

Who has a recent graphics card on the desktop/laptop?

Are you using it for something else than games? ;-)



Part I: Multicore processors

Part II: GPUs



New dense linear algebra libraries for multicore processors

- Scalability for manycore
- Data locality
- Heterogeneity?



LAPACK (*Linear Algebra Package*)

- Fortran-77 codes
 - One routine (algorithm) per operation in the library
 - Storage in column major order
-
- Parallelism extracted from calls to multithreaded BLAS
 - Extracting parallelism only from BLAS limits the amount of parallelism and, therefore, the scalability of the solution!
 - Column major order does hurt data locality



FLAME (*Formal Linear Algebra Methods Environment*)

- Libraries of algorithms, not codes
 - Notation reflects the algorithm
 - APIs to transform algorithms into codes
 - Systematic derivation procedure (automated using MATHEMATICA)
 - Storage and algorithm are independent
-
- Parallelism dictated by data dependencies, extracted at execution time
 - Storage-by-blocks



Part I: Multicore processors

- 1 Motivation
- 2 Cholesky factorization (Overview of FLAME)
- 3 Parallelization
- 4 Other matrix factorizations: LU & QR
- 5 Experimental results
- 6 Concluding remarks

Part II: GPUs



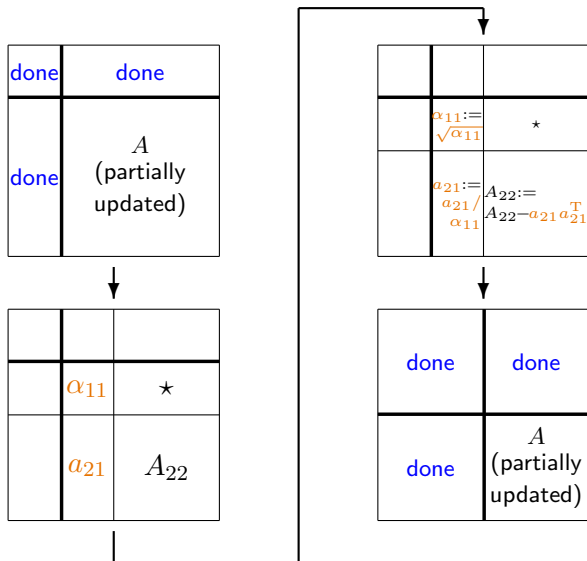
Definition

Given $A \rightarrow n \times n$ symmetric positive definite, compute

$$A = L \cdot L^T,$$

with $L \rightarrow n \times n$ lower triangular

The Cholesky Factorization: Whiteboard Presentation



done	done
done	A (partially updated)



	α_{11}	a_{12}^T
	a_{21}	A_{22}

Repartition

$$\left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right)$$

$$\rightarrow \left(\begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$$

where α_{11} is a scalar

Algorithm: $[A] := \text{CHOL_UNB}(A)$

Partition $A \rightarrow \left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right)$

where A_{TL} is 0×0

while $n(A_{BR}) \neq 0$ **do**

Repartition

$\left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$

where α_{11} is a scalar

$$\alpha_{11} := \sqrt{\alpha_{11}}$$

$$a_{21} := a_{21}/\alpha_{11}$$

$$A_{22} := A_{22} - a_{21}a_{21}^T$$

Continue with

$\left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left(\begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$

endwhile

From algorithm to code...

FLAME notation

Repartition

$$\left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$$

where α_{11} is a scalar

FLAME/C code

```
FLA_Repart_2x2_to_3x3(
    ATL, /**/ ATR,          &A00, /**/ &a01,          &A02,
    /* ***** */ /* ***** */
    &a10t, /**/ &alpha11, &a12t,
    ABL, /**/ ABR,          &A20, /**/ &a21,          &A22,
    1, 1, FLA_BR );
```



```
int FLA_Cholesky_unb( FLA_Obj A )
{
  /* ... FLA_Part_2x2( ); ... */
  while ( FLA_Obj_width( ATL ) < FLA_Obj_width( A ) ){

    FLA_Repart_2x2_to_3x3(
      ATL, /**/ ATR,          &A00, /**/ &a01,      &A02,
      /* ***** */      /* ***** */
      &a10t, /**/ &alpha11, &a12t,
      ABL, /**/ ABR,          &A20, /**/ &a21,      &A22,
      1, 1, FLA_BR );

    /*-----*/
    FLA_Sqrt( alpha11 );          /* a21 := sqrt( alpha11 ) */
    FLA_Inv_Scal( alpha11, a21 ); /* a21 := a21 / alpha11 */
    FLA_Syr      ( FLA_MINUS_ONE,
                  a21, A22 );    /* A22 := A22 - a21 * a21t */
    /*-----*/
    /* FLA_Cont_with_3x3_to_2x2( ); ... */
  }
}
```



```

int FLA_Cholesky_blk( FLA_Obj A, int nb_alg )
{
  /* ... FLA_Part_2x2( ); ... */
  while ( FLA_Obj_width( ATL ) < FLA_Obj_width( A ) ){
    b = min( FLA_Obj_length( ABR ), nb_alg );
    FLA_Repart_2x2_to_3x3(
      ATL, /**/ ATR,          &A00, /**/ &A01, &A02,
      /* ***** */      /* ***** */
      ABL, /**/ ABR,          &A10, /**/ &A11, &A12,
      b, b, FLA_BR );      &A20, /**/ &A21, &A22,

    /*-----*/
    FLA_Cholesky_unb( A11 );          /* A21 := Cholesky( A11 ) */
    FLA_Trsm_rltn( FLA_ONE, A11,
                  A21 );          /* A21 := A21 * inv( A11 )' */
    FLA_Syrk_ln ( FLA_MINUS_ONE, A21,
                 A22 ); /* A22 := A22 - A21 * A21' */
    /*-----*/
    /* FLA_Cont_with_3x3_to_2x2( ); ... */
  }
}

```


Visit [http://www.cs.utexas.edu/users/flame/Spark/...](http://www.cs.utexas.edu/users/flame/Spark/)

The screenshot shows the Spark web interface. On the left is a form for generating code, and on the right is a preview of the generated C code.

