

L13: Application Case Studies

CS6963



Administrative Issues

- Next assignment, triangular solve
 - Due 5PM, Monday, March 8
 - handin cs6963 lab 3 <probfile>
- Project proposals
 - Due 5PM, Wednesday, March 17 (hard deadline)
 - handin cs6963 prop <pdffile>

CS6963

2
L13: Application Case Studies

Outline

- Discussion of strsm (more on Monday)
- How to approach your projects
- Application Case Studies
 - Advanced MRI Reconstruction
 - Reading: Kirk and Hwu, Chapter 7, <http://courses.ece.illinois.edu/ece498/al/textbook/Chapter7-MRI-Case-Study.pdf>

CS6963

3
L13: Application Case Studies

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

4
L13: Application Case Studies

Approaching Projects/STRSM/Case Studies

1. Parallelism?
 - How do dependences constrain partitioning strategies?
2. Analyze data accesses for different partitioning strategies
 - Start with global memory: coalesced?
 - Consider reuse: within a thread? Within a block? Across blocks?
3. Data Placement (adjust partitioning strategy)?
 - Registers, shared memory, constant memory, texture memory or just leave in global memory
4. Tuning
 - Unrolling, fine-tune partitioning, floating point, control flow,

CS6963

L13: Application Case Studies



5

Step 1. Simple Partition for STRSM

```
__global__ void strsm1( int n, float *A, float *B )
{
    int bx = blockIdx.x;
    int tx = threadIdx.x;
    int j = bx*THREADSPERBLOCK + tx; // one thread per column, columns work
    independently
    int JN = j * n;
    int i, k;

    for (k = 0; k < n; ++k) { // ROW
        int KN = k * n;
        for (i = k+1; i < n; ++i) { // ALSO row
            // B[i][j] := A[i][k] * B[k][j] element depends on elts in ROWS above it in same col
            B[ JN + i ] -= A[ KN + i ] * B[ JN + k ];
        }
    }
}
```

CS6963

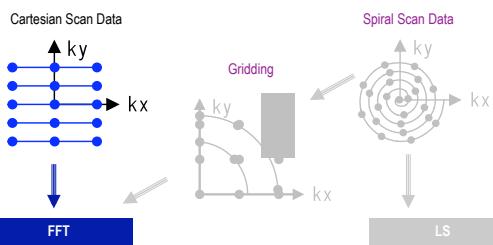
Slide source: Mark Hall

6

L13: Application Case Studies



Reconstructing MR Images

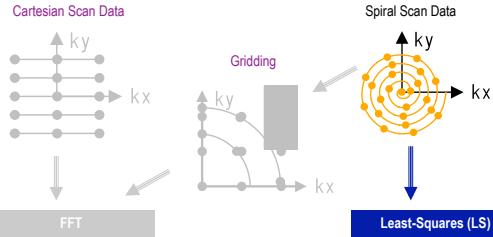


© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

7



Reconstructing MR Images



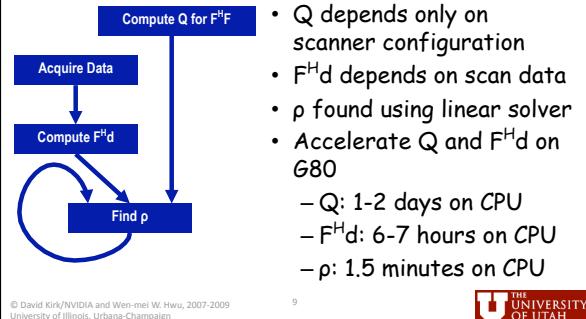
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

8



Least-Squares Reconstruction

$$F^H F \rho = F^H d$$



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

9



<pre> for (m = 0; m < M; m++) { phiMag[m] = rPhi[m]*rPhi[m] + iPhi[m]*iPhi[m]; for (n = 0; n < N; n++) { expQ = 2*PI*(kx[m]*x[n] + ky[m]*y[n] + kz[m]*z[n]); rQ[n] += phiMag[m]*cos(expQ); iQ[n] += phiMag[m]*sin(expQ); } } (a) Q computation </pre>	<pre> for (m = 0; m < M; m++) { rMu[m] = rPhi[m]*rD[m] + iPhi[m]*iD[m]; iMu[m] = rPhi[m]*iD[m] - iPhi[m]*rD[m]; for (n = 0; n < N; n++) { expFhD = 2*PI*(kx[m]*x[n] + ky[m]*y[n] + kz[m]*z[n]); rFhD[n] += rMu[m]*cArg - iMu[m]*sArg; iFhD[n] += iMu[m]*cArg + rMu[m]*sArg; } } (b) F^H d computation </pre>
---	--

