# Retargeting PLAPACK to Clusters with Hardware Accelerators

#### Manuel Fogué<sup>1</sup> Francisco Igual<sup>1</sup> Enrique S. Quintana-Ortí<sup>1</sup> Robert van de Geijn<sup>2</sup>

<sup>1</sup>Departamento de Ingeniería y Ciencia de los Computadores. University Jaume I. Castellón (Spain).

<sup>2</sup>Department of Computer Sciences. Texas Institute for Computational Engineering and Sciences. The University of Texas at Austin.



Motivation

### Motivation

#### Past, present and future of GPU computing

#### Systems with one GPU



#### Software:

- CUDA
- OpenCL
- OpenGL, Cg, ...
- ATI Stream

#### Multi-GPU systems



### Software:

- GPUSs
- Star-PU
- FLAME SuperMatrix

#### **Clusters of GPUs**



### Software:

- ??
- Our approach: CUPLAPACK

#### Future is today...

China's new Nebulae Supercomputer is No. 2 at Top500 (June 2010) Powered with Tesla C2050 GPUs Motivation

### Motivation

#### Past, present and future of GPU computing

#### Systems with one GPU



#### Software:

- CUDA
- OpenCL
- OpenGL, Cg, ...
- ATI Stream

#### Multi-GPU systems

### Software:

- GPUSs
- Star-PU
- FLAME SuperMatrix

#### **Clusters of GPUs**



### Software:

- ??
- Our approach: CUPLAPACK

#### Future is today...

China's new Nebulae Supercomputer is No. 2 at Top500 (June 2010) Powered with Tesla C2050 GPUs Motivation

### Motivation

#### Past, present and future of GPU computing

#### Systems with one GPU



#### Software:

- CUDA
- OpenCL
- OpenGL, Cg, ...
- ATI Stream

#### Multi-GPU systems



### Software:

- GPUSs
- Star-PU
- FLAME SuperMatrix

#### **Clusters of GPUs**



#### Software:

- ??
- Our approach: CUPLAPACK

#### Future is today...

China's new Nebulae Supercomputer is No. 2 at Top500 (June 2010) Powered with Tesla C2050 GPUs







- Porting PLAPACK to clusters of GPUs
- Experimental results
- 5 Conclusions and future work



### • PLAPACK:

• Parallel dense linear algebra package for message-passing architectures (i.e., clusters).

#### • CUPLAPACK:

- Retarget of PLAPACK to clusters with graphics processors.
- Why a PLAPACK-based approach?:
  - Modularity: layered design.
  - Programmability: FLAME methodology. Object-based approach.
- Goals:
  - Programmability: Transparent port to clusters of GPUs.
  - Performance: Clusters of GPUs with hundreds of nodes.
  - Scalability: Keep performance with large problems.

### **PLAPACK** overview



#### PLAPACK: Parallel Linear Algebra Package

- "Library infrastructure for parallel implementation of linear algebra algorithms on distributed memory machines".
- Natural transition from algorithms to distributed-memory codes. Focuses on programmability.
- Layered design. Focuses on extensibility and flexibility.
- Physically Based Matrix Distribution (PBMD) and templates.

### Cholesky factorization. Right-looking variant (I)

Cholesky factorization of a positive definite matrix A:

$$A = LL^T$$

Partition:

$$A = \begin{pmatrix} A_{11} & \star \\ A_{21} & A_{22} \end{pmatrix}, \ L = \begin{pmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{pmatrix},$$

where  $A_{11}$  and  $L_{11}$  are  $b \times b$  matrices.

From  $A = LL^T$ , we have that:

$$\begin{pmatrix} A_{11} & \star \\ A_{21} & A_{22} \end{pmatrix} = \begin{pmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{pmatrix} \begin{pmatrix} L_{11}^T & L_{21}^T \\ 0 & L_{22}^T \end{pmatrix}$$
$$= \begin{pmatrix} L_{11}L_{11}^T & \star \\ L_{21}L_{11}^T & L_{21}L_{21}^T + L_{22}L_{22}^T \end{pmatrix}.$$

### Cholesky factorization. Right-looking variant (I)

Cholesky factorization of a positive definite matrix A:

$$A = LL^{T}$$

Partition:

$$A = \begin{pmatrix} A_{11} & \star \\ A_{21} & A_{22} \end{pmatrix}, \ L = \begin{pmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{pmatrix},$$

where  $A_{11}$  and  $L_{11}$  are  $b \times b$  matrices.

