

GPU Teaching Kit

Accelerated Computing



Module 15 - Application Case Study – Advanced MRI Reconstruction Lecture 15.2 – Kernel Optimizations

Objective

- To learn how to apply parallel programming techniques to an application
 - Determining parallelism structure
 - Loop transformations
 - Memory layout considerations
 - Validation



First Version of the FHD Kernel.

```
__global___ void cmpFhD(float* rPhi, iPhi, rD, iD,
     kx, ky, kz, x, y, z, rMu, iMu, rFhD, iFhD, 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 (int n = 0; n < N; n++) {
   float expFhD = 2*PI*(kx[m]*x[n] + ky[m]*y[n] + kz[m]*z[n]);
   float cArq = cos(expFhD); float sArq = sin(expFhD);
  rFhD[n] += rMu[m]*cArq - iMu[m]*sArq;
   iFhD[n] += iMu[m]*cArg + rMu[m]*sArg;
```

Loop Fission

```
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]*cArq -
                iMu[m]*sArg;
    iFhD[n] += iMu[m]*cArq +
                rMu[m]*sArg;
         (a) F<sup>H</sup>D 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 (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]*sArq;
    iFhD[n] +=
               iMu[m]*cArq +
                rMu[m]*sArq;
       (b) after loop fission
```

cmpMu Kernel

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

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

Loop Interchange of the FHD Computation

```
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]*sArq;
    iFhD[n] +=
               iMu[m]*cArq +
                rMu[m]*sArq;
   (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]*cArq +
                rMu[m]*sArq;
   (b) after loop interchange
```

Second Option of the FHD Kernel

```
_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;
for (int m = 0; m < M; m++) {
  float expFhD = 2*PI*(kx[m]*x[n]+ky[m]*y[n]+kz[m]*z[n]);
  float cArq = cos(expFhD);
  float sArg = sin(expFhD);
  rFhD[n] += rMu[m]*cArg - iMu[m]*sArg;
  iFhD[n] += iMu[m]*cArg + rMu[m]*sArg;
```

Using Registers to Reduce Memory Accesses

```
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 (int 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);
  rFhDn_r += rMu[m]*cArg - iMu[m]*sArg;
  iFhDn r += iMu[m]*cArq + rMu[m]*sArq;
rFhD[n] = rFhD r; iFhD[n] = iFhD r;
```

Chunking 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,
                    cudaMemCpyHostToDevice);
  cudaMemcpyToSymbol(ky_c,&ky[i*CHUNK_SIZE],4*CHUNK_SIZE,
                    cudaMemCpyHostToDevice);
  cudaMemcpyToSymbol(kz c,&kz[i*CHUNK SIZE],4*CHUNK SIZE,
                    cudaMemCpyHostToDevice);
  cmpFhD<<<FHD THREADS PER BLOCK, N/FHD THREADS PER BLOCK>>>
     (rPhi, iPhi, phiMag, x, y, z, rMu, iMu, CHUNK_SIZE);
/* Need to call kernel one more time if M is not */
 /* perfect multiple of CHUNK SIZE */
```

Revised FHD Kernel – Constant Memory

```
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*(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]*cArq + rMu[m]*sArq;
 rFhD[n] = rFhD_r; iFhD[n] = iFhD r;
```

Using Hardware __sin() and __cos()

```
_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 (int 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 cArq = cos(expFhD);
  float sArg = __sin(expFhD);
  rFhDn r += rMu[m]*cArq - iMu[m]*sArq;
  iFhDn_r += iMu[m]*cArg + rMu[m]*sArg;
rFhD[n] = rFhD_r; iFhD[n] = iFhD r;
```

Validating Reconstructed Image Using Peak Signal-to-Noise Ratio

A.N. Netravali and B.G. Haskell, Digital Pictures: Representation, Compression, and Standards (2nd Ed), Plenum Press, New York, NY (1995).

$$MSE = \frac{1}{mn} \sum_{i} \sum_{j} (I(i, j) - I_0(i, j))^2 \qquad PSNR = 20 \log_{10} \left(\frac{\max(I_0(i, j))}{\sqrt{MSE}} \right)$$

 I_0 is a known, "perfect" answer. This is typically done by creating k-space samples for a known image, producing a reconstructed image, and compare the two.

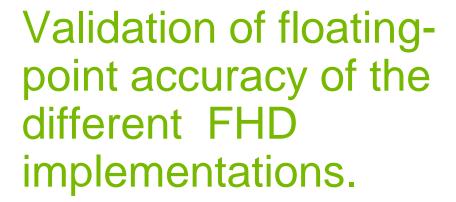
(1) True



(2) Gridded 41.7% error PSNR = 16.8 dB



(3) CPU.DP 12.1% error PSNR = 27.6 dB





(4) CPU.SP 12.0% error PSNR = 27.6 dB

12.1 % error

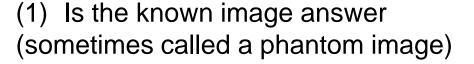
PSNR = 27.6 dB



12.1 % error PSNR = 27.6 dB



(6) GPU.RegAlloc 12.1 % error PSNR = 27.6 dB





(8) GPU.ConstMem (7) GPU.Coalesce 12.1% error PSNR = 27.6 dB

(9) GPU.FastTrig 12.1 % error PSNR = 27.5 dB

Note that all GPU optimized versions have comparable PSNR as (3) the CPU double-precision version

Component and Whole-Application Speedup

| | Q | | F ^H D | | Total | |
|--------------------------------|-----------------|-------------|----------------------|-----------|-------------------------|--------------------|
| Reconstruction | Run Time (m) | GFL OP | Run Time (m) | GFLO P | Linear Solver (m) | Recon. Time (m) |
| Gridding + FFT (CPU, DP) | N/A | N/A | N/A | N/A | N/A | 0.39 |
| LS (CPU, DP) | 4009.0 | 0.3 | 518.0 | 0.4 | 1.59 | 519.59 |
| LS (CPU, SP) | 2678.7 | 0.5 | 342.3 | 0.7 | 1.61 | 343.91 |
| LS (GPU, Naïve) | 260.2 | 5.1 | 41.0 | 5.4 | 1.65 | 42.65 |
| LS (GPU, CMem) | 72.0 | 18.6 | 9.8 | 22.8 | 1.57 | 11.37 |
| LS (GPU, CMem, SFU, Exp) | 7.5 357X |) 178. 9 | 1.5 ₂₂₈ X | 144.5 | 1.69 | 3.19 108X |



GPU Teaching Kit

Accelerated Computing





The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under the Creative Commons Attribution-NonCommercial 4.0 International License.