Q v.s. $F^H d$

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

10



Algorithms to Accelerate

```

for (m = 0; m < M; m++) {
    rMu[m] = rPhi[m]*rD[m] +
              iPhi[m]*iD[m];
    iMu[m] = rPhi[m]*iD[m] -
              iPhi[m]*rD[m];

    for (n = 0; n < N; n++) {
        expFhD = 2*PI*(kx[m]*x[n] +
                        ky[m]*y[n] +
                        kz[m]*z[n]);
        cArg = cos(expFhD);
        sArg = sin(expFhD);

        rFhD[n] += rMu[m]*cArg -
                    iMu[m]*sArg;
        iFhD[n] += iMu[m]*cArg +
                    rMu[m]*sArg;
    }
}
      
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

- Scan data
 - M = # scan points
 - kx, ky, kz = 3D scan data
- Pixel data
 - N = # pixels
 - x, y, z = input 3D pixel data
 - rFhD, iFhD= output pixel data
- Complexity is O(MN)
- Inner loop
 - 13 FP MUL or ADD ops
 - 2 FP trig ops
 - 12 loads, 2 stores



Step 1. Consider Parallelism to Evaluate Partitioning Options

```

for (m = 0; m < M; m++) {
    rMu[m] = rPhi[m]*rD[m] +
              iPhi[m]*iD[m];
    iMu[m] = rPhi[m]*iD[m] -
              iPhi[m]*rD[m];

    for (n = 0; n < N; n++) {
        expFhD = 2*PI*(kx[m]*x[n] +
                        ky[m]*y[n] +
                        kz[m]*z[n]);
        cArg = cos(expFhD);
        sArg = sin(expFhD);

        rFhD[n] += rMu[m]*cArg -
                    iMu[m]*sArg;
        iFhD[n] += iMu[m]*cArg +
                    rMu[m]*sArg;
    }
}
      
```

What about M total threads?

Note: M is O(millions)

(Step 2) What happens to data accesses with this strategy?



One Possibility

```
__global__ void cmpFHD(float* rPhi, iPhi, phiMag,
    kx, ky, kz, x, y, z, rMu, iMu, int N) {

    int m = blockIdx.x * FHD_THREADS_PER_BLOCK + threadIdx.x;

    rMu[m] = rPhi[m]*rD[m] + iPhi[m]*iD[m];
    iMu[m] = rPhi[m]*iD[m] - iPhi[m]*rD[m];

    for (n = 0; n < N; n++) {
        expFHD = 2*PI*(kx[m]*x[n] + ky[m]*y[n] + kz[m]*z[n]);

        cArg = cos(expFHD); sArg = sin(expFHD);

        rFHD[n] += rMu[m]*cArg - iMu[m]*sArg; This code does not
        iFHD[n] += iMu[m]*cArg + rMu[m]*sArg; work correctly! The
    }                                         accumulation needs to
                                                use atomic operation.
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

13



Back to the Drawing Board - Maybe map the n loop to threads?

```
for (m = 0; m < M; m++) {

    rMu[m] = rPhi[m]*rD[m] + iPhi[m]*iD[m];
    iMu[m] = rPhi[m]*iD[m] - iPhi[m]*rD[m];

    for (n = 0; n < N; n++) {
        expFHD = 2*PI*(kx[m]*x[n] + ky[m]*y[n] + kz[m]*z[n]);
        cArg = cos(expFHD);
        sArg = sin(expFHD);

        rFHD[n] += rMu[m]*cArg - iMu[m]*sArg;
        iFHD[n] += iMu[m]*cArg + rMu[m]*sArg;
    }
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

14



for (m = 0; m < M; m++) {

 rMu[m] = rPhi[m]*rD[m] +
 iPhi[m]*iD[m];
 iMu[m] = rPhi[m]*iD[m] -
 iPhi[m]*rD[m];

 for (n = 0; n < N; n++) {
 expFHD = 2*PI*(kx[m]*x[n] +
 ky[m]*y[n] +
 kz[m]*z[n]);

 cArg = cos(expFHD);
 sArg = sin(expFHD);

 rFHD[n] += rMu[m]*cArg -
 iMu[m]*sArg;
 iFHD[n] += iMu[m]*cArg +
 rMu[m]*sArg;
 }
} (a) FHD computation

for (m = 0; m < M; m++) {

 rMu[m] = rPhi[m]*rD[m] +
 iPhi[m]*iD[m];
 iMu[m] = rPhi[m]*iD[m] -
 iPhi[m]*rD[m];

 for (n = 0; n < N; n++) {
 expFHD = 2*PI*(kx[m]*x[n] +
 ky[m]*y[n] +
 kz[m]*z[n]);

 cArg = cos(expFHD);
 sArg = sin(expFHD);

 rFHD[n] += rMu[m]*cArg -
 iMu[m]*sArg;
 iFHD[n] += iMu[m]*cArg +
 rMu[m]*sArg;
 }
} (b) after loop fission

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

15



A Separate cmpMu Kernel