From  $A = LL^T$ , we have that:

$$\begin{pmatrix} A_{11} & \star \\ A_{21} & A_{22} \end{pmatrix} = \begin{pmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{pmatrix} \begin{pmatrix} L_{11}^T & L_{21}^T \\ 0 & L_{22}^T \end{pmatrix}$$
$$= \begin{pmatrix} L_{11}L_{11}^T & \star \\ L_{21}L_{11}^T & L_{21}L_{21}^T + L_{22}L_{22}^T \end{pmatrix}$$

# Cholesky factorization. Right-looking variant (II)

This yields to equations:

$$L_{11}L_{11}^{T} = A_{11},$$
  

$$L_{21} = A_{21}L_{11}^{-T},$$
  

$$A_{22} - L_{21}L_{21}^{T} = L_{22}L_{22}^{T}.$$

Concluding with the algorithm:

# Cholesky factorization. Right-looking variant (II)

This yields to equations:

$$L_{11}L_{11}^{T} = A_{11},$$
  

$$L_{21} = A_{21}L_{11}^{-T},$$
  

$$A_{22}-L_{21}L_{21}^{T} = L_{22}L_{22}^{T}.$$

Concluding with the algorithm:

### Using PLAPACK as a library

#### From algorithms to code

partition 
$$A = \left(\frac{A_{TL} \mid \star}{A_{BL} \mid A_{BR}}\right)$$
  
where  $A_{TL}$  is  $0 \times 0$   
do until  $A_{BR}$  is  $0 \times 0$   
determine block size  $b$   
repartition  
 $\left(\frac{A_{TL} \mid \star}{A_{BL} \mid A_{BR}}\right) = \left(\frac{A_{00} \mid \star \mid \star}{A_{20} \mid A_{21} \mid A_{22}}\right)$   
where  $A_{11}$  is  $b \times b$   
 $A_{11} \leftarrow L_{11} = \text{Chol. fact.}(A_{11})$   
 $A_{21} \leftarrow L_{21} = A_{21}L_{11}^{-T}$   
 $A_{22} \leftarrow A_{22} - L_{21}L_{21}^{T}$   
continue with  
 $\left(\frac{A_{TL} \mid \star}{A_{BL} \mid A_{BR}}\right) = \left(\frac{A_{00} \mid \star \mid \star}{A_{20} \mid A_{21} \mid A_{22}}\right)$   
enddo

### Using PLAPACK as a library

#### From algorithms to code

$$partition A = \left(\frac{A_{TL}}{A_{BL}} \mid \frac{\star}{A_{BR}}\right)$$
where  $A_{TL}$  is  $0 \times 0$   
do until  $A_{BR}$  is  $0 \times 0$   
determine block size  $b$   
repartition  

$$\left(\frac{A_{TL}}{A_{BL}} \mid \frac{\star}{A_{BR}}\right) = \left(\frac{A_{00} \mid \star \mid \star}{A_{10} \mid A_{11} \mid \star}\right)$$
where  $A_{11}$  is  $b \times b$   
where  $A_{11}$  is  $b \times b$   

$$A_{11} \leftarrow L_{11} = Chol. fact.(A_{11})$$

$$A_{21} \leftarrow L_{21} = A_{21}L_{11}^{-T}$$

$$A_{22} \leftarrow A_{22} - L_{21}L_{21}^{T}$$
continue with  

$$\left(\frac{A_{TL} \mid \star}{A_{BL} \mid A_{BR}}\right) = \left(\frac{A_{00} \mid \star \mid \star}{A_{20} \mid A_{11} \mid \star}\right)$$
enddo

