Use of GPUs in Dense Linear Algebra

Enrique S. Quintana-Ortí

Universitat Jaume I Spain

April, 2008

Use of GPUs in Dense Linear Algebra

- < 同 > < 三 > < 三 >

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

▲ 同 ▶ → 三 ▶

Motivation



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

イロト イポト イヨト イヨト

Outline



2 Evaluation and Tuning of Level 3 (CU)BLAS

Oesign and Implementation of LAPACK: linear systems

4 Multiple GPUs

・ 同 ト ・ 三 ト ・ 三

Introduction

Evaluation and Tuning of Level 3 (CU)BLAS Design and Implementation of LAPACK: linear systems Multiple GPUs

Introduction

 Graphics processors (GPUs) have emerged as an interesting hardware accelerator for "general-purpose" computations (GPGPU):



2 Low cost

• Applied to physical simulations, real-time image processing, linear algebra, signal processing, . . .

▲ 同 ▶ → 三 ▶

CUDA hardware

- 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 can be seen 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

イロト イポト イヨト イヨト

 Local 32-bit registers per processor

CUDA software

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

CUDA libraries

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

イロト イポト イヨト イヨト

Introduction

Evaluation and Tuning of Level 3 (CU)BLAS Design and Implementation of LAPACK: linear systems Multiple GPUs

CUBLAS. Example

```
int main( void ){
                                        A typical CUDA (and
float * h_vector, * d_vector;
                                        CUBLAS) program has 3
h_vector = (float *) malloc (M*sizeof (float));
                                         phases:
cublasAlloc(M. sizeof(float).
                                          Allocation and transfer of
           (void **) &d_vector);
                                              data to GPU
cublasSetVector(M, sizeof(float), h_vector,
               d_vector. 1):
                                          Execution of the kernel
cublasSscal(M, ALPHA, d_vector, 1);
cublasGetVector(M, sizeof(float), d_vector,
                                              (or BLAS routine)
               h_vector. 1):
                                          Transfer of the results
cublasFree(d_vector);
                                              back to main memory
```

(日) (同) (日) (日)

Introduction

Evaluation and Tuning of Level 3 (CU)BLAS Design and Implementation of LAPACK: linear systems Multiple GPUs

Experimental setup

The nVIDIA 8800 Ultra is a CUDA architecture

	CPU	GPU
Processor	Intel Core 2 Duo	NVIDIA 8800 Ultra
Codename	Crusoe E6320	G80
Clock frequency	1.86 GHz	575 MHz
Memory speed	$2 \times 333 \text{ MHz}$	$2 \times 900 \text{ MHz}$
Bus width	64 bits	384 bits
Max. bandwidth	5.3 GB/s	86.4 GB/s
Memory	1024 MB DDR2	768 MB GDDR3
Bus	PCI Express x16 (4 GB/s)	

Use of MKL on the CPU, and CUDA and CUBLAS 1.0 on the GPU

- 4 同 6 4 日 6 4 日 6

Evaluation and tuning of Level 3 (CU)BLAS

GEMM

$$C := \beta \cdot C + \alpha \cdot op(A) \cdot op(B)$$

SYRK

$$C := \beta \cdot C + \alpha \cdot A \cdot A^{T} \text{ or} C := \beta \cdot C + \alpha \cdot A^{T} \cdot A$$

where op(X) = X or X^T

Others

Similar results observed for TRSM, and are to be expected for TRMM and SYMM

Image: A (1) → A (

GEMM evaluation. Experimental results



Use of GPUs in Dense Linear Algebra

- 4 同 6 4 日 6 4 日 6

3

GEMM evaluation. Main remarks

Main remarks

• Peak performance is $\sim \! 116$ Gflops for GEMM

(日) (同) (三) (三)

GEMM evaluation. Main remarks

Main remarks

- $\bullet\,$ Peak performance is ${\sim}116$ Gflops for GEMM
- Results attained for big matrices are much better than those for small-medium sized matrices

Stream-oriented architecture

GEMM evaluation. Main remarks

Main remarks

- $\bullet\,$ Peak performance is ${\sim}116$ Gflops for GEMM
- Results attained for big matrices are much better than those for small-medium sized matrices

\Downarrow

Stream-oriented architecture

- The peak observed for *m* = 4000 is also observed for all dimensions multiple of 32
- Proposal: apply padding!

< 🗇 🕨 < 🖃 🕨

GEMM padding. Experimental results



Use of GPUs in Dense Linear Algebra

- 4 同 6 4 日 6 4 日 6

SYRK evaluation. Experimental results





(日) (同) (三) (三)

SYRK evaluation. Main remarks

Main remarks

• Peak performance is ${\sim}40$ Gflops for SYRK \Rightarrow Suboptimal implementation

Use of GPUs in Dense Linear Algebra

(日) (同) (三) (三)

SYRK evaluation. Main remarks

Main remarks

- Peak performance is \sim 40 Gflops for SYRK \Rightarrow Suboptimal implementation
- As for GEMM, results attained for big matrices are better

A (1) > A (1) > A

SYRK evaluation. Main remarks

Main remarks

- Peak performance is \sim 40 Gflops for SYRK \Rightarrow Suboptimal implementation
- As for GEMM, results attained for big matrices are better
- Proposal: build SYRK on top of GEMM and use padding

Image: A = A

SYRK tuning

• Partition the SYRK operands:



• The second block of columns of *C* is given by:

$$C_{11} := \beta \cdot C_{11} + \alpha \cdot A_1 \cdot A_1^T$$
$$C_{21} := \beta \cdot C_{21} + \alpha \cdot A_2 \cdot A_1^T$$

◆ 同 → ◆ 三

SYRK tuning. Experimental results



Use of GPUs in Dense Linear Algebra

A (1) > A (1) > A

Partitioning for larger matrices

- Transfer times GPU <---> CPU are an important bottleneck
- Our blocked version of GEMM allows to overlap:
 - The partial multiplication A_p and B_p and
 - The transference of the next pair of blocks A_{p+1} and B_{p+1} , being A_p and B_p blocks of columns of A and B, respectively.
- Nor CUDA 1.0 neither G80 allow overlapping of communication and calculation (supported by more recent systems!)
- An orthogonal benefit: enables computation with larger matrices

< (1) > < (2) > <

Hybrid computation

- While GPU is performing a calculation, CPU can also be executing part of the computation.
- We implement a hybrid GEMM implementation following the decomposition:



• Values for N' and N'' must be selected carefully to balance the computation load.

▲ 同 ▶ → 三 ▶

Hybrid computation. Experimental results



SGEMM

Use of GPUs in Dense Linear Algebra

- < 同 > < 三 > < 三 >

Design and implementation of LAPACK: linear systems

Cholesky factorization

$$A = LL^T \quad (\text{or } A = U^T U)$$

with L lower triangular (or U upper triangular)

Others

Similar results obtained for the LU factorization with partial pivoting, and are to be expected for the QR factorization (linear least-squares problems)

- < 同 > < 三 > < 三 >

Cholesky factorization. Blocked variants

Algorithm: $A := CHOL_BLK(A)$ Partition where ... while $m(A_{TI}) < m(A)$ do Determine block size b Repartition $\begin{pmatrix} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{pmatrix} \rightarrow \begin{pmatrix} A_{00} & A_{01} & A_{02} \\ \hline A_{10} & A_{11} & A_{12} \\ \hline A_{10} & A_{11} & A_{12} \end{pmatrix}$ where A_{11} is $b \times b$ Variant 2: Variant 1: Variant 3: $\begin{vmatrix} \mathbf{A}_{10} := \mathbf{A}_{10} \operatorname{TRIL} (\mathbf{A}_{00})^{-\mathrm{T}} \\ \mathbf{A}_{11} := \mathbf{A}_{11} - \mathbf{A}_{10} \mathbf{A}_{10}^{\mathrm{T}} \\ \mathbf{A}_{11} := \mathbf{C}_{11} - \mathbf{A}_{10} \mathbf{A}_{10}^{\mathrm{T}} \\ \mathbf{A}_{11} := \operatorname{CHoL_UNB}(\mathbf{A}_{11}) \\ \mathbf{A}_{21} := \mathbf{A}_{21} - \mathbf{A}_{20} \mathbf{A}_{10}^{\mathrm{T}} \\ \mathbf{A}_{$ $A_{11} := CHOL_UNB(A_{11})$ $A_{21} := A_{21} \operatorname{TRIL} (A_{11})^{-T}$ $A_{22} := A_{22} - A_{21}A_{21}^{\mathrm{T}}$ $A_{21} := A_{21} \operatorname{TRIL} (A_{11})^{-\mathrm{T}}$ Continue with endwhile

Cholesky factorization. Experimental results



Use of GPUs in Dense Linear Algebra

◆ 同 → ◆ 三

э

Cholesky padding. Experimental results



Image: A = A

э

Cholesky hybrid and recursive. Experimental results

Cholesky factorization. Recursive and hybrid variants



Use of GPUs in Dense Linear Algebra

◆ 同 → ◆ 三

Iterative refinement for extended precision

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

 \rightarrow 32 bits of accuracy!



Iterative refinement. Experimental results



Use of GPUs in Dense Linear Algebra

◆ 同 → ◆ 三

э

What if multiple GPUs are available?

Already here:

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

Image: A = A

What if multiple GPUs are available?

Already here:

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

How are we going to program these?

A (1) > (1) > (1)

Experimental setup

	CPU	GPU
Processor	16 x Intel Itanium2	NVIDIA Tesla s870 (4 G80)
Clock frequency	1.5 GHz	575 MHz
Use of MKL on the	Itanium vs. CUDA	and CUBLAS 1.1 on the
Tesla		

<ロ> <同> <同> < 回> < 回>

2

SuperMatrix for multithreaded architectures

- 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

SuperMatrix for multithreaded architectures

The FLAME runtime system "pre-executes" the code:

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

- 4 同 6 4 日 6 4 日 6

SuperMatrix for multithreaded architectures



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

・ロト ・同ト ・ヨト ・ヨト

SuperMatrix for multithreaded architectures





▲ 同 → - ▲ 三

э

Porting SuperMatrix to multiple GPUs

- Run-time system (scheduling), storage, and code are independent
- No significative modification to the FLAME codes: Interfacing to CUBLAS

A software effort of two hours!

Porting SuperMatrix. Experimental results



Cholesky factorization on 4 GPU of NVIDIA S870

Hope you enjoyed

Thanks!

Use of GPUs in Dense Linear Algebra

(日) (同) (三) (三)

æ