GPU Supercomputing – From Blue Waters to Exascale

Wen-mei Hwu

Professor, University of Illinois at Urbana-Champaign (UIUC)

CTO, MulticoreWare Inc.



New BW Configuration



Today's and Tomorrow's Intellectual Challenges

- Challenges that will be faced by the Science Teams include:
 - Scaling applications to large processor counts
 - In the face of limited bandwidth
 - Effective using of many core and accelerator components
 - Using both general purpose and accelerated nodes in single applications
 - Application based resiliency
- NCSA establishing a focused effort in NCSA/UIUC Enhanced Intellectual Services for Petascale Performance (NEIS-P2) to work directly with the Science Teams and also the general community to enable teams to take full advantage of the extraordinary capabilities of Petascale+ scale systems.
 - Due in February 2012
 - GTC Asia, Beijing, 2011

Science Area	Number of Teams	Codes	Structured Grids	Unstructured Grids	Dense Matrix	Sparse Matrix	N- Body	Monte Carlo	FF T	Significan t I/O
Climate and Weather	3	CESM, GCRM, CM1, HOMME	X	X		Х		X		
Plasmas/Magnet osphere	2	H3D(M), OSIRIS, Magtail/UPIC	X				X		X	X
Stellar Atmospheres and Supernovae	2	PPM, MAESTRO, CASTRO, SEDONA	Х			X		X		X
Cosmology	2	Enzo, pGADGET	Х			Х	Х			
Combustion/Tur bulence	1	PSDNS	X						X	
General Relativity	2	Cactus, Harm3D, LazEV	X			Х				
Molecular Dynamics	4	AMBER, Gromacs, NAMD, LAMMPS			Х		Х		X	
Quantum Chemistry	2	SIAL, GAMESS, NWChem			Х	Х	X	х		Х
Material Science	3	NEMOS, OMEN, GW, QMCPACK			Х	Х	X	Х		
Earthquakes/Sei smology	2	AWP-ODC, HERCULES, PLSQR, SPECFEM3D	X	Х			Х			X
Quantum Chromo Dynamics	1	Chroma, MILD, USQCD	X		X	X	Х		X	
Social Networks	1	EPISIMDEMICS								
Evolution	1	Eve								
Computer Science	1			X	X	Х			X	X

Current Science Team GPU Plans and Results

- Nearly 1/3 of PRAC projects have active GPU efforts, including
 - AMBER
 - LAMMPS
 - USQCD/MILC
 - GAMESS
 - NAMD
 - QMCPACK
 - PLSQR/SPECFEM3D
- Others are investigating use of GPUs (e.g., Cactus, PPM, AWP-ODC)
- Some examples follow

There is a critical need for scalable kernel libraries

- Both CPUs and GPUs require scalable parallel kernel libraries
 - GPU needs are more urgent
- Only a small percentage of the Intel Math Kernel Library (MKL) functions have scalable forms.
- Software lasts through many hardware generations and needs to be scalable to be economically viable

Solid Scalable Kernels

- Dense SGEMM/DGEMM, LU, Triangular solvers (CUBLAS, CULA, MAGMA)
- Sparse Matrix Vector Multiplication, Tridiagonal solvers (CUSP, QUDA, PARBOIL)
- FFTs, Convolutions (CUFFT, Parboil)
- N-Body (NAMD/VMD, FMM BU, PARBOIL)
- Histograms (PARBOIL)
- Some PDE solvers (CURRENT, PARBOIL)

Example: Tridiagonal Solver

- Implicit finite difference methods, cubic spline interpolation, preconditioners
- An algorithm to find a solution of Ax = d, where A is an n-by-n tridiagonal matrix and d is an n-element vector



GTC Asia, Beijing, 2011

Thomas Algorithm

Special two-way Gaussian elimination



Forward reduction

Parallel Cyclic Reduction(PCR)

 Simultaneous reduction of odd and even rows – a.k.a. forward reduction only CR

Summary of Previous Approaches

	Complexity	Number of operations	Number of processing steps with n- parallelism machine	
Thomas	O(n)	2n	2n	
PCR	O(n log n)	12 n log n	log n	