Form Fields:

- Name of the function to be generated: LU
- Type of function: blocked
- Variant Name: (empty)
- Number of operands: 1
- Operand Tag: A
- Type: matrix
- Direction: TL→BR
- Input/Output: input/output
- Pick an output language: FLAME@lab

Generated C Code:

```

s Copyright 2003, 2004, 2005, 2006 The University of Texas at Austin
s
s For licensing information see
s   http://www.cs.utexas.edu/users/flame/license.html
s
s Programmed by: Name of author
s               Email of author
s
function [ A_out ] = LU_blk_var1( A, nb_blk )
[ ATL, ATR, ...
  ABL, ABR ] = FLA_Pack_2x2( A, ...
  0, 0, 'FLA_TL' );
while ( size( ATL, 1 ) < size( A, 1 ) )
  b = min( size( ABR, 1 ), nb_blk );
  [ ADD, AD1, AD2, ...
    AID, A11, A12, ...
    A2D, A21, A22 ] = FLA_Repack_2x2_to_3x3( ATL, ATR, ...
    ABL, ABR, ...
    b, b, 'FLA_BB' );
  -----
  v           update line 1           v
  v           1                       v
  v           update line n           v
  -----
  [ ATL, ATR, ...
    ABL, ABR ] = FLA_Coat_with_3x3_to_2x2( ADD, AD1, AD2, ...
    AID, A11, A12, ...
    A2D, A21, A22, ...
  
```

- C code: FLAME/C
- M-script code for MATLAB: FLAME@lab
- Other APIs:
 - L^AT_EX
 - Fortran-77
 - LabView
 - Message-passing parallel: PLAPACK
 - FLAG: GPUs



- 1 Motivation
- 2 Cholesky factorization (Overview of FLAME)
- 3 Parallelization
- 4 Other matrix factorizations: LU & QR
- 5 Experimental results
- 6 Concluding remarks

Part II: GPUs



LAPACK parallelization: kernels in multithread BLAS

$$\left(\begin{array}{c|c} \hline & \hline A_{11} & * \\ \hline A_{21} & A_{22} \\ \hline \end{array} \right), \quad A_{11} \text{ is } b \times b$$

- Advantage: Use legacy code
- Drawbacks:
 - Each call to BLAS is a synchronization point for threads
 - As the number of threads increases, serial operations with cost $O(nb^2)$ are no longer negligible compared with $O(n^2b)$



FLAME parallelization: SuperMatrix

- Traditional (and pipelined) parallelizations are limited by the control dependencies dictated by the code
- The parallelism should be limited only by the data dependencies between operations!
- In dense linear algebra, imitate a superscalar processor:
dynamic detection of data dependencies



```

int FLA_Cholesky_blk( FLA_Obj A, int nb_alg )
{
  /* ... FLA_Part_2x2( ); ... */
  while ( FLA_Obj_width( ATL ) < FLA_Obj_width( A ) ){
    b = min( FLA_Obj_length( ABR ), nb_alg );
    /* ... FLA_Repart_2x2_to_3x3( ); ... */

    /*-----*/
    FLA_Cholesky_unb( A11 );          /* A21 := Cholesky( A11 ) */
    FLA_Trsm_rltn( FLA_ONE, A11,
                  A21 );          /* A21 := A21 * inv( A11 )' */
    FLA_Syrk_ln ( FLA_MINUS_ONE, A21,
                 A22 ); /* A22 := A22 - A21 * A21' */
    /*-----*/
    /* FLA_Cont_with_3x3_to_2x2( ); ... */
  }
}

```

The *FLAME runtime system* “pre-executes” the code:

- Whenever a routine is encountered, a pending task is annotated in a global task queue



$$\left(\begin{array}{c|c|c} A_{00} & * & * \\ \hline A_{10} & A_{11} & * \\ \hline A_{20} & A_{21} & A_{22} \end{array} \right) \quad \text{Runtime} \quad \rightarrow$$

- ① FLA_Cholesky_unb(A_{00})
- ② $A_{10} := A_{10} \text{TRIL}(A_{00})^{-T}$
- ③ $A_{20} := A_{20} \text{TRIL}(A_{00})^{-T}$
- ④ $A_{11} := A_{11} - A_{10}A_{10}^T$
- ⑤ ...

SuperMatrix

- Once all tasks are annotated, the real execution begins!
- Tasks with all input operands available are runnable; other tasks must wait in the global queue
- Upon termination of a task, the corresponding thread updates the list of pending tasks



- Algorithm and storage are independent
- Matrices stored by blocks are viewed as matrices of matrices
- No significant modification to the FLAME codes



- 1 Motivation
- 2 Cholesky factorization (Overview of FLAME)
- 3 Parallelization
- 4 Other matrix factorizations: LU & QR
- 5 Experimental results
- 6 Concluding remarks

Part II: GPUs

- Pivoting for stability limits the amount of parallelism

$$\left(\begin{array}{c|c|c} & & \\ \hline & A_{11} & A_{12} \\ \hline & A_{21} & A_{22} \\ \hline \end{array} \right), \quad A_{11} \text{ is } b \times b$$

All operations on A_{22} must wait till $\begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix}$ is factorized

