GPU Computing

Axel Koehler Sr. Solution Architect HPC



NVIDIA: Parallel Computing Company



ARM SoCs: Tegra







Continued Demand for Ever Faster Supercomputers



Comprehensive Earth System Model at 1KM scale, enabling modeling of cloud convection and ocean eddies. First-principles simulation of combustion for new high-efficiency, lowemision engines.





Coupled simulation of entire cells at molecular, genetic, chemical and biological levels. Predictive calculations for thermonuclear and core-collapse supernovae, allowing confirmation of theoretical models.



Power Crisis in Supercomputing



Multi-core CPUs



- Industry has gone multi-core as a first response to power issues
 - Performance through parallelism, not frequency
- But CPUs are fundamentally designed for single thread performance rather than energy efficiency
 - Fast clock rates with deep pipelines
 - Data and instruction caches optimized for latency
 - Superscalar issue with out-of-order execution
 - Lots of predictions and speculative execution
 - Lots of instruction overhead per operation

Less than 2% of chip power today goes to flops.



Accelerated Computing

Add GPUs: Accelerate Applications

CPUs: designed to run a few tasks quickly.



GPUs: designed to run many tasks *efficiently*.

Energy efficient GPU Performance = Throughput

Fixed function hardware

- Transistors are primarily devoted to data processing
- Less leaky cache

SIMT thread execution

Groups of threads formed into warps which always executing same instruction

Cooperative sharing of units with SIMT

- eg. fetch instruction on behalf of several threads or read memory location and broadcast to several registers
- Lack of speculation reduces overhead
- **Minimal Overhead**
 - Hardware managed parallel thread execution and handling of divergence



Supercomputing Weather / Climate Modeling Molecular Dynamics Computational Physics		
Life Sciences Biochemistry Bioinformatics Material Science	Tesla K20/K20X	
Manufacturing Structural Mechanics Comp Fluid Dynamics (CFD) Electromagnetics	Kepler GK110	
Defense / Govt Signal Processing Image Processing Video Analytics	Tesla K10	
Oil and Gas Reverse Time Migration Kirchoff Time Migration	Kepler GK104	

Product Name	К10	K20	K20X		
GPU Architecture	Kepler: GK104	GK110	GK110		
# of GPUs	2	1	1		
Peak Single Flops Peak SGEMM	4.58 TF (2.3 TF per GPU) 2.98 TF	3.52 TF 2.61 TF	3.95 TF 2.90 TF		
Peak Double Flops Peak DGEMM	0.19 TF (0.095 TF per GPU) 0.12 TF	1.17 TF 1.10 TF	1.32 TF 1.22 TF		
Memory size	8 GB (4GB per GPU)	5 GB	6 GB		
Memory BW (ECC off)	320 GB/s (160GB/s per GPU)	208 GB/s	250 GB/s		
New CUDA Features	GPUDirect w/ RDMA	GPUDirect (RDMA), Hyper-Q, Dynamic Parallelism			
ECC Features	External DRAMs only	DRAM, Caches & Reg Files			
# CUDA Cores	3072 (1536 per GPU)	2496	2688		
Total Board Power	225W	225W	235W		
Board Type	PCI-e Passive	PCI-e Passive, Active, SXM	PCI-e Passive SXM		

Kepler GK110 Block Diagram

- 7.1B Transistors
- 15 SMX units
- 1.3 TFLOP FP64
- 1.5 MB L2 Cache
- 384-bit GDDR5
- PCI Express Gen3 compliant



				PCI Express 3	0 Hoet Niterlace						
				- Gga Din	ad Engine						
	SMX		504X	EMX	SMX	BMX	SMX	EAXX			
Munory Controller M									Munnery Controller M		
Annuary Controller;											
Memory Controller	SMX-								Mannory Controller		

Kepler GK110 SMX vs Fermi SM



3x sustained perf/W

Ground up redesign for perf/W 6x the SP FP units 4x the DP FP units Slower FU clocks

~4x the overall instruction throughput 2x register file size (64K regs)
2x threadblocks (16) & 1.33x threads (2K)

