

## L11: Dense Linear Algebra on GPUs, cont.

CS6963



### Administrative Issues

- Next assignment, triangular solve
  - Due 5PM, Friday, March 5
  - handin cs6963 lab 3 <probfile>
- Project proposals (discussed today)
  - Due 5PM, Wednesday, March 17 (hard deadline)
  - handin cs6963 prop <pdffile>
  - A few minutes at end of class to help form groups

CS6963

2  
L11: Dense Linear Algebra

### Outline

- Project
  - Suggested topics
- Triangular solve assignment
- More on dense linear algebra
  - LU

CS6963

3  
L11: Dense Linear Algebra

### Project Ideas

- Suggestions from Nvidia (Nathan Bell via Pete Shirley, Steve Parker):
  - Sparse solvers
    - Tailor computation to matrix structure
    - Optimized SpMV for multiple vectors
    - Block Krylov methods
    - Block eigensolvers
  - Parallel preconditions
  - Parallel Graph Algorithms
    - (Maximal) Independent sets
    - Graph colorings
    - Partitioning of graphs that minimize a given metric (e.g. edge cut)
    - Optimal network flows
    - Graph clustering

4  
L11: Dense Linear Algebra

## More Project Ideas

- Suggestions from Steve Parker:
  - Non-uniform, unstructured
  - Double precision, cluster-based

5  
L11: Dense 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

6  
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

7  
L11: Dense Linear Algebra

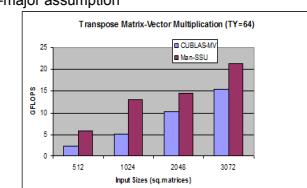


## SGEMV (single-precision)

sgemv (matrix vector multiplication) and transposed sgemv  
 Recall transpose vs. column-major assumption

```
// This is transposed due to
// Fortran assumption
// Follows coalesced accesses
// Place b in shared memory
for (i = 0; i < n; i++)
  for (j = 0; j < n; j++)
    c[i] += a[i * n + j] * b[j];
```

```
// Non-transposed
// Store a in shared memory
// and access in a different order
for (i = 0; i < n; i++)
  for (j = 0; j < n; j++)
    c[i] += a[i * n + j] * b[j];
```



Results from Malik Khan

Non-transposed more than 2X slower for 8K matrix!

Reference: An Optimizing Compiler for GPGPU Programs with Input-Data Sharing, Yang et al., POPL 2010 Poster, Jan. 2010.



## Dependences 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

9  
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 2.0 library

CS6963

10  
L11: Dense Linear Algebra

## Performance Issues?

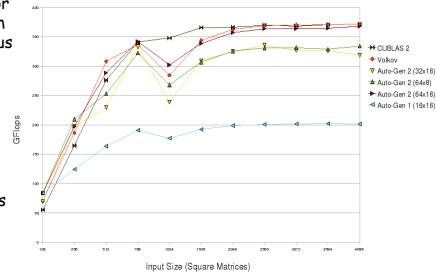
- + Abundant data reuse
- - Difficult edge cases
- - Different amounts of work for different  $\langle j, k \rangle$  values
- - Complex mapping or load imbalance

CS6963

11  
L11: Dense Linear Algebra

## Latest Research: CUDA-CHILL

- Automatically generate CUDA for NVIDIA GPU from sequential code plus script
- Abstraction permits parallel threads & staging of data
- **Heterogeneous code generation:** Alternative scripts generate CUDA, OpenMP or sequential code tuned for memory hierarchy



Results provided by Malik Khan, Gabe Rudy, Chun Chen

CS6963

12  
L11: Dense Linear Algebra

## Recent Results

Suppose your matrix is not evenly divided by the tile sizes (64 and 16)?