#### int PLA Chol( int b, PLA Obi A ) PLA Obi ABR = NULL. A11 = NULL. A21 = NULL: /\* ... \*/ /\* View ABR = A \*/ PLA Obj view all( A, &ABR ); while (TRUE) { /\* Partition ABR = / A11 || \* \ A21 // ABR / \* where A11 is b x b \*/ PLA Obj split 4( ABR, b, b, &A11, PLA DUMMY, &A21, &ABR ); /\* A11 := L11 = Cholesky Factor( A11 ) \*/ PLA Local chol ( PLA LOWER TRIANGULAR, A11 ); /\* Update A21 := L21 = A21 \* inv( L11' ) \*/ PLA Trsm( PLA SIDE RIGHT, PLA LOWER TRIANGULAR, PLA TRANSPOSE, PLA NONUNIT DIAG, one, A11, A21); /\* Update A22 := A22 - L21 \* L21 ' \*/ PLA Svrk( PLA LOWER TRIANGULAR, PLA NO TRANS, minus one. A21. one. ABR ): /\* ... \*/

### Using PLAPACK as a library. Cholesky code

```
int PLA Chol ( int b, PLA Obj A )
 2
3
      PLA_Obj ABR = NULL.
 4
             A11 = NULL, A21 = NULL;
 5
      /* ... */
 6
 7
      /* View ABR = A */
8
      PLA Obj view all( A, &ABR );
9
      while (TRUE) {
10
        /* Partition ABR = / A11 // * \
11
12
        *
                           / ===== /
                           \ A21 || ABR /
13
14
        * where A11 is b x b
                                       */
15
        PLA Obi split 4 (ABR, b, b,
16
                            &A11, PLA_DUMMY,
17
                            &A21. &ABR ):
18
        /* A11 := L11 = Cholesky Factor( A11 ) */ /* A11 := L11 = Cholesky Factor( A11 ) */
19
20
21
22
        /* Update A21 := L21 = A21 * inv( L11') */
23
        PLA_Trsm( PLA_SIDE_RIGHT, PLA_LOWER_TRIANGULAR,
24
                 PLA TRANSPOSE. PLA NONUNIT DIAG.
25
                 one. A11.
26
                      A21 ):
27
28
        /* Update A22 := A22 - L21 * L21 ' */
29
30
                 minus one, A21,
31
                 one, ABR);
32
      /* ... */
33
34
```

```
PLA Local chol (PLA LOWER TRIANGULAR, A11); CUPLA Local chol (PLA LOWER TRIANGULAR, A11);
/* Update A22 := A22 - L21 * L21 ' */

PLA Syrk( PLA LOWER TRIANGULAR, PLA NO TRANS, CUPLA Syrk( PLA LOWER TRIANGULAR, PLA NO TRANS,
```

### Using PLAPACK as a library. Cholesky code

j

```
int PLA Chol ( int b, PLA Obj A )
 2
3
      PLA_Obj ABR = NULL.
 4
             A11 = NULL, A21 = NULL;
 5
      /* ... */
 6
 7
      /* View ABR = A */
8
      PLA Obj view all( A, &ABR );
9
10
      while (TRUE) {
        /* Partition ABR = / A11 //
12
        *
                            / ===== /
                            A21 // ABR /
13
14
        * where A11 is b x b
                                        */
15
        PLA Obi split 4 (ABR, b, b,
16
                             &A11, PLA_DUMMY,
17
                             &A21. &ABR ):
18
        /* A11 := L11 = Cholesky Factor( A11 ) */ /* A11 := L11 = Cholesky Factor( A11 ) */
19
        PLA_Local_chol ( PLA_LOWER_TRIANGULAR, A11 );
20
21
22
        /* Update A21 := L21 = A21 * inv( L11') */
23
24
                  PLA TRANSPOSE. PLA NONUNIT DIAG.
25
                  one. A11.
26
                       A21 ):
27
28
        /* Update A22 := A22 - L21 * L21 ' */
        /* Update A22 := A22 - L21 * L21 ' */

PLA Syrk (PLA LOWER TRIANGULAR, PLA NO TRANS,
29
30
                 minus one, A21,
31
                  one, ABR ):
32
33
      /* ... */
34
```

```
int CUPLA Chol( int b, PLA Obj A )
                                                  PLA Obj ABR = NULL,
                                                        A11 = NULL, A21 = NULL;
                                                  /* ... */
                                                  /* View ABR = A */
                                                  PLA Obj view all( A, &ABR );
                                                  while (TRUE) {
                                                  /* Partition ABR = / A11 // *
                                                                     / ===== /
                                                                      \ A21 || ABR /
                                                 * where A11 is b x b
                                                                        */
                                                  PLA Obi split 4 (ABR, b, b,
                                                                        &A11. PLA DUMMY.
                                                                        &A21. &ABR ):
                                                    CUPLA Local chol ( PLA LOWER TRIANGULAR, A11 );
                                                 /* Update A21 := L21 = A21 * inv(L11') */
PLA_Trsm(PLA_SIDE_RIGHT, PLA_LOWER_TRIANGULAR, CUPLA_Trsm(PLA_SIDE_RIGHT, PLA_LOWER_TRIANGULAR,
                                                               PLA TRANSPOSE. PLA NONUNIT DIAG.
                                                               one. A11.
                                                                   A21 ):
                                                 /* Update A22 := A22 - L21 * L21 ' */
                                                CUPLA Syrk ( PLA LOWER TRIANGULAR, PLA NO TRANS,
                                                               minus one, A21,
                                                               one, ABR ):
                                                /* ... */
```