SMX			_	_						-				_				_
	ry: Not	entraint .		R.	w	arge Dischard	-		-	w	ingi Dath	with the			W	igi Dehee		-
The second second	.a 11	Diseasch	ile.	COnversion Contract (Conversion Contract)			C Desired to Date of Contracts Date of					Disease Line Consults lives				-		
							egister	File (65,536	* 32-4	xH3	-			-			
		-				+			+	+	-	100				+	-	
Care Gran	9-1	DEDUC	600	Gare	(inter	TP see	Links	SFU	<u>Gene</u>	Core	Gunn	Diff. Later	Gree	Care	Otre	par terra	LOW	SPU
Sine Com	e i	107 (Date	Gine	Gum	Gum	III Ind	LUST	aru.	Ginie	Gan	Gum	1-100	Gare	cine	Gen	207 1040	LOW	1 90
Game Gone	C:m	BELINE	Gare	Gore	Cure	Diff Local	1.0-11	seu	See	Gom	Guis	C. Unit	Core	Gam	Cen	E	LOWT	884
Care Care	-		-	Core	dam	CP Loss	LUST	siru	Sme	G arry	dun	Diff. Lines	6-m	Gare	C	ter tite	LOWF	39U
Cure Coine	C in t	ar tint	1011	Garr	Cure	-	LDIST	SFU	6.6	Core	Circ	P*100	Con	Cân	Circ	an Line	10 H	BFU
Cone Cone	Cate	BP-More	Con	Com	Chre	-	LIST	sru	See	Con	Core	Der sone	Core	Cam	Core	and the second	LDWT	sru
Care Com	dim	Dir Luis	Cons	Chie	Com		(nat	SFU	Corre	60m	Core	Diff Line	C	CÓN	Circ	and the second	win	seu
Care Care	8-1	DP.Det	Gerr	Gon	Gum		LDBT	sru	Goog	Con	Gami	D# HIM	G 100	Gam	Curs	an lass	LOST	5751
Case Care	d and	DP Unit	Cite	Gire	Gum	TT Line	UNT	1 70	tini.	Core	Guin	Dir Line	Cire	Cim	Circ	Der Alma	LUNT	1P U
Gun Gun	Q.ext	111 1.000	Core	Cine	Core	-	1011	3FU	Con	C (m)	Core	Der Unit	Com	Cons	Corr	127 LINE	LOWT	804
Cure Cure	6.m	BP LINE	Core	Gim	C um	-	LOST	SFU	84 s	Core	Q uite	DP-UM	Core	Core	Core	Der Land	LOST	SRU
Gitte Guite	684	ar linn	Gain	Gire	Gire	-	i dest	SFU	Gon	Gyre	Gore	DP-Unit	Gire	C(m	Gire	Der Lind	LDAT	SFU
Care Com	-		Cine	Gam	Gime	-	ATEST	sru	See	Core	Gune	Diff. Laws	Con	Chris	Con	Lan Line	LOW	srú
Care Com	Gara	DP 100	Coin	60m	Com		1057	SFU	Con	Com	Core	-	Core	6ĝin	Cere.		Linet	88U
Care Care	t a	million	<u>Gun</u>	Gare	ψų.	CP Line	LDET	SFU	G 111	B are	Gun		Case .	Care	<u>Corr</u>	Er land	LOST	seu
Cite Care	č.	Der Losit	Č. H	Gare	Guns		trat	8FU	Con	Core	Eore	Dir Line	Cim	Gim	6 int	DertAines	Quer	BFU
	Character Methods																	
							40.00	a News	-					_				
Tex		Tex	4		Тех	8	Tex	E.	1	Тех	J	Tex	6 - 1		Тех	1	Tex	8
Tex		Tex			Tex		Tex	6		Тех		Tex		Tex			Tex	

Hyper-Q



Proxy - A Multi-Process Runtime for MPI



Why

Speedups for MPI programs with low-GPU utilization

How

- multiple CPU processes on a single GPU simultaneously
- client-server architecture
- client processes share the same CUDA context

When

Currently on Cray; Production on Linux with CUDA 5.5

What is Dynamic Parallelism?

The ability to launch new kernels from the GPU

- Dynamically based on run-time data
- Simultaneously from multiple threads at once
- Independently each thread can launch a different grid



Fermi: Only CPU can generate GPU work

Kepler: GPU can generate work for itself

Dynamic Parallelism

Simpler Code, More General, Higher Performance





Dynamic Parallelism



Familiar Syntax and Programming Model



Simpler Code: LU Example

LU decomposition (Fermi) dgetrf(N, N) { for j=1 to N for i=1 to 64idamax<<<>>>> idamax(); memcpy dswap(); dswap<<<>>>> memcpy dscal(); dscal<<<>>>> dger<<<>>>> dger(); next i memcpy dlaswap(); dlaswap<<<>>> dtrsm<<<>>>> dtrsm(); dgemm<<<>>>> dgemm(); next i **CPU Code** GPU Code