- Algorithms-by-blocks for the Cholesky factorization do not present this problem
- Is it possible to design an algorithm-by-blocks for the LU factorization while maintaining pivoting?



$$\left(\begin{array}{c|cc|c} & & & \\ \hline & A_{11} & A_{12} & A_{13} \\ \hline & A_{21} & A_{22} & A_{23} \\ \hline & A_{31} & A_{32} & A_{33} \\ \hline \end{array} \right), \quad A_{ij} \text{ is } t \times t$$

- 1 Factorize $P_{11}A_{11} = L_{11}U_{11}$
- 2 Apply permutation P_{11} and factor L_{11} :

$$L_{11}^{-1}P_{11}A_{12} \mid L_{11}^{-1}P_{11}A_{13}$$

- 3 Factorize $P_{21} \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix} = L_{21}U_{21}$,

- 4 Apply permutation P_{21} and factor L_{21} :

$$L_{21}^{-1}P_{21} \begin{pmatrix} A_{12} \\ A_{22} \end{pmatrix} \mid L_{21}^{-1}P_{21} \begin{pmatrix} A_{13} \\ A_{23} \end{pmatrix}$$

- 5 Repeat steps 2–4 with A_{31}

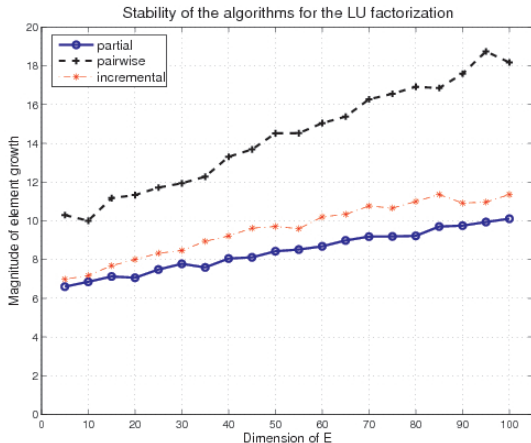


$$\left(\begin{array}{c|c|c|c} & & & \\ \hline & A_{11} & A_{12} & A_{13} \\ \hline & A_{21} & A_{22} & A_{23} \\ \hline & A_{31} & A_{32} & A_{33} \\ \hline \end{array} \right), \quad A_{ij} \text{ is } t \times t$$

Different from LU factorization with column pivoting

- To preserve structure, permutations only applied to blocks on the right!
- To obtain high performance a blocked algorithm with block size $b \ll t$, is used in the factorization and application of factors
- To maintain the computational cost, the upper triangular structure of A_{11} is exploited during the factorization

Stability? Element growth with random matrices:





- Same problem as with LU: proceeding by blocks of columns limits the amount of parallelism

$$\left(\begin{array}{c|c} \hline & \hline A_{11} & A_{12} \\ \hline & \hline A_{21} & A_{22} \\ \hline \end{array} \right), \quad A_{11} \text{ is } b \times b$$

All operations on A_{22} must wait till $\begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix}$ is factorized

- Is it possible to design an algorithm-by-blocks for the QR factorization while maintaining pivoting?



$$\left(\begin{array}{c|c|c} & & \\ \hline & A_{11} & A_{12} & A_{13} \\ \hline & A_{21} & A_{22} & A_{23} \\ \hline & A_{31} & A_{32} & A_{33} \\ \hline \end{array} \right), \quad A_{ij} \text{ is } t \times t$$

1 Factorize $Q_{11}A_{11} = R_{11}$

2 Apply factor Q_{11} :

$$Q_{11}^T A_{12} \mid Q_{11}^T A_{13}$$

3 Factorize $Q_{21} \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix} = R_{21}$,

4 Apply factor Q_{21} :

$$Q_{21}^T \begin{pmatrix} A_{12} \\ A_{22} \end{pmatrix} \mid Q_{21}^T \begin{pmatrix} A_{13} \\ A_{23} \end{pmatrix}$$

5 Repeat steps 2–4 with A_{31}



- 1 Motivation
- 2 Cholesky factorization (Overview of FLAME)
- 3 Parallelization
- 4 Other matrix factorizations: LU & QR
- 5 **Experimental results**
- 6 Concluding remarks

Part II: GPUs



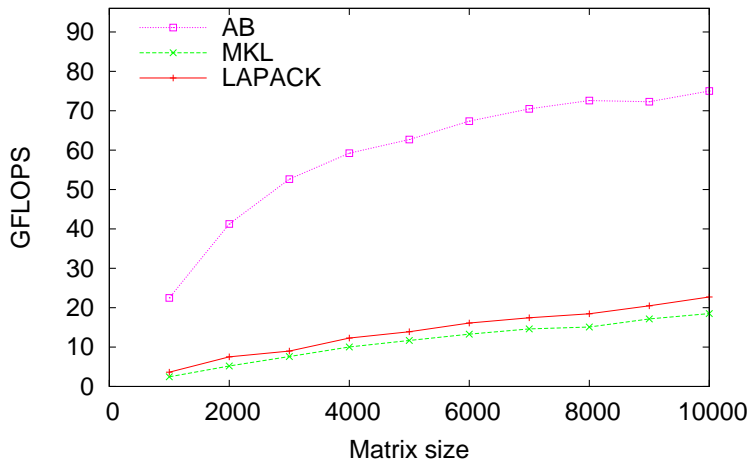
General

Platform	Specs.
SET	CC-NUMA with 16 Intel Itanium-2 processors
NEUMANN	SMP with 8 dual-core Intel Pentium4 processors

