



Approaching Projects/STRSM/Case Studies
1. Parallelism?
 How do dependences constrain partitioning strategies?
 Analyze data accesses for different partitioning strategies
 Start with global memory: coalesced?
 Consider reuse: within a thread? Within a block? Across blocks?
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,
6963 L14: Application Case Studies II UNIVERSITY OF UTAH













One Possibility	,
global void cmpFHd(float* rPhi, iPhi, kx, ky, kz, x, y, z, rMu, iMu, int N) {	phiMag,
int m = blockIdx.x * FHD_THREADS_PER_BLO	CK + threadIdx.x;
rMu[m] = rPhi[m]*rD[m] + iPhi[m]*iD[m]; iMu[m] = rPhi[m]*iD[m] - iPhi[m]*rD[m];	
<pre>for (n = 0; n < N; n++) { expFhD = 2*PI*(kx[m]*x[n] + ky[m]*y[n]</pre>	+ kz[m]*z[n]);
cArg = cos(expFhD); sArg = sin(expFhD);
<pre>rFhD[n] += rMu[m]*cArg - iMu[m]*sArg; iFhD[n] += iMu[m]*cArg + rMu[m]*sArg; } }</pre>	This code does not work correctly! The accumulation needs to use atomic operation.
id Kirk/NVIDIA and Wen-mel W. Hwu, 2007-2009 11	



3





<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; m < M; m++) { expFhD = 2*PI*(kx[m]*x[n] + ky[m]*y[n] + kz[m]*z[n]);</pre>				
<pre>cArg = cos(expFhD); sArg = sin(expFhD);</pre>	cArg = cos(expFhD); sArg = sin(expFhD);				
rFhD[n] += rMu[m]*cArg - iMu[m]*sArg; iFhD[n] += iMu[m]*cArg + rMu[m]*sArg;	rFhD[n] += rMu[m]*cArg - iMu[m]*sArg; iFhD[n] += iMu[m]*cArg + rMu[m]*sArg;				
<pre>} } (a) before loop interchange</pre>	<pre>} } (b) after loop interchange</pre>				
Figure 7.9 Loop interchange of the F ^H D computation					
© David Kirk/MVIDIA and Wen-mel W. Hwu, 2007-2009 University of Illinois, Urbana-Champaign	5 UNIVERSITY				



Step 3. Using Registers to Reduce							
Global Memory Trattic							
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						
<pre>rFhDn_r += rMu[m]*cArg - iMu[m]*sArg; iFhDn_r += iMu[m]*cArg + rMu[m]*sArg;</pre>	kx, ky and kz are read- only and based on m						
<pre> / rFhD[n] = rFhD_r; iFhD[n] = iFhD_r; }</pre>							
© David Kirk/NVIDIA and Wen-mel W. Hwu, 2007-2009 17 University of Illinols, Urbana-Champaign	UNIVERSITY OF UTAH						



Tiling k-space data to fit into constant memory





















Summary of Results									
	Q		F ^H d						
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, SEU Exp)	7.5	178.9	1.5	144.5	1.69	3.19			
51 O, Exp)	<u> </u>		228X -			108X —			
David Kirk/NVIDIA and Wen-mei V University of Illinois, Urbana-	V. Hwu, 2007-2009 Champaign				- E	UNIVERSITY OF UTAH			