```
__global__ void cmpMu(float* rPhi, iPhi, rD, iD, rMu, iMu)
{
    int m = blockIdx.x * MU_THREADS_PER_BLOCK + threadIdx.x;

    rMu[m] = rPhi[m]*rD[m] + iPhi[m]*iD[m];
    iMu[m] = rPhi[m]*iD[m] - iPhi[m]*rD[m];
}
```

© David Kirk/NVIDIA and Wen-
mei W. Hwu, 2007-2009
University of Illinois, Urbana-
Champaign

16



```

for (m = 0; m < M; m++) {
    for (n = 0; n < N; n++) {
        expFhD = 2*PI*(kx[m]*x[n] +
                      ky[m]*y[n] +
                      kz[m]*z[n]);
        cArg = cos(expFhD);
        sArg = sin(expFhD);

        rFhD[n] += rMu[m]*cArg -
                    iMu[m]*sArg;
        iFhD[n] += iMu[m]*cArg +
                    rMu[m]*sArg;
    }
} (a) before loop interchange

```



```

for (n = 0; n < N; n++) {
    for (m = 0; m < M; m++) {
        expFhD = 2*PI*(kx[m]*x[n] +
                      ky[m]*y[n] +
                      kz[m]*z[n]);
        cArg = cos(expFhD);
        sArg = sin(expFhD);

        rFhD[n] += rMu[m]*cArg -
                    iMu[m]*sArg;
        iFhD[n] += iMu[m]*cArg +
                    rMu[m]*sArg;
    }
} (b) after loop interchange

```

Figure 7.9 Loop interchange of the F^HD computation

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

17



Step 2. New FHd kernel

```

__global__ void cmpFHd(float*
    kx, ky, kz, x, y, z, rMu, iMu, int M) {
    int n = blockIdx.x * FHD_THREADS_PER_BLOCK + threadIdx.x;
    for (m = 0; m < M; n++) {
        float expFhD = 2*PI*(kx[m]*x[n]+ky[m]*y[n]+kz[m]*z[n]);

        float cArg = cos(expFhD);
        float sArg = sin(expFhD);

        rFhD[n] += rMu[m]*cArg - iMu[m]*sArg;
        iFhD[n] += iMu[m]*cArg + rMu[m]*sArg;
    }
}

```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

18



Step 3. Using Registers to Reduce Global Memory Traffic

```

__global__ void cmpFHd(float* rPhi, iPhi, phiMag,
    kx, ky, kz, x, y, z, rMu, iMu, int M) {
    int n = blockIdx.x * FHD_THREADS_PER_BLOCK + threadIdx.x;
    float xn_r = x[n]; float yn_r = y[n]; float zn_r = z[n];
    float rFhDn_r = rFhD[n]; float iFhDn_r = iFhD[n];

    for (m = 0; m < M; m++) {
        float expFhD = 2*PI*(kx[m]*xn_r+ky[m]*yn_r+kz[m]*zn_r);

        float cArg = cos(expFhD);
        float sArg = sin(expFhD);
        Still too much stress
        on memory! Note that
        kx, ky and kz are read-
        only and based on m

        rFhDn_r += rMu[m]*cArg - iMu[m]*sArg;
        iFhDn_r += iMu[m]*cArg + rMu[m]*sArg;
    }
    rFhD[n] = rFhD_r; iFhD[n] = iFhD_r;
}

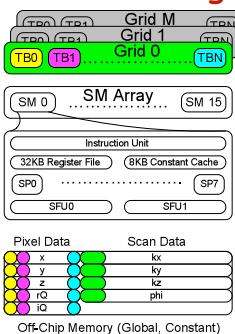
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