Parallelization of Tiled PCR (Owens)

• Exponential growth of halo elements and subsequent reductions



Hierarchical Tiling



Mapping onto GPU							
1 st level tile	One per thread block						
2 nd level tile	Collaborative streaming within thread block						
2 nd level tile buffer	Shared memory						

Tiled PCR using HT Sliding Window



Scalable in size and # systems (8 million elements) M=2048





M=16

M=1



M := number of systems / x-axis := number of unknowns

A Scalable Kernel Requires

- Massive parallelism in application algorithms
 - Data parallelism
- Regular computation and data accesses
 Similar, balanced work for parallel threads
- Avoidance of conflicts in critical resources
 - Off-chip DRAM (Global Memory) bandwidth
 - Conflicting parallel updates to memory locations

Eight Techniques for Scalable Kernels (so far)

Techniques	Memory Bandwidth	Update Contention	Load Balance	Regularity	Efficiency
Scatter to Gather		Х			
Privatization		Х			
Tiling	Х				Х
Coarsening	Х	Х			Х
Data Layout	Х	Х			Х
Input Binning	Х				Х
Regularization			Х	Х	Х
Compaction	Х		Х	Х	Х

GTC Asia, Beijing, 2011

http://courses.engr.illinois.edu/ece598/hk/

Use of Optimizations in Parboil

Benchmark	Unoptimized Im-	Optimizations Applied	Optimized Implementation	Primary Limit of Effi-		
	plementation Bot-		Bottleneck	ciency		
	tleneck					
cutcp	Contention, Local-	Scatter-to-Gather, Binning, Regular-	Instruction Throughput	Reads/Checks of Irrel-		
	ity	ization, Coarsening		evant Bin Data		
mri-q	Poor Locality	Data Layout Transformation, Tiling,	Instruction Throughput	N/A (true bottleneck)		
		Coarsening				
gridding	Contention, Load	Scatter-to-Gather, Binning, Com-	Instruction Throughput	Reads/Checks of Irrel-		
	Imbalance	paction, Regularization, Coarsening		evant Bin Data		
sad	Locality	Tiling, Coarsening	Memory Bandwidth/Latency	Register Capacity		
stencil	Locality	Coarsening, Tiling	Bandwidth	Local Memory, Regis-		
				ter Capacity		
tpacf	Locality,	Tiling, Privatization, Coarsening	Instruction Throughput	N/A (true bottleneck)		
	Contention					
lbm	Bandwidth	Data Layout Transformation	Bandwidth	N/A (true bottleneck)		
dmm	Bandwidth	Coarsening, Tiling	Instruction Throughput	N/A (true bottleneck)		
spmv	Bandwidth	Data Layout Transformation	Bandwidth	N/A (true bottleneck)		
bfs	Contention, Load	Privatization, Compaction, Regular-	Bandwidth	Whole-Device Local		
	Imbalance	ization		Memory Capacity		
histogram	Contention, Band-	Privatization, Scatter-to-Gather	Bandwidth	Reads of Irrelevant		
	width			Input (alleviated by		
				cache)		

How a mathematician writes matrix multiplication

$$(MN)_{j,i} = \sum_{k} M_{j,k} N_{k,i}$$

GTC Asia, Beijing, 2011

How a smart CUDA programmer writes matrix multiplication

```
float c[TILE_N];
 for (int i=0; i < TILE_N; i++) c[i] = 0.0f;
 int mid = threadIdx.y * blockDim.x + threadIdx.x;
 int m = blockIdx.x * TILE M + mid;
 int n = blockIdx.y * TILE N + threadIdx.x;
 __shared__ float b_s[TILE_TB_HEIGHT][TILE_N];
 for (int i = 0; i < k; i+=TILE_TB_HEIGHT) {
  float a;
  b_s[threadIdx.y][threadIdx.x]=B[n + (i+threadIdx.y)*Idb];
  syncthreads();
  for (int j = 0; j < TILE_TB_HEIGHT; j++) {
   a = A[m + (i+i)*Ida];
   for (int kk = 0; kk < TILE N; kk++) c[kk] += a * b s[i][kk];
    _syncthreads();
 int t = ldc*blockldx.y * TILE_N + m;
 for (int i = 0; i < TILE N; i++) {
  C[t+i*ldc] = C[t+i*ldc] * beta + alpha * c[i];
dim3 grid( m/TILE M, n/TILE N), threads( TILE N, TILE TB HEIGHT);
```

mysgemmNT<<<grid, threads>>>(A, Ida, B, Idb, C, Idc, k, alpha, beta);

...

