

## L11: Sparse Linear Algebra on GPUs

CS6963



## Administrative Issues

- Next assignment, triangular solve
  - Due 5PM, Tuesday, March 15
  - handin cs6963 lab 3 <probfile>
- Project proposals
  - Due 5PM, Wednesday, March 7 (hard deadline)
  - handin cs6963 prop <pdffile>

CS6963

2  
L12: Sparse Linear Algebra

## Triangular Solve (STRSM)

```
for (j = 0; j < n; j++)
  for (k = 0; k < n; k++)
    if (B[j*n+k] != 0.0f) {
      for (i = k+1; i < n; i++)
        B[j*n+i] -= A[k*n+i] * B[j*n+k];
    }
```

Equivalent to:  
`cublasStrsm('l' /* left operator */, 'l' /* lower triangular */,  
 'N' /* not transposed */, 'u' /* unit triangular */,  
 N, N, alpha, d_A, N, d_B, N);`

See: <http://www.netlib.org/blas/strsm.f>

CS6963

3  
L11: Dense Linear Algebra

## A Few Details

- C stores multi-dimensional arrays in row major order
- Fortran (and MATLAB) stores multi-dimensional arrays in column major order
  - **Confusion alert:** BLAS libraries were designed for FORTRAN codes, so column major order is implicit in CUBLAS!

CS6963

4  
L11: Dense Linear Algebra

## Dependencies in STRSM

```

for (j = 0; j < n; j++)
  for (k = 0; k < n; k++)
    if (B[j*n+k] != 0.0f) {
      for (i = k+1; i < n; i++)
        B[j*n+i] -= A[k * n + i] * B[j * n + k];
    }

```

Which loop(s) "carry" dependences?  
Which loop(s) is(are) safe to execute in parallel?

CS6963

5  
L11: Dense Linear Algebra

## Assignment

- Details:
  - Integrated with simpleCUBLAS test in SDK
  - Reference sequential version provided
- 1. Rewrite in CUDA
- 2. Compare performance with CUBLAS library

CS6963

6  
L11: Dense Linear Algebra

## Performance Issues?

- + Abundant data reuse
- - Difficult edge cases
- - Different amounts of work for different  $\langle j, k \rangle$  values
- - Complex mapping or load imbalance

CS6963

7  
L11: Dense Linear Algebra

## Outline

- Next assignment
- For your projects:
  - "Debunking the 100X GPU vs. CPU Myth: An Evaluation of Throughput Computing on CPU and GPU", Lee et al., ISCA 2010.
- Sparse Linear Algebra
- Readings:
  - "Implementing Sparse Matrix-Vector Multiplication on Throughput Oriented Processors," Bell and Garland (Nvidia), SC09, Nov. 2009.
  - "Model-driven Autotuning of Sparse Matrix-Vector Multiply on GPUs", Choi, Singh, Vuduc, PPOPP 10, Jan. 2010.
  - "Optimizing sparse matrix-vector multiply on emerging multicore platforms," Journal of Parallel Computing, 35(3):178-194, March 2009. (Expanded from SC07 paper.)

CS6963

8  
L11: Sparse Linear Algebra

## Overview: CPU and GPU Comparisons

- Many projects will compare speedup over a sequential CPU implementation
  - Ok for this class, but not for a research contribution
- Is your CPU implementation as "smart" as your GPU implementation?
  - Parallel?
  - Manages memory hierarchy?
  - Minimizes synchronization or accesses to global memory?