LU decomposition (Kepler)



GPU Callable Libraries Direct Access to BLAS and other Libraries from GPU Code



Source files compiled separately to create independent object files Linker creates GPU Callable Libraries and links with CUDA code

NVIDIA® GPUDirect™ Support for RDMA Direct Communication Between GPUs and PCIe devices



Server 1

Server 2

GPU-aware MPI

- Support GPU to GPU communication through standard MPI interfaces without exposing low level details to the programmer (make MPI implementations aware of GPU pointers)
 - e.g. enable MPI_Send, MPI_Recv from/to GPU memory
 - Made possible by Unified Virtual Addressing (UVA) in CUDA 4.0
 - MVAPICH2, OpenMPI, Platform MPI

Code without MPI integration

At Sender:

cudaMemcpy(s_buf, s_device, size, cudaMemcpyDeviceToHost); MPI_Send(s_buf, size, MPI_CHAR, 1, 1, MPI_COMM_WORLD);

At Receiver:

MPI_Recv(r_buf, size, MPI_CHAR, 0, 1, MPI_COMM_WORLD, &req); cudaMemcpy(r_device, r_buf, size, cudaMemcpyHostToDevice);

Code with MPI integration

At Sender:

MPI_Send(s_device, size, ...); At Receiver: MPI_Recv(r_device, size, ...);

CUDA Compiler Contributed to Open Source LLVM

Developers want to build front-ends for Java, Python, R, DSLs

Target other processors like ARM, FPGA, GPUs, x86





Minimum Change, Big Speed-up

Application Code



Ways to Accelerate Applications



MATLAB Parallel Computing On-ramp to GPU Computing

Most popular math functions on GPUs

 \bullet

- Random number generation
- FFT
- Matrix

- multiplications
- Solvers
- Convolutions
- Min/max

- SVD
 - Cholesky and LU factorization



GPU features in Communications Systems Toolbox





GPU Accelerated Libraries

"Drop-in" Acceleration for Your Applications



OpenACC Directives



Your original Fortran or C code

Easy, Open, Powerful

- Simple Compiler hints
- Works on multicore CPUs & many core GPUs
- Compiler Parallelizes code
- Future Integration into OpenMP standard planned

http://www.openacc.org



OpenACC Directives Example !\$acc data copy (A, Anew) Copy arrays into GPU memory within data region

iter = iter +1

err=0. fp kind

iter=0

!\$acc kernels

```
do j=1,m
   do i=1,n
    Anew(i,j) = .25 fp kind *( A(i+1,j) + A(i-1,j) \&
                             +A(i, j-1) + A(i, j+1))
    err = max(err, Anew(i,j)-A(i,j))
   end do
  end do
!$acc end kernels
```

do while (err > tol .and. iter < iter max)

```
IF(mod(iter,100)==0 .or. iter == 1) print *, iter, err
A= Anew
```