512x512 matrix sgemm:  
 CUBLAS 258 Gflops  
 Our compiler 284 Gflops

500x500 matrix sgemm:  
 CUBLAS 167 Gflops  
 Our compiler 223 Gflops

CS6963

13  
L11: Dense Linear Algebra

## Breaking it Down: Tiling

```

for (i = 0; i < n; i++)
    for (j = 0; j < n; j++)
        for (k = 0; k < n; k++)
            c[i][j] += a[k][i]*b[j][k];

1. tile_control({"i","j"}, {TI,TJ},
    {l1_control="ii", l2_control="jj"},
    {"ii", "jj", "i", "j"})
tile_control({"k"}, {TK}, {l1_control="kk"}, 
    {"ii", "jj", "kk", "i", "j", "k"}, strided)
tile_control({"i"}, {TJ},
    {l1_control="ty", l1_tile="tx"}, 
    {"ii", "jj", "kk", "tx", "ty", "j", "k"})

--Assign loop levels to thread space and name the kernel
cudaize("mm_GPU",
    {a=N*N, b=N*N, c=N*N}, --array sizes for data copying
    {block={"ii", "jj"}, thread={"tx", "ty"}})
    
```

CS6963

15  
L11: Dense Linear Algebra

## CUDA-CHILL: Higher-Level Interface (in progress)

```

init("mm.sp2", "MarkedLoop")
tile_control({"i","j"}, {TI,TJ},
    {l1_control="ii", l2_control="jj"}, 
    {"ii", "jj", "i", "j"})
tile_control({"k"}, {TK}, {l1_control="kk"}, 
    {"ii", "jj", "kk", "i", "j", "k"}, strided)
tile_control({"i"}, {TJ},
    {l1_control="ty", l1_tile="tx"}, 
    {"ii", "jj", "kk", "tx", "ty", "j", "k"})

--Assign loop levels to thread space and name the kernel
cudaize("mm_GPU",
    {a=N*N, b=N*N, c=N*N}, --array sizes for data copying
    {block={"ii", "jj"}, thread={"tx", "ty"}})
--Copy the "c" array usage to registers
copy_to_registers("kk", "c", {"tx", "ty"})
copy_to_shared("ty", "b")
--Unroll two innermost loop levels fully
unroll_to_depth(2)

Gabe Rudy Master's thesis
    
```

CS6963

14  
L11: Dense Linear Algebra

## Breaking it Down: Tiling

```

for (i = 0; i < n; i++)
    for (j = 0; j < n; j++)
        for (k = 0; k < n; k++)
            c[i][j] += a[k][i]*b[j][k];

2. Tile k loop (strided, move kk above i&j)
for (ii=0; i<n/TI; i++)
    for (jj=0; j<n/TJ; j++)
        for (kk=0; kk<n; kk+=TK)
            for (i = ii*Tl; i < min(n,(ii+1)*Tl); i++)
                for (j = jj*TJ; j < min(n,(jj+1)*TJ); j++)
                    for (k = kk; k < min(n, kk+TK); k++)
                        c[i][j] += a[k][i]*b[j][k];

1. tile_control({"i","j"}, {TI,TJ},
    {l1_control="ii", l2_control="jj"}, 
    {"ii", "jj", "i", "j"})
2. tile_control({"k"}, {TK}, {l1_control="kk"}, 
    {"ii", "jj", "kk", "i", "j", "k"}, strided)
tile_control({"i"}, {TJ},
    {l1_control="ty", l1_tile="tx"}, 
    {"ii", "jj", "kk", "tx", "ty", "j", "k"})

--Assign loop levels to thread space and name the kernel
cudaize("mm_GPU",
    {a=N*N, b=N*N, c=N*N}, --array sizes for data copying
    {block={"ii", "jj"}, thread={"tx", "ty"}})
    
```

CS6963

16  
L11: Dense Linear Algebra

### Breaking it Down: Tiling

```

for (i = 0; i < n; i++)
  for (j = 0; j < n; j++)
    for (k = 0; k < n; k++)
      c[i][j] += a[k][i]*b[j][k];