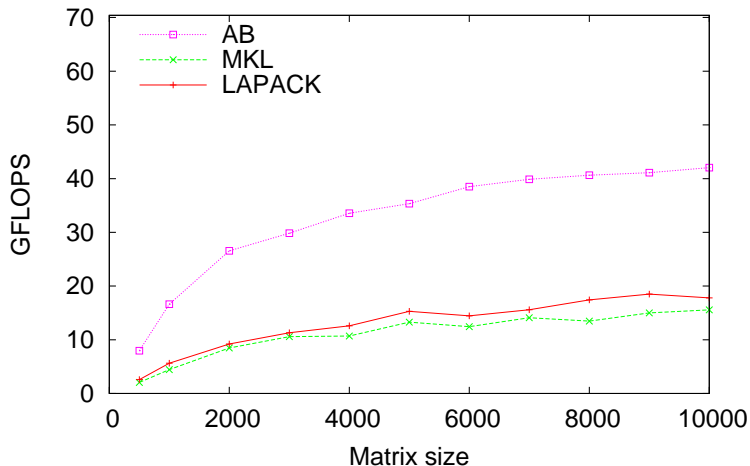
Implementations

- LAPACK 3.0 routine + multithreaded MKL
- Multithreaded routine in MKL
- AB + serial MKL
- AB + serial MKL + storage-by-blocks

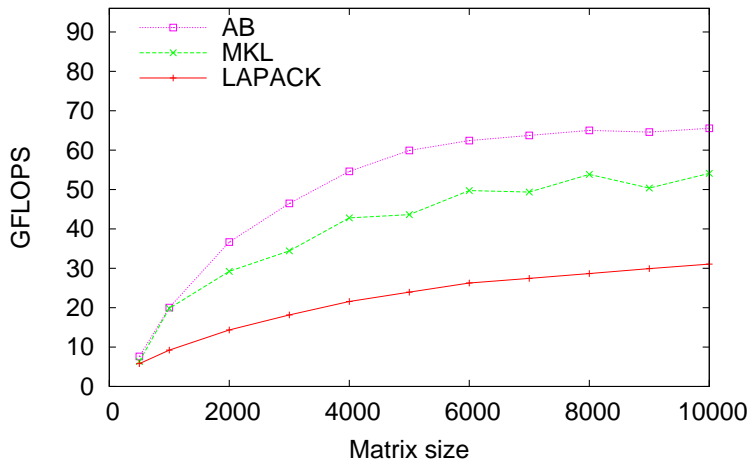
Cholesky factorization on 16 Intel Itanium 2@1.5GHz



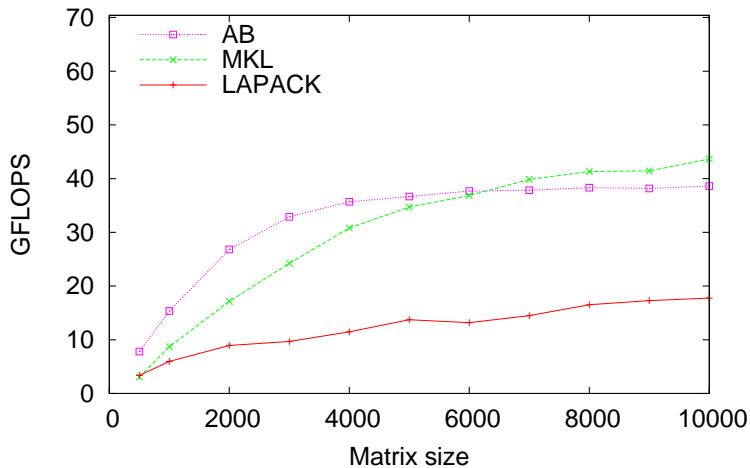
Cholesky factorization on 8 dual AMD Opteron@2.2GHz



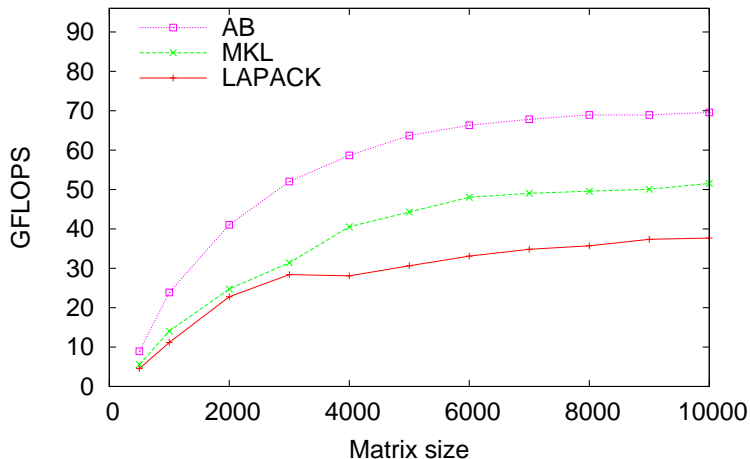
LU factorization on 16 Intel Itanium 2@1.5GHz



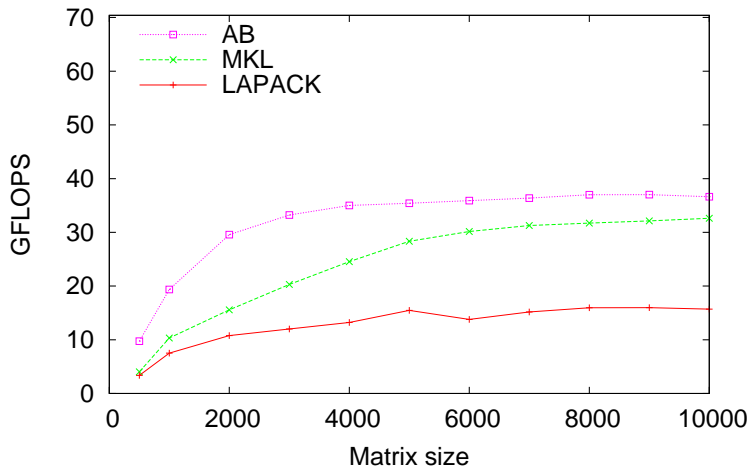
LU factorization on 8 dual AMD Opteron@2.2GHz