### GPU acceleration of PLAPACK

### How to accelerate PLAPACK?

- Naive implementation:
  - Objects created in RAM.
  - GPU 
     ⊂ CPU transfers bound to calculations.
  - Intercept calls to local BLAS and transfer data to/from GPU.
  - Straightforward implementation.
- Tuned implementation:
  - Objects created in GPU.
  - GPU 
     ⊂ CPU transfers bound to communications.
  - Only transfer to RAM when necessary (communication).
  - Creation of objects involves allocation on GPUs.
  - Modification of communication layer in PLAPACK (PLA\_Copy, PLA\_Reduce).

### GPU acceleration of PLAPACK

#### How to accelerate PLAPACK?

- Naive implementation:
  - Objects created in RAM.
  - GPU 
     ⊂ CPU transfers bound to calculations.
  - Intercept calls to local BLAS and transfer data to/from GPU.
  - Straightforward implementation.
- Tuned implementation:
  - Objects created in GPU.
  - GPU  $\rightleftharpoons$  CPU transfers bound to communications.
  - Only transfer to RAM when necessary (communication).
  - Creation of objects involves allocation on GPUs.
  - Modification of communication layer in PLAPACK (PLA\_Copy, PLA\_Reduce).

## **PLAPACK** infrastructure

#### PLAPACK layered design

Naive User Application						application layer
Application Program Interface	Higher Level Global LA Routines PLA_Global BLAS				library layer	
(PLA_API)	PLA_Copy/Red	Copy/Reduce LA object PLA_Local BLAS		A_Local BLAS	PLAPACK abstraction layer	
MMPI	PLA/MPI interface	PLA_malloe	PBMI Templat	) .es	PLA/ BLAS interface	machine/distribution independent layer
Message-Passing Interface		malloc	cartesian distributions		vendor BLAS	machine/distribution specific layer

#### Necessary changes to use GPU

- PLA\_Local\_BLAS: NVIDIA CUBLAS.
- IA object manipulation: management of objects in GPU memory.
- Organization layer: transfers to RAM prior to communications.

## PLAPACK infrastructure

#### PLAPACK layered design

Naive User Application						application layer
Application Program Interface	Higher Level Global LA Routines PLA_Global BLAS				library layer	
(PLA_API)	PLA_Copy/Rec	educe LA object manipulation PLA_Local BLAS		PLAPACK abstraction layer		
MMPI	PLA/MPI interface	PLA_malloe	PBMD Template	es	PLA/ BLAS interface	machine/distribution independent layer
Message-Passing Interface		malloc	cartesian distributions		vendor BLAS	machine/distribution specific layer

#### Necessary changes to use GPU

- PLA\_Local\_BLAS: NVIDIA CUBLAS.
- 2 LA object manipulation: management of objects in GPU memory.
- Ommunication layer: transfers to RAM prior to communications.

