# Solving Dense Linear Systems on Platforms with Multiple Hardware Accelerators

Francisco D. Igual Enrique S. Quintana-Ortí Gregorio Quintana-Ortí

Universidad Jaime I de Castellón (Spain)

Robert A. van de Geijn

The University of Texas at Austin

PPoPP09 - February, 2009



http://www.cs.utexas.edu/users/flame/

### Joint work with:

Universidad Jaime I (Spain)	UT-Austin
Sergio Barrachina	Ernie Chan
Maribel Castillo	Robert A. van de Geijn
Francisco D. Igual	Field G. Van Zee
Rafael Mayo	
Enrique S. Quintana-Ortí	
Gregorio Quintana-Ortí	
Rafael Rubio	

### Supported by:

National Science Foundation (NSF) Spanish Office of Science nVIDIA Anonymous corporation



## The sky is falling

- Parallelism is thrust upon the masses
- Popular libraries like LAPACK must be completely rewritten
- Give us money or the world as we know it will come to an end



### Evolution vs intelligent design

- Parallelism is thrust upon the masses
- Popular libraries like LAPACK must be completely rewritten (end of an evolutionary path)
- Great, let's finally toss them and start over
- Cheaper than trying to evolve?



## 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

## Problems

- Extracting parallelism increases synchronization and thus limits performance
- Column major order hurts data locality

# Motivation



## FLAME (Formal Linear Algebra Methods Environment)

- Notation for expressing algorithms
- Systematic derivation procedure (automated using MATHEMATICA)
- Families of algorithms for each operation
- APIs to transform algorithms into codes
- Storage and algorithm are independent

### We will show FLAME elegantly supports

- Storage-by-blocks
- Parallelism with data dependencies
- High performance even on "exotic" architectures like multiGPUs









# Motivation

- Cholesky factorization (Overview of FLAME)
- Parallelization
- Targeting platforms with multiple accelerators
- Section 2 Experimental results
- Oncluding remarks



## Definition

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

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

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

Algorithms should ideally be represented in a way that captures how we reason about them

# The Cholesky Factorization: On the Whiteboard





# **FLAME** Notation





# **FLAME** Notation





# FLAME/C Code



#### From algorithm to code...

## FLAME notation

#### Repartition

$$\begin{pmatrix} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{pmatrix} \rightarrow \begin{pmatrix} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^{\mathsf{T}} & \alpha_{11} & a_{12}^{\mathsf{T}} \\ \hline A_{20} & a_{21} & A_{22} \end{pmatrix}$$

where  $\alpha_{11}$  is a scalar

#### FLAME/C representation in code

http://www.cs.utexas.edu/users/flame/

# (Unblocked) FLAME/C Code



```
int FLA_Cholesky_unb( FLA_Obj A )
Ł
 /* ... FLA_Part_2x2( ); ... */
 while ( FLA_Obj_width( ATL ) < FLA_Obj_width( A ) ){</pre>
  FLA_Repart_2x2_to_3x3(
       ATL, /**/ ATR, &A00, /**/ &a01, &A02,
     &a10t, /**/ &alpha11, &a12t,
       ABL, /**/ ABR, &A20, /**/ &a21, &A22,
       1. 1, FLA_BR );
  /*----*/
  FLA_Syr (FLA_MINUS_ONE,
           a21, A22 ); /* A22 := A22 - a21 * a21t */
  /*-----
                            ----*/
  /* FLA_Cont_with_3x3_to_2x2( ); ... */
 }
}
```

# Blocked FLAME/C Code



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

Elle Edit Yew Go Booksnarks Iools Window Help						
🔾 🗿 🕥 🔕 🗈 http://www.cs.utevas.obd/aem/	fane/Sark/	3.5				
New Tab (Sy Spark						
	1					
Generate Code and/or Update Form Reset Form	* Copyright 2003, 2004, 2005, 2006 The University of Texas at	Austin				
	4					
an about this section in Soule	<pre>% For licensing information see</pre>	1000				
ALCOLOGICAL STREET	http://www.cs.utexas.edw/users/flowe/licease.	RCm1				
lame of the function to be generated	h Programmed by: Name of author					
U	b Email of author					
and the second	formation of a sure of an All holds served of the starts					
ype of function blocked 💌	Landeron ( * one ) - no ore refer( *, no ard )					
Sector	[ ATL, ATR,					
ariant Name 1 👻	ABL, ABR ] = FLA_Part_2x2( A,					
	0, 0, 'FLA_TL' );					
even should this section extendention to Soudy	while ( size( ATL, 1 ) < size( A, 1 ) )					
Sumber of operands 1 .	b = min( size( ADP, 1 ), sb_alg );					
and the second	[ 100 101 102					
ick properties of the operands	A10, A11, A12,					
Operand Tag Type Direction Input/Output	A20, A21, A22 ] = FLA_Repart_2x2_to_3x3 ( ATL, ATR,					
1 A M matter M TI-SDD M monthstard M	ABL, ABR,					
T A R HOLK R LEADL R HOLEODO R	D, D, 'YLA_BR'	11				
	\	·····				
earn about this section introduction to Speck						
test of fell the second of the second	supdate line 1	2				
Ack an output language: FLAME@iob 🛛 💌	and ate line a					
sam about this section introduction to 2pack	· · · · · · · · · · · · · · · · · · ·	h				
	1 477 4779					
Additional Information	ABL, ABR 1 = FLA Cost with 3x3 to 2x2 / A00, A01, A02,					
of the second	A10, A11, A12,	<i>i</i>				
Name of Author Name of author	A20, A21, A22,	á – 1				
		-				