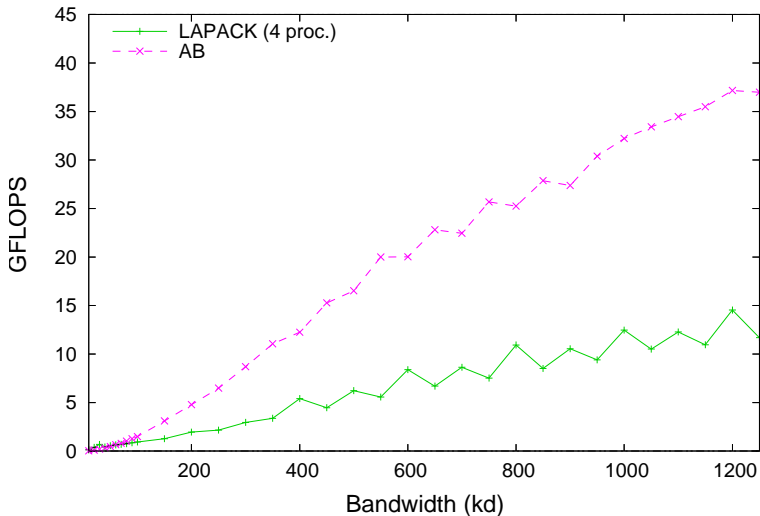
QR factorization on 16 Intel Itanium 2@1.5GHz



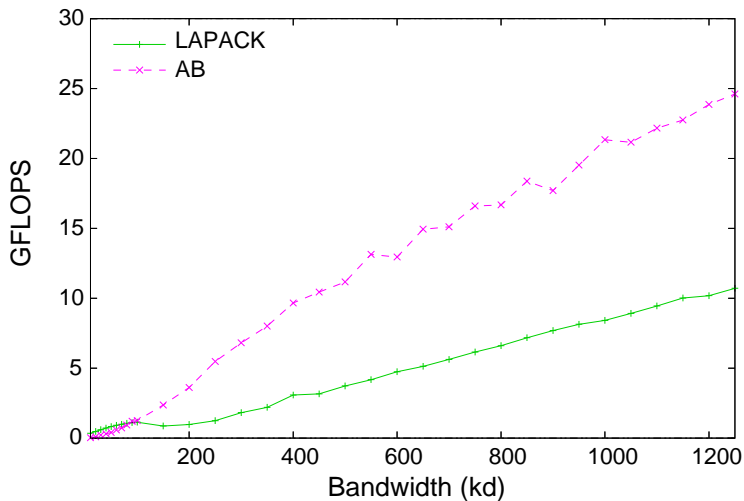
QR factorization on 8 dual AMD Opteron@2.2GHz



Band Cholesky factorization on 16 Intel Itanium 2@1



Band Cholesky factorization on 8 dual AMD Opteron@2.2





- 1 Motivation
- 2 Cholesky factorization (Overview of FLAME)
- 3 Parallelization
- 4 Experimental results
- 5 Concluding remarks

Part II: GPUs

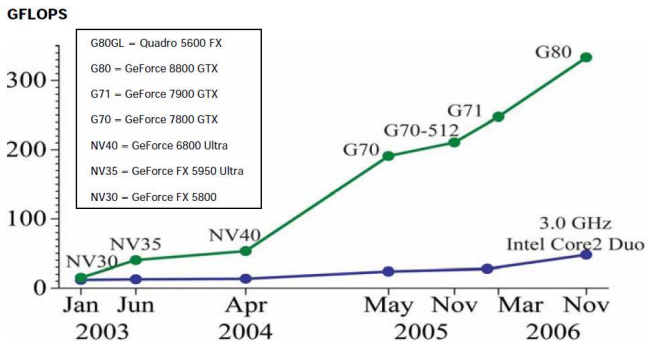


- More parallelism is needed to deal with the large number of cores of future architectures and data locality issues: traditional dense linear algebra libraries will have to be rewritten
- Some operations require new algorithms to better expose parallelism: LU with incremental pivoting, tiled QR,...
- The FLAME infrastructure (FLAME/C API, FLASH, and SuperMatrix) reduces the time to take an algorithm from whiteboard to high-performance parallel implementation



Part I: Multicore processors

Part II: GPUs



The power and versatility of modern GPU have transformed them into the first widely extended HPC platform



Part I: Multicore processors

Part II: GPUs

- 1 Motivation
- 2 Introduction
- 3 LAPACK on 1 GPU
- 4 LAPACK on multiple GPUs
- 5 FLAG@lab
- 6 Concluding remarks

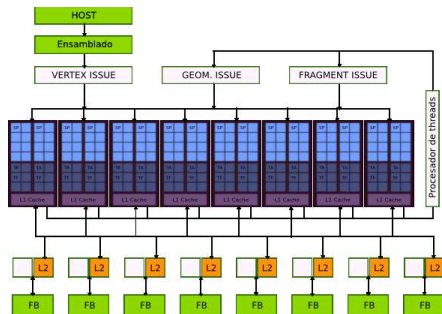


Part I: Multicore processors

Part II: GPUs

- 1 Motivation
- 2 Introduction
- 3 LAPACK on 1 GPU
- 4 LAPACK on multiple GPUs
- 5 FLAG@lab
- 6 Concluding remarks

- A CUDA-enabled device is seen as a coprocessor to the CPU, capable of executing a very high number of threads in parallel
- Example: nVIDIA G80 as a set of SIMD Multiprocessors with On-Chip Shared Memory



- Up to 128 *Streaming Processors (SP)*, grouped in clusters
- SP are SIMD processors
- Small and fast Shared Memory shared per SP cluster
- Local 32-bit registers per processor



- The CUDA API provides a simple framework for writing C programs for execution on the GPU
- Consists of:
 - A minimal set of extensions to the C language
 - A runtime library of routines for controlling the transfers between video and main memory, run-time configuration, execution of device-specific functions, handling multiple GPUs,...