### LA object manipulation. Cholesky driver using GPU M

```
int main( void ){
 1
 2
       /* ... */
 3
       // Object creation on GPU
       CUPLA Matrix create (MPI FLOAT, size, size, templ,
 4
 5
                           PLA ALIGN FIRST. PLA ALIGN FIRST.
 6
                           &A GPU):
 7
8
       // Object initialization on GPU
 9
       CUPLA Obi set to SPD random( A GPU ):
10
       /* ... */
11
12
13
       // Cholesky factorization on GPU
       CUPLA Chol ( nb alg. A GPU ):
14
15
       /* ... */
16
17
       // Object destruction on GPU
       CUPLA Obi free ( &A GPU ):
18
19
```

- CUPLA\_Matrix\_create creates a L. A. object on GPU.
- CUPLA\_Obj\_free destroys an object on GPU.
- CUPLA\_Chol operates with objects created on GPU.

### **Communication layer**

ĴΊ

Communication is restricted to routines PLA\_Copy and PLA\_Reduce:

int PLA\_Copy ( PLA\_Obj obj\_from, PLA\_Obj obj\_to )

Purpose: Copy contents between linear algebra objects.

IN obj\_from Object to be copied IN/OUT obj\_to Object into which to copy

Purpose: Reduce using op the contents in the duplicated object given by obj\_from and overwrite obj\_to with the result.

IN	obj_from	Object to be reduced
IN		Reduce operator to be used
IN/OUT	obj_to	Object into which to reduce result

### **Communication layer**

j'i

Communication is restricted to routines PLA\_Copy and PLA\_Reduce:

int PLA_Copy ( PLA_Obj ob	j_from, PLA_Obj obj_to )
---------------------------	--------------------------

Purpose: Copy contents between linear algebra objects.

IN obj\_from Object to be copied IN/OUT obj\_to Object into which to copy

int PLA_Reduce ( PLA_Obj obj_from, MPI_OP op, PLA_Obj obj_to )						
Purpose:	Purpose: Reduce using op the contents in the duplicated object given by obj_from and overwrite obj_to with the result.					
	IN IN IN/OUT	obj_from op obj_to	Object to be reduced Reduce operator to be used Object into which to reduce result			

### Modifications to PLA\_Copy (I)

```
1
     /* PLAPACK PLA Copy between column aligned matrices */
 2
 3
     while ( TRUE ){
 4
       PLA Obj split size ( from cur, PLA SIDE TOP, & size from, & owner from );
 5
       PLA Obj split size( to cur, PLA SIDE TOP, & size to, & owner to);
 6
 7
       if ( 0 == ( size = min( size from, size to ) ) ) break;
8
9
       PLA Obj horz split 2( from cur, size, &from 1, &from cur );
10
       PLA Obj horz split 2( to cur, size, &to 1, &to cur );
11
12
       if ( myrow == owner from && owner from == owner to ){
13
         PLA Local copy( from 1, to 1 );
14
15
       else{
16
         if ( myrow == owner from ) {
17
           PLA Obj get local contents (from 1, PLA NO TRANS, &dummy, &dummy,
18
                                        buffer temp, size, 1 );
19
20
           MPI Send( BF( buffer temp ), size * local width, datatype,
                     owner to, 0, comm col );
21
22
23
         if (myrow == owner to ) {
24
           MPI Recv( BF( buffer temp ), size * local width, datatype.
25
                     owner from, MPI ANY TAG, comm col, &status );
26
27
           PLA Obj set local contents (PLA NO TRANS, size, local width,
28
                                        buffer temp, size, 1, to 1):
29
30
31
32
     PLA free( buffer temp ):
```

Ϊī

### Modifications to PLA\_Copy (II)

1	/* PLAPACK PLA	Obi get local contents */		
2				
3 4 5 6 7 8 9	int PLA_Obj_get PLA_Obj int int int void int int	_local_contents ( obj, trans, *rows_in_buf, *cols_in_buf, *buf, ldim_buf, stride_but)	int CUPLA_Obj_ge PLA_Obj int int int void int int	
11		011100_0017		
12	/**********	******		
13				
14 15 16 17 18 19	<pre>for ( j=0; j<n =="" memcpy(="" pre="" tmp_="" tmp_local="tmp_obj" }<=""></n></pre>	; j++ ){ buf_local + j*ldim_buf*typesize; buf_obj + j*ldim*typesize; local, tmp_obj, m*typesize );	cublasGetMatrix (	
20	/***********	************/		

#### Necessary GPU $\rightleftharpoons$ CPU transfers

Data is transferred to RAM prior to a communication. Data is transferred to GPU after each communication.