end do

!\$acc end data

Parallelize code inside region

Close off parallel region

Wide Adoption of Tesla GPUs



Over 200 GPU-Accelerated Applications

POPULAR OPU-ACCELERATED APPLICATIONS

writed Features

Hulti-SPU Release State

http://www.nvidia.com/teslaapps/

Malauriar Dynamic		POPULAR GR	PU-ACCELERATED APPLICA	TIONS, Continue	d		55							
Abatime	Models molecular dynamics of tropol simulations of proteins, DNA and liga	Application	Description	Suppo	ted Fostures Expected H	lulti-GPU Release Sta	tun							
ACEMD	Simulation of mechanics force Selds, & explicit setwart on CUDA	Weather & Elimite	Forecasting	POPULAR GP	U-ACCELERATED APPLICA	TIONS, Cantinue	đ	Service Company						
AMBER	Suite of programs to simulate molecu dynamics on biumolecules	ASUCA	Weather forecasting model hilly optimized for SPLIs	Application	Description	Support	rted Features Expected M	ulti-GPU Release						
DL-POLY	Simulate macromolecules, polymers, systems, etc on a distributed memory earolist convector	CAM/SE	Community Atmospheric Model is a g atmosphere model for weather and c research	R + 3 And the set of the										
6804425	55 Simulation of biochamical molecules 00		Weather modeling and forecasting ap used by NASA	Technologies ADS	and high speed digital circuits	- Apple and	(111)25.44(J)	DODIN AR CO						
WOMO Date	Darticle donarran our king eritten in	HIRLAM	Weather forecasting model fully optim for GPUs	Technologies	analyzing 30 EM effects of high speed RF/Microwave components	Accelonary SITM	Seismic Processing	Application	Description	Supported Features	Expected	Matti-GPU	Release Statu	
1 AMAIDS	up for GPUs	HOMME	Weather mediating tool for atmospher scientists	ANSYS Nacion	Circuit sehulation engine for RE/anals mixed-signal IC design; IBIS-AMI ana	CGGVentas RTH RASHLENE	Seamic Processing Seamic interpretation	11410	125071.00077		Speed Up*	Sepport		
NAMO	Designed for high-performance simular filame melacular systems	INCOM .	Weather forecasting model using icos horsontal gnd	CST Microwave	speedup with GPU computing High frequency electromagnetic	Headerano Suite	Seismic Imaging	Chroma	General purpose LOCD application	Wilson-clover fermions, Krylov	- 5-åx	19s	Available now	
Quantum Chemistry	hadred a databases and a family of	Hilgen	Numerical model designed for study- atmosphere, acean, and climate	Studio (MWS) Gauda OPC, OPV	Reld simulation Collection of several software tools fo	Paradigm EarthStudy260	Reservoir Moduling	MILC	General purpose LOCO application	Staggered fermions, Krytev solvers, Gauge Ltik futtening	5-dx	Yes	Available now	
GAMESS-US	Computational chemistry suite used to simulate atomic and molecular electric	NIM	Weather forecasting model using icos horizontal grid	Contestion - 2 million	Computational lithography running of Bauda hardware platform	Paradigm Eches RTM	Seamic Proceesing	Computational Flat	d Dynamics					
	structure	WHE	Weather and Dozon modeling application	Harmonia Ar-dis	Julie EM modeling and simulation	Paradigm SKUA	Reservoir Modeling	Altan AcuSelve	General purpose CFD flow solver	Linear equation solver	2x	Tes	Available now	
MAChem	Computational chemistry package de for HPC clusters	Telting and Effects		RECARGE MCAD	30 EM modeling and smulation	Paradigm	Seismic Interpretation	Autodesk Moldflow	Optimize design of plastic parts and injection molds	Linear equation solver	2x	Single Only	Available now	
D-CHEM	Computational chemistry package de	Adabe Pranters Pro	Video editing	CAD		Vaselien Schlunderger WesterGern	Seismic Processing	FEFLO IGNU-Lateori	Navier-Stokes flow solver based on unstructured grids for modeling both compressible and incontressible flows	Explicit solver	10x	Ves	br: Devolopment	
TeraChem	for HPC clusters Quantum chemistry software designe	Avid Media Composer	Video editing	CATIA V& Line Rendering	Photorealistic rendering	Groega2 RTM Seismic City	Seemic Processing	FluiDyna LButtra	Computing physical flows in and around solid bodies	LBM, particle CFD	20x	Yes	Available now	
	19 run et MVDIA GPU	GenArts Sapphere Sony Vecas Pro	Effects plug-in for odes editing Video editing	Bunkspeed Pro Suite	Easy to use photorealistic rendering s	Prestack Interpretation		FluiDyna Culsas- OpenF0AM	Computing physical flaws with Guises — a software library with special algorithms for	Linear equation solvers	3i Sober	Single Only	Available now	
Materiale Science	Materials code for investigation the of	Animation	Research Control of	HTT DeltaGen 10.x	Photorsatiotic rendering used for des	Spectra5ets	Seismic Processing / Imaging	Pronotech Participantes	Fluid simulation for free surface flow like Tsuriami, material processing and liquids	MPS, Particle CFD	48.98	Yes	Analiable now	
ONCIACK	of temperature on magnetism Solves the many-body Schundioger at	Autodesk 30s Mar	3D modeling, animation, and rendering			Sinneridge Reservair	Reservoir Sknutation	530 Sandia NL SIGI	Masswely parallel direct numerical selver (DNS) for the full compressible Narier-Stokes	Chamistry kernel	Rx 5P, 5x DP karnel	Yes	In Developmen	
entres entres s.	for electronic structures using a guan	Autodesk Maya	30 moteling, animation, and	HTT Demand	TeamCenter and RTT formats	Simulation		Turbestream	Uttrafast CFD solver for turbomachines	Explicit salver	19x.	Tes	Available now	
Quantum-	An integrated suits of computer coder distances characterized and computer coder	and the second sec	cheese & Table Organize			Summe all H	Seismic Processing	Vratis SpeedIT- OpenFOAM Solver	Set of accelerated solvers for sparse linear systems of equations	Linear equation solvers	3x Solver	Yes	Araitable now	
ten Material Material	modeling at the nanoscale	Deferre & Willing			GPU acceleration for MATLAB	Littered a	Real time entities analytical engine b	Computational Dis	ctural Hechasles					
WASP	First principles materials code that or electronic structures and quantum-	Advanced Onthe Serves	Contraction of Contract and Con-	Mathematica	Symbolic math analysis	Association MATLAB	Data parallel mathematics (MATLAB	AbaşınyStanlari	Simulation and analysis tool for structural mechanics	Linear equation solver	15-25#	Single Only	Available now	
Concernance and	mechanical molecular synamics	Eterrix Blaze Terra	Geospatial Visualization	Woldram	Technical computing language and	Mathworks Manax	PCT, MDCSI Risk analytics (MACSI	AND/S Hechanical	Simulation and analysis tool for structural mechanics	Linear equation solver	21	Sogle Only	Available now	
Amira S	A multifaceted software platform for	Exalls (117) ENVI	Geospatial Visualization	Mathemata	Integrated development environment (MATLAB PCT, NDCS)	Numerical Algorithms Group	Random Number Generators	Impetus Alva	Predicts large deformations of structures, and components exposed to extreme loading	Linear equation colver, SPH	10x SPH, 2x Total	Yes.	Asiaitable now	
	life sciences and bio-modical data	GenEye Analytics Signature Analysi	Geosparial Visualization			SciComp, Inc	Derivative pricing ISciFinancel	15-09NA Implicit	Kubichesis simulation tarkate used	Linear equation solver	31	Vis	In Development	
Core Hopping	Rapid screening at novel cores to imp drug properties	Geo/Web3d Desktop	Geospatial Visualization			Welfram Mathematica	Mathematical Development Environm	MSC Nastran	Simulation and analysis tool for structural mechanics	Linkar equation solver	1.4-21	Yes	Available new	
FaitROCS VMD	30 molecular shape comparison Vexations and analyzing large bio-m	Incogna GIS	Geospatual Visualization			*SPU performance of terms to kenter and	emparted against multi-come x86 CPU secket. I Isomande combanteur, Partismande resulte a	Hart	Simulation and analysis tool for structural mechanics	Linear equation solver	1.5x	Res	In Development	
	systems in 3-D graphics	Intergraph Motion Video Analyst	Video fitters and mosaic ing Geo-Is FMV analytics with intelligence data Video Exolution					RADIOSS Implica	Used to maximize durability, NVH, crash, talkity, manufacturability and fluid-structure interaction performance	Linear equation solver	24	Single Only	In Developmen	
		Paneptes 1.0	THE AUTOM											
		the second se	the second se											