CUDA libraries

On top of CUDA, nVIDIA provides two optimized libraries:
CUFFT and CUBLAS



```
int main( void ){
    ...
    float* h_vector , * d_vector;

    h_vector = ( float* ) malloc ( M * sizeof ( float ) );

    ...
    cublasAlloc ( M , sizeof ( float ) ,
                 ( void** ) & d_vector );

    cublasSetVector ( M , sizeof ( float ) , h_vector ,
                    d_vector , 1 );
    cublasSscal ( M , ALPHA , d_vector , 1 );
    cublasGetVector ( M , sizeof ( float ) , d_vector ,
                    h_vector , 1 );

    cublasFree ( d_vector );
    ...
}
```

A typical CUDA (and CUBLAS) program has 3 phases:

- 1 Allocation and transfer of data to GPU
- 2 Execution of the BLAS kernel
- 3 Transfer of results back to main memory



Part I: Multicore processors

Part II: GPUs

- 1 Motivation
- 2 Introduction
- 3 LAPACK on 1 GPU
- 4 LAPACK on multiple GPUs
- 5 FLAG@lab
- 6 Concluding remarks



Algorithm: $A := \text{CHOL_BLK}(A)$

Partition ...

where ...

while $m(A_{TL}) < m(A)$ **do**

Determine block size b

Repartition

$$\left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c|c|c} A_{00} & A_{01} & A_{02} \\ \hline A_{10} & A_{11} & A_{12} \\ \hline A_{20} & A_{21} & A_{22} \end{array} \right)$$

where A_{11} is $b \times b$

Variant 1:

$$\begin{aligned} A_{11} &:= \text{CHOL_UNB}(A_{11}) \\ A_{21} &:= A_{21} \text{TRIL}(A_{11})^{-T} \\ A_{22} &:= A_{22} - A_{21} A_{21}^T \end{aligned}$$

Variant 2:

$$\begin{aligned} A_{10} &:= A_{10} \text{TRIL}(A_{00})^{-T} \\ A_{11} &:= A_{11} - A_{10} A_{10}^T \\ A_{11} &:= \text{CHOL_UNB}(A_{11}) \end{aligned}$$

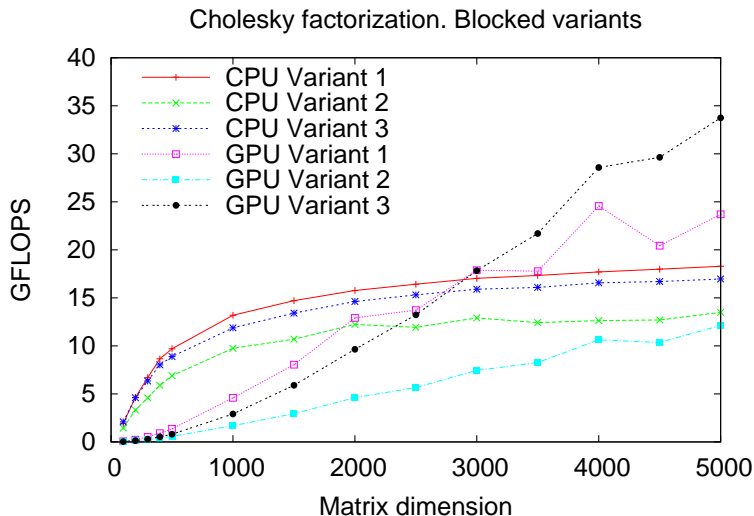
Variant 3:

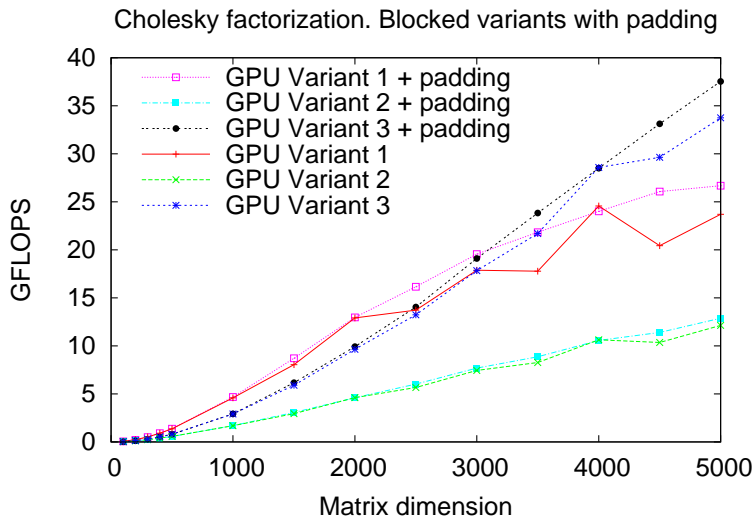
$$\begin{aligned} A_{11} &:= A_{11} - A_{10} A_{10}^T \\ A_{11} &:= \text{CHOL_UNB}(A_{11}) \\ A_{21} &:= A_{21} - A_{20} A_{10}^T \\ A_{21} &:= A_{21} \text{TRIL}(A_{11})^{-T} \end{aligned}$$

Continue with

...