```

- tile\_control({{"i", "j"}, {TJ, TJ}}, {l1\_control="ii", l2\_control="jj"}, {"ii", "jj", "i", "j"})
- tile\_control({{"k"}, {TK}}, {l1\_control="kk"}, {"ii", "jj", "kk", "i", "j", "k"}, strided)
- tile\_control({{"i"}, {TJ}}, {l1\_control="ty"}, {"ii", "jj", "kk", "tx", "ty", "j", "k"})

--Assign loop levels to thread space and name the kernel  
 cudaize("mm\_GPU", {a=N\*N, b=N\*N, c=N\*N}, --array sizes for data copying  
 {block={"ii", "jj"}, thread={"tx", "ty"}})

CS6963

17

L11: Dense Linear Algebra



THE UNIVERSITY OF UTAH

### Breaking it Down: "cudaize"

```

for (i = 0; i < n; i++)
  for (j = 0; j < n; j++)
    for (k = 0; k < n; k++)
      c[i][j] += a[k][i]*b[j][k];

```

- tile\_control({{"i", "j"}, {TJ, TJ}}, {l1\_control="ii", l2\_control="jj"}, {"ii", "jj", "i", "j"})
- tile\_control({{"k"}, {TK}}, {l1\_control="kk"}, {"ii", "jj", "kk", "i", "j", "k"}, strided)
- tile\_control({{"i"}, {TJ}}, {l1\_control="ty"}, {ii, jj, kk, tx, ty, j, k})

--Assign loop levels to thread space and name the kernel  
 4. cudaize("mm\_GPU", {a=N\*N, b=N\*N, c=N\*N}, --array sizes for data copying  
 {block={"ii", "jj"}, thread={"tx", "ty"}})

CS6963

18

L11: Dense Linear Algebra



THE UNIVERSITY OF UTAH

### Final Result: Copy C to registers, B to shared memory and unroll

```

--Copy the "c" array usage to registers
5. copy_to_registers("kk", "c", {"tx", "ty"})
6. copy_to_shared("ty", "b")
--Unroll two innermost loop levels fully
7. unroll_to_depth(2)

float P1[16];
__shared__ float P2[16][17];
bx = blockIdx.x, by = blockIdx.y;
tx = threadIdx.x, ty = threadIdx.y;
P1[0:15] = [16*by:16*by+15][tx+64*bx+16*ty];
for (t6 = 0; t10 <= 1088; t6+=16) {
  P2[tx][4*ty:4*ty+3] = b[16*by+4*ty:16*by+4*ty+3]
  [tx+t6];
  __syncthreads();
  P1[0:15] += a[t6][64*bx+16*ty+tx]*P2[0][0:15];
  P1[0:15] += a[t6+1][64*bx+16*ty+tx]*P2[1][0:15];
  ...
  P1[0:15] += a[t6+15][64*bx+16*ty+tx]*P2[15][0:15];
  __syncthreads();
}
c[16*by:16*by+15][tx+64*bx+16*ty] = P1[0:15];

```

B goes into shared memory

C goes into registers and is copied back at end

A is used from registers, but not reused

CS6963

19

L11: Dense Linear Algebra



### LU Decomposition and Other Matrix Factorizations

- Volkov talks about matrix factorization algorithms
- What are these:
  - LU Decomposition example
  - Also, Cholesky and QR in paper

CS6963

20

L11: Dense Linear Algebra



THE UNIVERSITY OF UTAH

### LU Decomposition (no pivoting)

```
DO K=1,N-1
DO I=K+1,N
  A(I,K)=A(I,K)/A(K,K)
DO I=K+1,N
  DO J=K+1,N
    A(I,J)=A(I,J)-A(I,K)*A(K,J)
```

- LAPACK version decomposes this into three parts: mini-LU, (S/D)TRSM, & (S/D) GEMM
- Row pivoting (used for numeric stability) reorders matrix and computation

CS6963

21  
L11: Dense Linear Algebra

### What's Coming

- See Khan/Rudy poster on Friday!
- Before Spring Break
  - Two application case studies from newly published Kirk and Hwu
  - Sparse linear algebra
  - Review for Midterm

CS6963

22  
L11: Dense Linear Algebra