19



Tiling of Scan Data



- LS recon uses multiple grids
- Each grid operates on all pixels
- Each grid operates on a distinct subset of scan data
- Each thread in the same grid operates on a distinct pixel

Thread n operates on pixel n:

```

for (m = 0; m < M/32; m++) {
    exQ = 2*PI*(kx[m]*x[n] +
                 ky[m]*y[n] +
                 kz[m]*z[n])
    rQ[n] += phi[m]*cos(exQ)
    iQ[n] += phi[m]*sin(exQ)
}

```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

20



Tiling k-space data to fit into constant memory

```
__constant__ float    kx_c[CHUNK_SIZE],
              ky_c[CHUNK_SIZE], kz_c[CHUNK_SIZE];
...
__ void main() {

    int i;
    for (i = 0; i < M/CHUNK_SIZE; i++) {
        cudaMemcpyToSymbol(kx_c, &kx[i*CHUNK_SIZE], 4*CHUNK_SIZE);
        cudaMemcpyToSymbol(ky_c, &ky[i*CHUNK_SIZE], 4*CHUNK_SIZE);
        cudaMemcpyToSymbol(kz_c, &ky[i*CHUNK_SIZE], 4*CHUNK_SIZE);
        ...
        cmpFHD<<<N/FHD_THREADS_PER_BLOCK, FHD_THREADS_PER_BLOCK>>>
            (rPhi, iPhi, phiMag, x, y, z, rMu, iMu, int M);
    }
    /* Need to call kernel one more time if M is not */
    /* perfect multiple of CHUNK SIZE */
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

21



Revised Kernel for Constant Memory

```
__global__ void cmpFHD(float*
    x, y, z, rMu, iMu, int M) {

    int n = blockIdx.x * FHD_THREADS_PER_BLOCK + threadIdx.x;

    float xn_r = x[n]; float yn_r = y[n]; float zn_r = z[n];
    float rfHdN_r = rfHd[n]; float iFhDn_r = iFhD[n];

    for (m = 0; m < M; m++) {
        float expFhd = 2*PI*(kx_c[m]*xn_r+ky_c[m]*yn_r
        +kz_c[m]*zn_r);

        float cArg = cos(expFhd);
        float sArg = sin(expFhd);

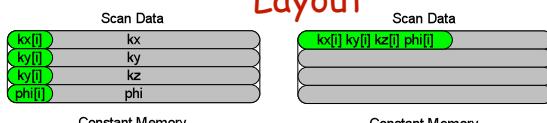
        rfHdN_r += rMu[m]*cArg - iMu[m]*sArg;
        iFhDn_r += iMu[m]*cArg + rMu[m]*sArg;
    }
    rfHd[n] = rfHd_r; iFhD[n] = iFhD_r;
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009

22



Sidebar: Cache-Conscious Data Layout



- kx, ky, kz, and phi components of same scan point have spatial and temporal locality
 - Prefetching
 - Caching
- Old layout does not fully leverage that locality
- New layout does fully leverage that locality

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

23



Adjusting K-space Data Layout

```
struct kdata {
    float x, float y, float z;
} k;

__constant__ struct kdata k_c[CHUNK_SIZE];
...

__ void main() {

    int i;

    for (i = 0; i < M/CHUNK_SIZE; i++) {
        cudaMemcpyToSymbol(k_c, k, 12*CHUNK_SIZE);

        cmpFHD<<<FHD_THREADS_PER_BLOCK,N/FHD_THREADS_PER_BLOCK>>>
            ();
    }
}
```

© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

24



```

__global__ void cmpFHD(float* rPhi, iPhi, phiMag,
    x, y, z, rMu, iMu, int M) {

    int n = blockIdx.x * FHD_THREADS_PER_BLOCK + threadIdx.x;

    float xn_r = x[n]; float yn_r = y[n]; float zn_r = z[n];
    float rFhDn_r = rFhD[n]; float iFhDn_r = iFhD[n];

    for (m = 0; m < M; m++) {
        float expFhd = 2*PI*(k[m].x*xn_r+k[m].y*yn_r+k[m].z*zn_r);

        float cArg = cos(expFhd);
        float sArg = sin(expFhd);

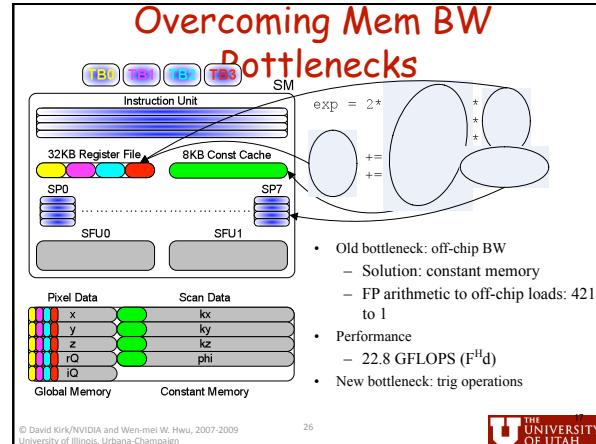
        rFhDn_r += rMu[m]*cArg - iMu[m]*sArg;
        iFhDn_r += iMu[m]*cArg + rMu[m]*sArg;
    }
    rFhD[n] = rFhD_r; iFhD[n] = iFhD_r;
}

```

Figure 7.16 Adjusting the k-space data memory layout in the F^HD kernel

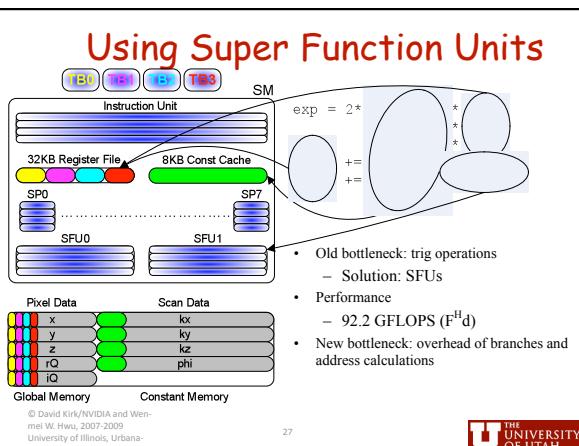
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

25



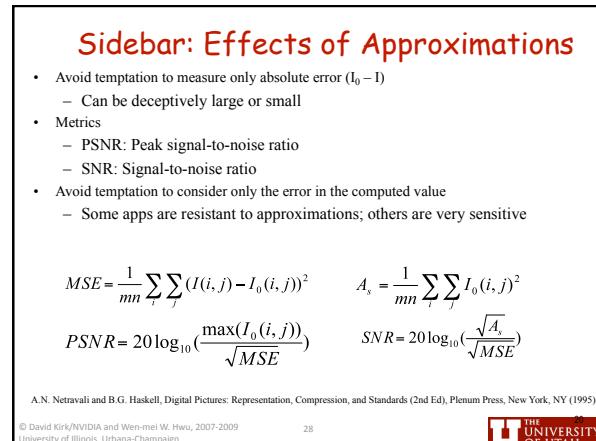
© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

26



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

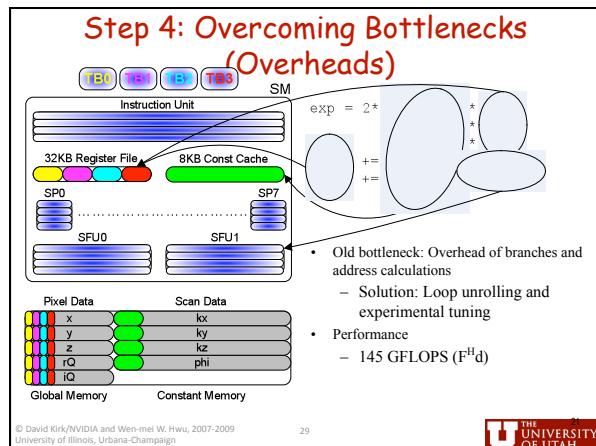
27



© David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009
University of Illinois, Urbana-Champaign

28





What's Coming

- Before Spring Break
 - Complete MRI application study + more STRSM
 - MPM/GIMP application study from last year (read)

GPU Acceleration of the Generalized Interpolation Material Point Method Wei-Fan Chiang, Michael DeLisi, Todd Hummel, Tyler Prete, Kevin Tew, Mary Hall, Phil Wallstedt, and James Guilkey (http://saahpc.ncsa.illinois.edu/09/papers/Chiang_paper.pdf)

- Review for Midterm (week after spring break)

CS6963

30
L13: Application Case Studies