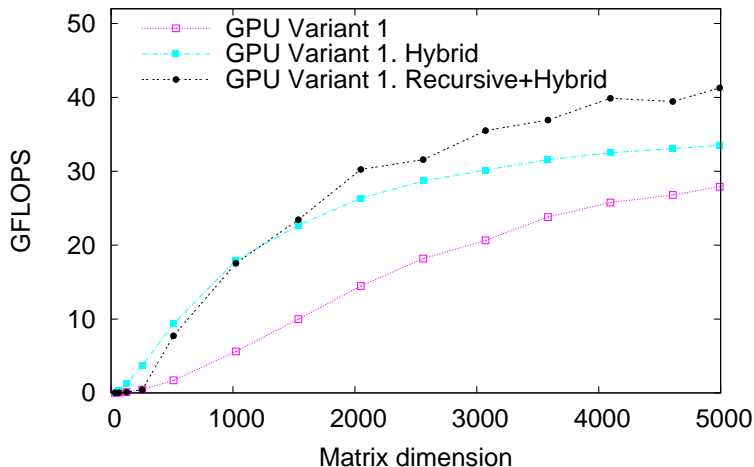
endwhile







Cholesky factorization. Recursive and hybrid variants





Compute the Cholesky factorization $A = LL^T$ and solve $(LL^T)x = b$ in the GPU

→ 32 bits of accuracy!

$i \leftarrow 1$

repeat

$$r^{(i)} \leftarrow b - A \cdot x^{(i)}$$

$$r_{(32)}^{(i)} \leftarrow r^{(i)}$$

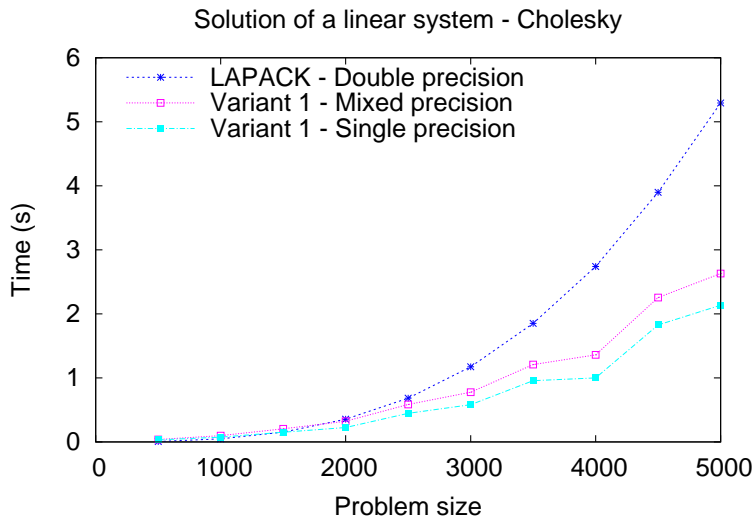
$$z_{(32)}^{(i)} \leftarrow L_{(32)}^{-T} (L_{(32)}^{-1} r_{(32)}^{(i)})$$

$$z^{(i)} \leftarrow z_{(32)}^{(i)}$$

$$x^{(i+1)} \leftarrow x^{(i)} + z^{(i)}$$

$i \leftarrow i + 1$

until $\|r^{(i)}\| < \sqrt{\varepsilon} \|x^{(i)}\|$





Part I: Multicore processors

Part II: GPUs

- 1 Motivation
- 2 Introduction
- 3 LAPACK on 1 GPU
- 4 LAPACK on multiple GPUs
- 5 FLAG@lab
- 6 Concluding remarks



Already here:

- Multiple ClearSpeed boards
- Multiple NVIDIA cards
- nVIDIA Tesla series

How are we going to program these?



Already here:

- Multiple ClearSpeed boards
- Multiple NVIDIA cards
- nVIDIA Tesla series

How are we going to program these?

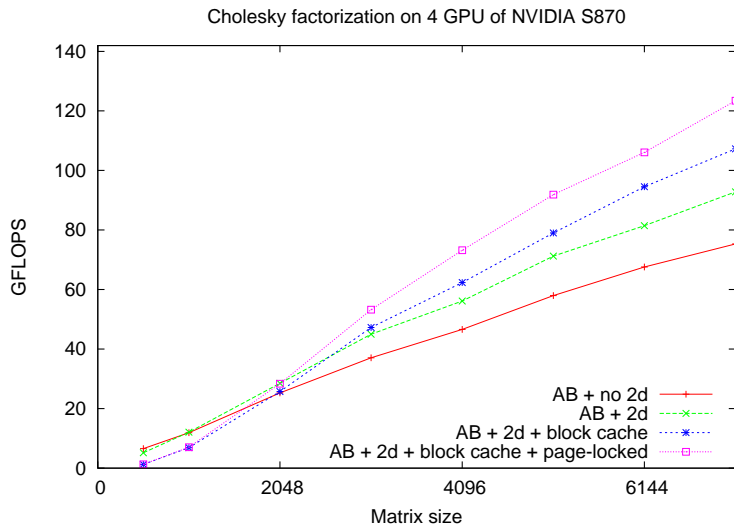


- Employ the equivalence: 1 core \equiv 1 GPU
- Difference: Transference from RAM to video memory
- Run-time system (scheduling), storage, and code are independent
- No significant modification to the FLAME codes: Interfacing to CUBLAS

A software effort of two hours!



