

Challenges of Energy Efficient Scientific Computing

John Shalf

National Energy Research Supercomputing Center Lawrence Berkeley National Laboratory

1st Symposium on Energy Efficient Electronics

Berkeley, June 12, 2009









Part I

Power Crisis in HPC







New Design Constraint: POWER

- Transistors still getting smaller
 - Moore's Law is alive and well
- But Dennard scaling is dead!
 - No power efficiency improvements with smaller transistors
 - No clock frequency scaling with smaller transistors
 - All "magical improvement of silicon goodness" has ended
- Cannot continue with business as usual
 - DARPA study extrapolated current design trends and found brick wall at end of exponential curves
 - Can only accelerate existing research prototypes (not "magic" new disruptive technology)!





ORNL Computing Power and Cooling 2006 - 2011



Computer Center Power Projections

- Immediate need to add 8 MW to prepare for 2007 installs of new systems
- NLCF petascale system could require an additional 10 MW by 2008
- Need total of 40-50 MW for projected systems by 2011
- Numbers just for computers: add 75% for cooling
- Cooling will require 12,000 15,000 tons of chiller capacity

Cost estimates based on \$0.05 kW/hr

sources used at each DOE site. Information is entered into EMS4 by the site and

Annual Average Electrical Power Rates \$//////									
Site	FY 2005	FY 2006	FY 2007	FY 2008	FY 2009	FY 2010			
BNL	43.70	50.23	53.43	57.51	58.20	56.40 *			
ANL .	44.92	53.01		<i>.</i> –					
ORNL	46.34	51.33	Data taken from Energy Management System-4 (EMS4). EMS4 is the DOE corporate system for collecting energy information from the sites. EMS4 is a web-based						
PNNL	49.82	N/A	system t	hat collects ener	rgy consumption	and cost informa	ation for all energy		

Puel Average Fleetrical Dever Dates (1/1/h

reviewed at Headquarters for accuracy.

OAK RIDGE NATIONAL LABORATORY U. S. DEPARTMENT OF ENERGY





The New York Times

🖸 TalkBack 🗹 E-mail 🚔 Print

"Hiding in Plain Sight, Google Seeks More Power", by John Markoff, June 14, 2006



New Google Plant in The Dulles, Oregon, from NYT, June 14, 2006

Relocate to Iceland?





HPC Power: It will only get worse

- Recent Baltimore Sun Article on NSA system in Maryland
 - Consuming 75MW and growing up to 15MW/year
 - Not enough power left for city of Baltimore!
- LBNL IJHPCA Study for ~1/5 Exaflop for Climate Science in 2008
 - Extrapolation of Blue Gene and AMD design trends
 - Estimate: 20 MW for BG and 179 MW for AMD
- DOE E3 Report

NATIONAL ENERGY RESEARCH SCIENTIFIC COMPUTING CENTER

- Extrapolation of existing design trends to exascale in 2016
- Estimate: 130 MW
- DARPA Study
 - More detailed assessment of component technologies
 - Estimate: 20 MW just for memory alone, 60 MW aggregate extrapolated from current design trends



The current approach is not sustainable!





DARPA Exascale Study

- Commissioned by DARPA to explore the challenges for Exaflop computing
- Two model for future performance growth
 - Simplistic: ITRS roadmap; power for memory grows linear with #of chips; power for interconnect stays constant
 - Fully scaled: same as simplistic, but memory and router power grow with peak flops per chip







We won't reach Exaflops with the current approach









... and the power costs will still be staggering



From Peter Kogge, DARPA Exascale Study







Primary Design Constraint: POWER

- Power Efficiency and clock rates no longer improving at historical rates
- Demand for supercomputing capability is accelerating!
- DOE is targeting an exaflop system for 2016



- Exascale (10¹⁸ FLOP/s) cannot be built by simply scaling petascale systems
- -Power requirements for incremental approach are *profoundly* impractical
- -And we have finite \$'s for development costs







The Challenge

Where do we get a 1000x improvement in performance with only a 10x increase in power?

How do you achieve this in 10 years with a finite development budget?







- **1. Processor**
- **2. Interconnect**
- 3. Memory
- 4. Software tuning (auto-tuning)
- **5. Algorithms**
- 6. Power/Cooling/facilities (ask Bill & Dale)







Hardware: What are the problems?