## The Comparison

- Architectures
  - Intel i7, quad-core, 3.2GHz, 2-way hyper-threading, SSE, 32KB L1, 256KB L2, 8MB L3
  - Same i7 with Nvidia GTX 280
- Workload
  - 14 benchmarks, some from the GPU literature



## Architectural Comparison

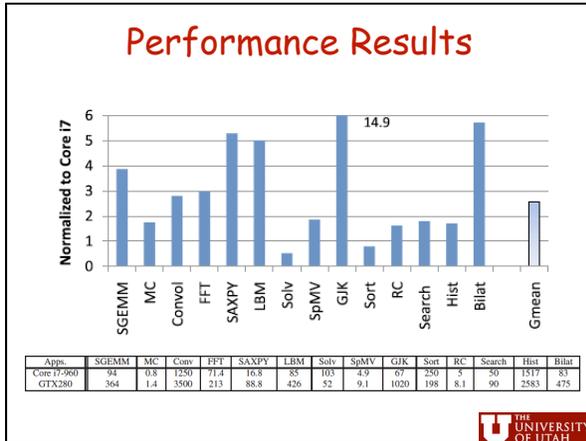
	Core i7-960	GTX280
Number PEs	4	30
Frequency (GHz)	3.2	1.3
Number Transistors	0.7B	1.4B
BW (GB/sec)	32	141
SP SIMD width	4	8
DP SIMD width	2	1
Peak SP Scalar FLOPS (GFLOPS)	25.6	116.6
Peak SP SIMD Flops (GFLOPS)	102.4	311.1/933.1
Peak DP SIMD Flops (GFLOPS)	51.2	77.8



## Workload Summary

Kernel	Application	SIMD	TLP	Characteristics
SGEMM (SGEMV) [148]	Linear algebra	Regular	Across 2D tiles	Compute bound after tiling
Monte Carlo (MC) [184, 9]	Computational Finance	Regular	Across paths	Compute bound
Convolution (Conv) [16, 19]	Image Analysis	Regular	Across pixels	Compute bound; BW bound for small filters
PFT (FFT) [17, 21]	Signal Processing	Regular	Across smaller PFTs	Compute/BW bound depending on size
SAXPY (SAXPY) [46]	Dot Product	Regular	Across vector	BW bound for large vectors
LBM (LBM) [32, 45]	Time Migration	Regular	Across cells	BW bound
Constraint Solver (Solv) [14]	Rigid body physics	Gather/Scatter	Across constraints	Synchronization bound
SpMV (SpMV) [80, 8, 47]	Sparse Solver	Gather	Across non-zero	BW bound for typical large matrices
GJK (GJK) [38]	Collision Detection	Gather/Scatter	Across objects	Compute Bound
Sort (Sort) [18, 39, 40]	Database	Gather/Scatter	Across elements	Compute bound
Ray Casting (RC) [43]	Volume Rendering	Gather	Across rays	4-8MB first level working set, over 500MB last level working set
Search (Search) [27]	Database	Gather/Scatter	Across queries	Compute bound for small tree, BW bound at bottom of tree for large tree
Histogram (Hist) [53]	Image Analysis	Requires conflict detection	Across pixels	Reduction/synchronization bound
Bilateral (Bilat) [52]	Image Analysis	Regular	Across pixels	Compute Bound





### CPU optimization

- Tile for cache utilization
- SIMD execution on multimedia extensions
- Multi-threaded, beyond number of cores
- Data reorganization to improve SIMD performance

### Sparse Linear Algebra

- Suppose you are applying matrix-vector multiply and the matrix has lots of zero elements
  - Computation cost? Space requirements?
- General sparse matrix representation concepts
  - Primarily only represent the nonzero data values
  - Auxiliary data structures describe placement of nonzeros in "dense matrix"

CS6963

15  
L11: Sparse Linear Algebra

### GPU Challenges

- Computation partitioning?
- Memory access patterns?
- Parallel reduction

BUT, good news is that sparse linear algebra performs TERRIBLY on conventional architectures, so poor baseline leads to improvements!

CS6963

16  
L12: Sparse Linear Algebra

### Some common representations

$$A = \begin{bmatrix} 1 & 7 & 0 & 0 \\ 0 & 2 & 8 & 0 \\ 5 & 0 & 3 & 9 \\ 0 & 6 & 0 & 4 \end{bmatrix}$$

$$\text{ptr} = [0 \ 2 \ 4 \ 7 \ 9]$$

$$\text{indices} = [0 \ 1 \ 1 \ 2 \ 0 \ 2 \ 3 \ 1 \ 3]$$

$$\text{data} = [1 \ 7 \ 2 \ 8 \ 5 \ 3 \ 9 \ 6 \ 4]$$

**Compressed Sparse Row (CSR):**  
Store only nonzero elements, with "ptr" to beginning of each row and "indices" representing column.