Titan: World's Fastest Supercomputer

18,688 Tesla K20X Accelerators

27 Petaflops Peak: 90% of Performance from GPUs

17.59 Petaflops Sustained Performance on Linpack

Tesla CUDA Architecture Roadmap



DP GFLOPS per Watt

How Will GPUs Evolve Over This Decade?

- Integration (memory, processor types, network)
- Further concentration on locality (both HW and SW)
- Reducing overheads (intra- and inter-GPU)
- Continued convergence with consumer technology









Echelon

NVIDIA's Extreme-Scale Computing Project DARPA UHPC Program











Fast Forward Program



Which Takes More Energy?

Performing a 64-bit floating-point FMA: 893,500.288914668 × 43.90230564772498 = 39,226,722.78026233027699 + 2.02789331400154 = 39,226,724.80815564 Or moving the three 64-bit operands 20

mm across the die:



This one takes over 4.7x the energy today (40nm)! It's getting worse: in10nm, relative cost will be 17x! Loading the data from off chip takes >> 100x the energy.

Summary

- Power is the main HPC constraint
 - Vast majority of work *must* be done by cores designed for efficiency
- NVIDIA GPU's are already designed for energy efficiency
- Data movement dominates the power
 - Locality at all levels and reduction of overhead is necessary
- GPU computing has a sustainable model
 - Aligned with technology trends, supported by consumer markets
- GPUs are the path to the tightly-coupled hybrid processor future



Thank you. Questions?

Axel Koehler Sr. Solution Architect HPC akoehler@nvidia.com