(Lessons from the Berkeley View)

- Current Hardware/Lithography Constraints
 - Power limits leading edge chip designs
 - Intel Tejas Pentium 4 cancelled due to power issues
 - Yield on leading edge processes dropping dramatically
 - IBM quotes yields of 10 20% on 8-processor Cell
 - Design/validation leading edge chip is becoming unmanageable
 - Verification teams > design teams on leading edge processors
- Solution: Small Is Beautiful
 - Simpler (5- to 9-stage pipelined) CPU cores
 - Small cores not much slower than large cores
 - Parallel is energy efficient path to performance:CV²F
 - Lower threshold and supply voltages lowers energy per op
 - Redundant processors can improve chip yield
 - Cisco Metro 188 CPUs + 4 spares; Sun Niagara sells 6 or 8 CPUs
 - Small, regular processing elements easier to verify





Low-Power Design Principles



 Cubic power improvement with lower clock rate due to V²F

- Slower clock rates enable use of simpler cores
- Simpler cores use less area (lower leakage) and reduce cost

Tailor design to application to REDUCE WASTE

This is how iPhones and MP3 players are designed to maximize battery life





ERSC Low-Power Design Principles

Tensilica XTensa Intel Atom Intel Core2 Power 5

- Power5 (server)
 - 120W@1900MHz
 - Baseline
- Intel Core2 sc (laptop) :
 - 15W@1000MHz
 - 4x more FLOPs/watt than baseline
- Intel Atom (handhelds)
 - 0.625W@800MHz
 - 80x more
- Tensilica XTensa DP (Moto Razor) :
 - 0.09W@600MHz
 - 400x more (80x-120x sustained)



NATIONAL ENERGY RESEARCH SCIENTIFIC COMPUTING CENTER





Low Power Design Principles



- Power5 (server)
 - 120W@1900MHz
 - Baseline
- Intel Core2 sc (laptop) :
 - 15W@1000MHz
 - 4x more FLOPs/watt than baseline
- Intel Atom (handhelds)
 - 0.625W@800MHz
 - 80x more
- Tensilica XTensa DP (Moto Razor) :
 - 0.09W@600MHz
 - 400x more (80x-100x sustained)

Even if each simple core is 1/4th as computationally efficient as complex core, you can fit hundreds of them on a single chip and still be 100x more coverignment of energy



Future HPC Technology Building Blocks

Previous Decade

- Optimization target: minimize price to buy more hardware
- COTS: Redirect off-the-shelf components designed for mass market
- This leveraged "Moore's Law" density improvements

Next Decade

- Optimization target: minimize power consumed for work performed
- Specialize and integrate: Embedded + SoC is proven design point
- This leverages "Bells Law" cost efficiency: Commodity not COTS







Future HPC Technology Building Blocks

Previous Decade

- Optimization target: minimize price to buy more hardware
- COTS: Redirect off-the-shelf components designed for mass market
- This leveraged "Moore's Law" density improvements

Next Decade

- Optimization target: minimize power consumed for work performed
- Specialize and integrate: Embedded + SoC is proven design point
- This leverages "Bells Law" cost efficiency: Commodity not COTS

Interim solution: Accelerators

- Demonstrate huge efficiency potential of manycore
- Demonstrate we have failed to learn from CM5 (PCIe)
- Stepping stone to convergence (merge manycore with host memory)
- But also points to benefits of some specialization







Conclusion

- Future HPC must move to simpler power-efficient core designs
 - Embedded/consumer electronics technology is central to the future of HPC
 - Convergence inevitable because it optimizes both cost and power efficiency







Interconnects







Interconnects: Leading Issues

- Cannot continue to scale fully-connected interconnect topologies
- Cannot continue to scale bandwidth using electrical networks

What technology be applied to address these constraints?







The problem with Wires:

Energy to move data proportional to distance

• Wire cost to move a bit:

- energy = bitrate * Length² / cross-section area
- On-Chip (1cm): ~1pJ/bit, 100Tb/s
- On-Module (5cm): ~2-5pJ/bit, 10Tb/s
- On-Board (20cm): ~10pJ/bit, 1Tb/s
- Intra-rack (1m): ~10-15pJ/bit, 1Tb/s
- Inter-cabinet(2-50m): 15-30pJ/bit, 5-10Tb/s aggregate
- To move a bit with optics: target ~1-2pJ/bit for all distance scales(but initial cost high)