---

$$\text{data} = \begin{bmatrix} 1 & 7 & * \\ 2 & 8 & * \\ 5 & 3 & 9 \\ 6 & 4 & * \end{bmatrix}$$

$$\text{indices} = \begin{bmatrix} 0 & 1 & * \\ 1 & 2 & * \\ 0 & 2 & 3 \\ 1 & 3 & * \end{bmatrix}$$

**ELL:** Store a set of K elements per row and pad as needed. Best suited when number non-zeros roughly consistent across rows.

**COO:** Store nonzero elements and their corresponding "coordinates".

THE UNIVERSITY OF UTAH

### CSR Example

```

for (j=0; j<nr; j++) {
  for (k = ptr[j]; k<ptr[j+1]-1; k++)
    t[j] = t[j] + data[k] * x[indices[k]];
}
    
```

CS6963
18  
L11: Sparse Linear Algebra

### Summary of Representation and Implementation

Kernel	Granularity	Coalescing	Bytes/Flop	
			32-bit	64-bit
DIA	thread : row	full	4	8
ELL	thread : row	full	6	10
CSR(s)	thread : row	rare	6	10
CSR(v)	warp : row	partial	6	10
COO	thread : nonz	full	8	12
HYB	thread : row	full	6	10

Table 1 from Bell/Garland: Summary of SpMV kernel properties.

THE UNIVERSITY OF UTAH

### Other Representation Examples

- **Blocked CSR**
  - Represent non-zeros as a set of blocks, usually of fixed size
  - Within each block, treat as dense and pad block with zeros
  - Block looks like standard matvec
  - So performs well for blocks of decent size
- **Hybrid ELL and COO**
  - Find a "K" value that works for most of matrix
  - Use COO for rows with more nonzeros (or even significantly fewer)

CS6963
20  
L11: Sparse Linear Algebra

## Stencil Example

What is a 3-point stencil? 5-point stencil?  
7-point? 9-point? 27-point?

Examples:

$$a[i] = (b[i-1] + b[i+1])/2;$$

$$a[i][j] = (b[i-1][j] + b[i+1][j] + b[i][j-1] + b[i][j+1])/4;$$

How is this represented by a sparse matrix?

CS6963

21  
L11: Sparse Linear Algebra

## Stencil Result (structured matrices)

See Figures 11 and 12, Bell and Garland

CS6963

22  
L11: Sparse Linear Algebra

## Unstructured Matrices

See Figures 13 and 14

Note that graphs can also be represented as sparse matrices. What is an adjacency matrix?

CS6963

23  
L11: Sparse Linear Algebra

## PPoPP paper

- What if you customize the representation to the problem?
- Additional global data structure modifications (like blocked representation)?
- Strategy
  - Apply models and autotuning to identify best solution for each application

CS6963

24  
L11: Sparse Linear Algebra

## Summary of Results

BELLPACK (blocked ELLPACK) achieves up to 29 Gflop/s in SP and 15.7 Gflop/s in DP

Up to 1.8x and 1.5x improvement over Bell and Garland.

CS6963

25  
L11: Sparse Linear Algebra

## This Lecture

- Exposure to the issues in a sparse matrix vector computation on GPUs
- A set of implementations and their expected performance
- A little on how to improve performance through application-specific knowledge and customization of sparse matrix representation

CS6963

26  
L11: Sparse Linear Algebra

## What's coming

- Next time: Application case study from Kirk and Hwu (Ch. 8, real-time MRI)
- Wednesday, March 2: two guest speakers from last year's class
  - BOTH use sparse matrix representation!
  - Shreyas Ramalingam: program analysis on GPUs
  - Pascal Grosset: graph coloring on GPUs

CS6963

27  
L11: Sparse Linear Algebra