

















or (m = 0; m	< M; m++) {
rMu[m] = rPh iMu[m] = rPh	i[m]*rD[m] + iPhi[m]*iD[m]; i[m]*iD[m] - iPhi[m]*rD[m];
<pre>for (n = 0; expFhD = 2 cArg = cos sArg = sin</pre>	<pre>n < N; n++) { *PI*(kx[m]*x[n] + ky[m]*y[n] + kz[m]*z[n]); (expFhD); (expFhD);</pre>
rFhD[n] += iFhD[n] += }	rMu[m]*cArg - iMu[m]*sArg; iMu[m]*cArg + rMu[m]*sArg;







Loop Interchange								
 Possible since all iterations of both level are independent 								
<pre>for (m = 0; m < M; m++) { for (n = 0; n < N; n++) { expFhD = 2*PI*(kx[m]*x[n] +</pre>	<pre>for (n = 0; n < N; n++) { for (m = 0; n < N; n++) { expFhD = 2*PI*(kx[m]*x[n] +</pre>							
<pre>cArg = cos(expFhD); sArg = sin(expFhD);</pre>	<pre>cArg = cos(expFhD); sArg = sin(expFhD);</pre>							
rFhD[n] += rMu[m]*cArg - iMu[m]*sArg; iFhD[n] += iMu[m]*cArg + rMu[m]*sArg;	<pre>rFhD[n] += rMu[m]*cArg -</pre>							
(a) before loop interchange (b) after loop interchange Figure 7.9 Loop interchange of the F ^B D computation								
2007–2009 University of Illinois, Urbana–Champaign								



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) {								
<pre>int n = blockIdx.x * FHD_THREADS_PER_BLOCK + threadIdx.x;</pre>								
<pre>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];</pre>								
<pre>for (m = 0; m < M; m++) { float expFhD = 2*PI*(kx[m]*xn_r+ky[m]*yn_r+kz[m]*zn_r);</pre>								
<pre>float cArg = cos(expFhD); float sArg = sin(expFhD);</pre>	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;								
<pre>} rFhD[n] = rFhD r; iFhD[n] = iFhD r;</pre>								
D ^D avid Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009 14 University of Illinois, Urbana-Champaign	UNIVERSITY OF UTAH							



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() {							
<pre>int i; for (i = 0; i < M/CHUNK_SIZE; i++); cudaMemcpyToSymbol(kx_c,&xx[i*CHUNK_SIZE],4*CHUNK_SIZE); cudaMemcpyToSymbol(kz_c,&xy[i*CHUNK_SIZE],4*CHUNK_SIZE); cudaMemcpyToSymbol(kz_c,&xy[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);</n></pre>							
<pre>} /* Need to call kernel one more time if M is not */ /* perfect multiple of CHUNK SIZE */</pre>							
David Kirk/NVIDIA and Wen-mei W. Hwu, 2007-2009 Iversity of Illinois, Urbana-Champaign 16							







global void cmpFHd(float* rPhi, iPhi, phiMag, x, y, z, rMu, iMu, int M) {				
<pre>int n = blockIdx.x * FHD_THREADS_PER_BLOCK + threadIdx.x;</pre>				
<pre>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];</pre>				
<pre>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);</pre>				
<pre>float cArg = cos(expFhD); float sArg = sin(expFhD);</pre>				
<pre>rFhDn_r += rMu[m]*cArg - iMu[m]*sArg; iFhDn_r += iMu[m]*cArg + rMu[m]*sArg;</pre>				
rFhD[n] = rFhD r; iFhD[n] = iFhD r;				
$^{\mathrm{J}}$ 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 20 University of Illinois, Urbana-Champaign 20				









Summary of Results										
	Q		F ^H d		Total					
Reconstruction	Run Time (m)	GFLOP	Run Time (m)	GFLOP	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)	13.6	98.2	2.4	92.2	1.60	4.00				
LS (GPU, CMem, SFU, Exp)	7.5 357X	178.9	1.5 228X	144.5	1.69	3.19 108X				
Javid Kirk (NVDIA and Wen-mei W. Hwu, 2007-2009 University of Illinois, Urbana-Champaign										