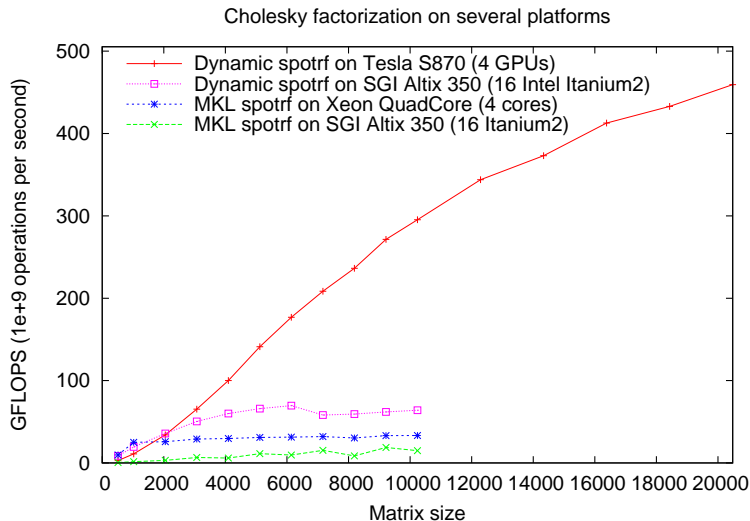
	CPU	GPU
Processor	16 x Intel Itanium2	NVIDIA Tesla s870 (4 G80)
Clock frequency	1.5 GHz	575 MHz





A more elaborate port required for high-performance:

- 2-D work distribution
- Memory/cache coherence techniques to reduce transferences between RAM and video memory: write-back and write-invalidate





Part I: Multicore processors

Part II: GPUs

- 1 Motivation
- 2 Introduction
- 3 LAPACK on 1 GPU
- 4 LAPACK on multiple GPUs
- 5 FLAG@lab
- 6 Concluding remarks



Assorted flavours:

- FLAG: A M-script API for GPU computing from MATLAB/OCTAVE
- FLAGOOC: A M-script API for Out-of-Core GPU computing from MATLAB/OCTAVE



```
function [ A ] = FLAG_Chol( A )

% [ ... ] = FLAG_Part_2x2( ... );

while ( FLAG_Obj_length( ATL, 1 ) < FLAG_Obj_length( A, 1 ) )
    b = min( FLAG_Obj_length( ABR, 1 ), nb_alg );

    % [ ... ] = FLAG_Repart_2x2_to_3x3( ... );
    %-----%
    FLAG_Chol_umb( A11 );          % A11 = Chol( A11 )
    FLAG_Trsm_rltn( 1.0, A11,
                   A21 ); % A21 = A21 * tril( A11 )^-T
    FLAG_Syrk_ln( -1.0, A21,
                 1.0, A22 ); % A22 = A22 - A21 * A21';
    %-----%
    % [ ... ] = FLAG_Cont_with_3x3_to_2x2( ... );
end
```

Just replace FLA_ in FLAME@lab by FLAG_!



- Simple precision may not be enough. Double precision is coming, but at the expense of speed?
- Overlap transferences and computation is also needed (close?)
- Programming dense linear algebra using CUBLAS on NVIDIA hardware is easy
- Programming at CUDA level?

I'll need to ask my student Francisco...



- Simple precision may not be enough. Double precision is coming, but at the expense of speed?
- Overlap transferences and computation is also needed (close?)
- Programming dense linear algebra using CUBLAS on NVIDIA hardware is easy
- Programming at CUDA level?
I'll need to ask my student Francisco...



- E. Chan, E.S. Quintana-Ortí, G. Quintana-Ortí, R. van de Geijn. SuperMatrix out-of-order scheduling of matrix operations for SMP and multicore architectures. *19th ACM Symp. on Parallelism in Algorithms and Architectures – SPAA'2007*.
- E. Chan, F. Van Zee, R. van de Geijn, E.S. Quintana-Ortí, G. Quintana-Ortí. Satisfying your dependencies with SuperMatrix. *IEEE Cluster 2007*.
- E. Chan, F.G. Van Zee, P. Bientinesi, E.S. Quintana-Ortí, G. Quintana-Ortí, R. van de Geijn. SuperMatrix: A multithreaded runtime scheduling system for algorithms-by-blocks. *Principles and Practices of Parallel Programming – PPOPP'2008*.
- E.S. Quintana-Ortí, R. van de Geijn. Updating an LU factorization with pivoting. *ACM Trans. on Mathematical Software, 2008*.



- S. Barrachina, M. Castillo, Francisco D. Igual, R. Mayo, E. S. Quintana-Ort' i. Evaluation and tuning of the level 3 CUBLAS for graphics processors. *Workshop on Parallel and Distributed Scientific and Engineering Computing, – PDSEC'2008.*
- S. Barrachina, M. Castillo, F. Igual, R. Mayo, E. S. Quintana. Solving dense linear systems on graphics processors. *Euro-Par'2008.*
- M. Castillo, F. Igual, R. Mayo, R. Rubio, E. S. Quintana, G. Quintana, R. van de Geijn. Out-of-Core Solution of Linear Systems on Graphics Processors. *Parallel/High-Performance Object-Oriented Scientific Computing – POOSC'08 .*



Cilk (MIT) and CellSs (Barcelona SuperComputing Center)

- **General-purpose** parallel programming
 - Cilk → irregular problems
 - CellSs → for the Cell B.E.
- High-level language based on OpenMP-like pramas + **compiler** + runtime system
- Moderate results for dense linear algebra

PLASMA (UTK – Jack Dongarra)

- **Traditional style** of implementing algorithms: Fortran-77
- **Complicated coding**
- Runtime system + ?

For more information...

Visit <http://www.cs.utexas.edu/users/flame>

Support...

- *National Science Foundation* awards CCF-0702714 and CCF-0540926 (ongoing till 2010)
- Spanish CICYT project TIN2005-09037-C02-02