### Modifications to PLA\_Copy (II)

		Ohi ant land antiput		
	/* FLAFAON FLA_	_ODJ_get_local_contents */	/* CUFLAFACK FLA_	_ODJ_get_local_contents */
2	int PLA_Obj_get	t_local_contents (	int CUPLA_Obj_ge	t_local_contents (
4	PLA_Obj	obj,	PLA_Obj	obj,
5	int	trans ,	int	trans,
6	int	∗rows_in_buf,	int	<pre>*rows_in_buf ,</pre>
7	int	*cols_in_buf,	int	∗cols_in_buf,
8	void	*buf,	void	*buf,
9	int	ldim_buf,	int	ldim_buf,
0	int	stride_buf)	int	stride_buf)
1				
2	/***********	**************/	/************	***********/
3				
4	for ( j=0; j <r< td=""><td>n; j++ ){</td><td></td><td></td></r<>	n; j++ ){		
5	tmp_local =	buf_local + j*ldim_buf*typesize;	cublasGetMatrix	(m, n, typesize,
6	tmp_obj =	buf_obj + j*ldim*typesize;		buf_obj, ldim,
7	memcpy(tmp	_local, tmp_obj, m∗typesize );		buf_local, ldim_buf );
8	}			
9				
20	/**********	**************/	/************	***********/

#### Necessary GPU $\rightleftharpoons$ CPU transfers

Data is transferred to RAM prior to a communication. Data is transferred to GPU after each communication.

### The LONGHORN visualization cluster

	Per Node	Per System (256 nodes)		
Number of cores	8	2,048		
CPU	Intel Xeon Nehalem @ 2.53 GHz			
Available memory	48 Gbytes	13.5 TBytes		
Interconnection network	QDR Infiniband			
Graphics system	128 NVIDIA Quadro Plex S4s			
GPU	2 x NVIDIA Quadro FX5800	512 x NVIDIA Quadro FX5800		
Interconnection bus	PCIExpress 2.0 (8x)			
Available video memory	8 Gbytes DDR3	2 TBytes DDR3		
Peak performance (SP)	161.6 GFLOPS	41.40 TFLOPS		
Peak performance (DP)	80.8 GFLOPS	20.70 TFLOPS		

Table: Detailed features of the LONGHORN cluster.

#### PLAPACK: version R32.

- BLAS: NVIDIA CUBLAS 2.2, MKL 10.1.
  - MPI: MVAPICH 1.4.

### **GEMM:** performance



Figure: Performance of CUPLAPACK for the matrix-matrix multiplication on LONGHORN.

### Cholesky factorization: performance



Figure: Performance of CUPLAPACK for the Cholesky factorization on LONGHORN.

### GEMM: scalability





Figure: Scalability of the PLAPACK-based codes for the matrix-matrix multiplication on LONGHORN. The reference is the performance of CUBLAS SGEMM on one GPU.

### GEMM: PLAPACK vs. CUPLAPACK



Figure: GEMM on 16 nodes of LONGHORN. PLAPACK vs CUPLAPACK.

### Cholesky factorization: PLAPACK vs. CUPLAPACK 🍱



Figure: Cholesky factorization on 16 nodes of LONGHORN. PLAPACK vs CUPLAPACK.

## Multiple GPUs per node vs. One GPU per node 🍡 🎽



Figure: Cholesky factorization on 32 GPUs of LONGHORN, using **one** or **two** GPUs per node.

### Conclusions and future work

#### Conclusions

- Retarget of PLAPACK to hybrid architectures.
- Transparent for programmers and library developers.
- Benefits from the layered design of PLAPACK.
- Remarkable performance and scalability on large GPU clusters.

#### Future work

- Port the Elemental framework to GPUs.
- Out-of-core + GPU support to deal with even larger problems.

### Acknowledgements



Thank you...

#### Questions?

We thank the Texas Advanced Computing Center for granting access to the platform where the experiments were conducted<sup>a</sup>.

<sup>a</sup>http://www.tacc.utexas.edu/resources/visualization/#longhorn