Writing efficient code is complicated. Tools can provide focused help or broad help

Planning how to execute an algorithm Implementing the plan



UIUC/MCW Tools for Heterogeneous Parallel Programming

Early

Writing optimizable, portable, kernels **Pyon** Task and data movement & multi-GPU copy support GMAC, TM Performance tools **PPA/ADAPT** TC/SM **DL** (Data Layout)

Higher-level Interfaces for Programmability, Portability, and Performance

Performance portability, Analysis, optimization (with collaboration with DSLs – Hanrahan/Keutzer)

Late

Use of UIUC/MCW Tools in Blue Waters

- Introduce compiler and library capabilities into the science team workflow to significantly reduce the programming effort and impact on code maintainabilitys:
 - Compiler based directives
 - GMAC a library that provides global shared memory and automates data transfer/coherence between the CPUs and the GPUs in a node
 - DL is a compiler-based memory layout transformation tool that uses a combination of compiler and runtime support to ease the task of adjusting memory layout to satisfy conflicting needs between the CPU and the GPU
 - TC is a compiler based tool for thread coarsening and data tiling.
- Provide expert support to the science teams through hand-on workshops, courses, and individualized collaboration programs.
 GTC Asia, Beijing, 2011

Example - DL (Data Layout)

- DRAM bursts are formed differently in a heterogeneous system
 - From last level cache misses on CPUs
 - From SIMD-ized memory accesses on many-core architectures like GPUs
- Data layout transformation can mitigate the gap
 - E.g.: Array-of-structure / Discrete-arrays
 - Bridging divergent layout requirements between CPU cores and GPU cores
 - Transparent and efficient marshaling

Data Layout Alternatives



GTC Asia, Bei Array of Structure of Tiled Array (ASTA) [z][y_{31:4}][x_{31:4}][e][y_{3:0}][x_{3:0}]

DL for OpenCL



UIUC/MCW solution ASTA

 Array-of-Structure-of-Tiled-Arrays: preserving locality while gaining coalesced memory access
 – A[x].foo → A[x/4].foo[x%4] for ASTA(4)



GTC Asia, Beijing, 2011



•	< ar2012:	INPAR Parb	inpar_chan	😨 St. Lucia Va	😨 To Rent a C	Reservation	Pending Inv	🔷 Create Your	Create Your	I IMPACT
	multicoreware	inc.com/index.php	option=com_con?	ent&view=article&	uid=26<emid=30				습 -	୯ 🚼 - ଜ
		MUL	TICO	RE	Company	Technology	Products	News C	Contact Us	Q SEARCH
			V WA	RE						

multicorewareinc.com

Download button at bottom

MulticoreWare Tools

Increasing Programmer Productivity and Performance Portability across Compute Platforms

The MulticoreWare tools framework drastically reduces the number of code versions and the time spent in Performance Optimization efforts, by providing key performance tracking and analysis tools coupled with micro-architecture optimized compile-time libraries. As a result, there is no need to have hand-optimized versions for each platform.



Login

Email
Password
Remember me
LOGIN

Forgot login? Register

Conclusion and Outlook

- Petascale and Exascale intellectual challenges
 - Scaling to large processor count with limited interconnect bandwidth
 - Effective use of massively parallel throughput oriented processors
- There is a critical need for scalable kernels
 - Algorithm design for scalable kernel libraries
 - Seamless use of kernels from major languages
 - Productivity tools for kernel development and deployment

THANK YOU!

Both Fusion and Discrete GPU Markets are Growing