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



# Motivation

- Cholesky factorization (Overview of FLAME)
- Parallelization
- Targeting platforms with multiple accelerators
- Section 2 Experimental results
- 6 Concluding remarks



## LAPACK parallelization: multithreaded BLAS

(	1	\			
	$A_{11}$	*			
	A <sub>21</sub>	A <sub>22</sub>	,	A <sub>11</sub> i	$s \ b  imes b$

- Pro?:
  - Evolve legacy code
- Con:
  - Continue to code in the LINPACK style (1970s)
  - Each call to BLAS (compute kernels) 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  $O(n^2b)$



#### Separation of concern

- Improve parallelism and data locality: algorithms-by-blocks
  - Matrix of matrix blocks
  - Matrix blocks as unit of data
  - Computation with matrix blocks as unit of computation
- Execute sequential code to generate DAG of tasks
- SuperMatrix
  - Runtime system for scheduling tasks to threads
- Sequential kernels to be executed by the threads
  - Always be sure to make the machine-specific part someone else's problem

# Matrix of blocks

$$A = \begin{pmatrix} A^{(0,0)} & \star & \star & \cdots & \star \\ A^{(1,0)} & A^{(1,1)} & \star & \cdots & \star \\ A^{(2,0)} & A^{(2,1)} & A^{(2,2)} & \cdots & \star \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ A^{(M-1,0)} & A^{(M-1,1)} & A^{(M-1,2)} & \cdots & A^{(M-1,N-1)} \end{pmatrix}$$

#### FLAME implementation of algorithm-by-blocks: no change

```
int FLA_Cholesky_blk( FLA_Obj A, int nb_alg )
ſ
 /* ... FLA_Part_2x2( ); ... */
 while ( FLA_Obj_width( ATL ) < FLA_Obj_width( A ) ){</pre>
   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

$(A^{(0,0)})$	*	*	•••	* )
$A^{(1,0)}$	$A^{(1,1)}$	*	• • •	*
A <sup>(2,0)</sup>	$A^{(2,1)}$	$A^{(2,2)}$		*
÷	÷	÷	·	÷
$A^{(M-1,0)}$	$A^{(M-1,1)}$	$A^{(M-1,2)}$		$A^{(M-1,N-1)}$

( A <sup>(0,0)</sup>	*	*		* )	
$A^{(1,0)}$	$A^{(1,1)}$	*		*	
A <sup>(2,0)</sup>	$A^{(2,1)}$	A <sup>(2,2)</sup>		*	
:	:	:	۰.	÷	
$A^{(M-1,0)}$	$A^{(M-1,1)}$	$A^{(M-1,2)}$		$A^{(M-1,N-1)}$ )	

# FLAME Parallelization: SuperMatrix





#### SuperMatrix

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

http://www.cs.utexas.edu/users/flame/



# Motivation

- Cholesky factorization (Overview of FLAME)
- Parallelization
- Targeting platforms with multiple accelerators
- Section 2 Experimental results
- 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



```
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:
```

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

# How to use multiple GPUs?





# How to use multiple GPUs?







- $\bullet$  Employ the equivalence: 1 core  $\equiv$  1 GPU
- Difference: Explicitly transfer data from RAM to video memory
- Run-time system (scheduling), storage, and code are independent
- No significative modification to the FLAME codes:
  - Interfacing to CUBLAS
- Run-time system can implement different scheduling policies

## Programming effort



- Employ the equivalence: 1 core  $\equiv 1 \mbox{ GPU}$
- Difference: Explicitly transfer data from RAM to video memory
- Run-time system (scheduling), storage, and code are independent
- No significative modification to the FLAME codes:
  - Interfacing to CUBLAS
- Run-time system can implement different scheduling policies

### Programming effort

### Two hours!



- Motivation
- Cholesky factorization (Overview of FLAME)
- Parallelization
- Targeting platforms with multiple accelerators
- Second Experimental results
- 6 Concluding remarks

# Porting SuperMatrix. Experimental results







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

For details, see PPoPP09 paper.















## Level-3 BLAS (matrix-matrix operations)

Follow very similarly

## LU/QR factorization

- New algorithms-by-blocks have been developed
- Require kernels for an individual GPU to be written



- Motivation
- Cholesky factorization (Overview of FLAME)
- Parallelization
- Targeting platforms with multiple accelerators
- Section 2 Experimental results
- 6 Concluding remarks



# Cilk (MIT) and SMPSs (Barcelona SuperComputing Center)

- General-purpose parallel programming
  - $\bullet \ Cilk \rightarrow irregular/recursive \ problems$
  - $\bullet~SMPSs \rightarrow$  more general, also manages dependencies
- High-level language based on OpenMP-like pragmas + compiler + runtime system
- Modest results for dense linear algebra

## **PLASMA** Project

- Next step in the LAPACK evolutionary path
- Traditional style of implementing algorithms
- Does not solve the programmability problem



- Motivation
- Cholesky factorization (Overview of FLAME)
- Parallelization
- Targeting platforms with multiple accelerators
- Section 2 Experimental results
- 6 Concluding remarks



#### A success story

- For the domain of dense linear algebra libraries, FLAME+SuperMatrix appears to solve the programmability problem
- Very good performance can be achieved with multiple accelerators in this domain

### For more information...

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