Photonics requires no redrive and passive switch little power





Copper requires to signal amplification even for on-chip connections







Interconnect Cost

(Scalable Topologies)

BERKELEY

- Fully-connected networks scale superlinearly in cost, but perform the best
- Limited-connectivity networks scale linearly in cost, but introduce new problems





Interconnect Design Considerations for Message Passing Applications

- Application studies provide insight to requirements for Interconnects (both on-chip and off-chip)
 - On-chip interconnect is 2D planar (crossbar won't scale!)
 - Sparse connectivity for most apps.; crossbar is overkill
 - No single best topology

NATIONAL ENERGY RESEARCH SCIENTIFIC COMPUTING CENTER

- Most point-to-point message exhibit sparse topology + often bandwidth bound
- Collectives tiny and primarily latency bound
- Ultimately, need to be aware of the on-chip interconnect topology in addition to the off-chip topology
 - Adaptive topology interconnects (HFAST)
 - Intelligent task migration?





Using Optical Circuit Switches to Make Fat-Trees into Fit-Trees



ERS

NATIONAL ENERGY RESEARCH

A 2-ary 4-tree with 16 nodes.



A (2, 2, 4)-TL fit-tree with 16 nodes.

- A Fit-tree uses OCS to prune unused (or infrequently used) connections in a Fat-Tree
- Tailor the interconnect to match application data flows







- Silicon photonics enables optics to be integrated with conventional CMOS
- Enables up to 27x improvement in communication energy efficiency!







Memory







According to the Environmental Protection Agency (EPA), data centers consumed about 60 billion kilowatt-hours (kWh) in 2006, roughly 1.5 percent of total U.S. electricity consumption.



Slide from Dean Klein (Micron Technology)







Technology Challenge

Our ability to sense, collect, generate and calculate on data is growing faster than our ability to access, manage and even "store" that data •Memory density is doubling every three years; processor logic is every two •Storage costs (dollars/Mbyte) are dropping gradually compared to logic costs





1Gbit DDR3 Architecture









1Gbit DDR3 Architecture









Assumptions

Cell Voltage	1.2	V				
Cell Capacitance	25	fF				
Bitline Capacitance	75	fF				
Memory System Bandwidth	1	EB/sec				
Simplified Results:						
Energy/bit	36	fJ				
Total Memory Cell Power	288	KW				
With Bitline	1150	KW				
With 512X Over-Fetch	590	MW				







Conclusions

- Memory technology requires major reorganization (if industry stays alive)
 - More ranks/banks, Less over-fetch, new drivers
 - Chip stacking or optical memory interfaces
- We will have to live with less memory / computational performance
- We will have lower memory bandwidth/ computational performance (< 0.001 bytes/ flop)







Algorithms

Scaling to Billion-way Parallelism







Office of

Science

U.S. DEPARTMENT OF ENERGY

The Future of HPC System Concurrency



Must ride exponential wave of increasing concurrency for forseeable future!

Fortunately, most of the concurrency growth is within a single socket 34





Climate Model New Approaches for Massive Parallelism

- Existing Latitude-longitude based algorithm advection algorithm breaks down significantly before 1km scale!
 - Grid cell aspect ratio at the pole is 10000!
 - Advection time step is problematic at this scale
- Ultimately requires new discretization for atmosphere model
 - Must expose sufficient parallelism to exploit power-efficient design
 - Partner with CSU/Randall Group to use the Icosahedral Code
 - Uniform cell aspect ratio across globe


Where to Find 12 Orders in 10 years? Jardin & Keyes

- 1. overs: increased processor speed and efficiency
- 1.5 orders: increased concurrency
- 1 order: higher-order discretizations
 - Same accuracy can be achieved with many fewer elements
- 1 order: flux-surface following gridding
 - Less resolution required along than across field lines
- 4 orders: adaptive gridding
 - Zones requiring refinement are <1% of ITER volume and resolution requirements away from them are ~10² less severe
- 3 orders: implicit solvers
 - Mode growth time 9 orders longer than Alfven-limited CFL





m

Hardware:

5

Software:



Conclusion

- Consequence: must find strong-scaling from explicit parallelism
 - That's a tall order!
 - Used in 1980's to argue against MPPs







Software Performance

Auto-tuning: Don't depend on a human to do a machine's job.







Performance Profiles

(maintaining system balance)



- Neither memory bandwidth nor FLOPs dominate runtime
- The "other" category dominated by memory latency stalls
- Points to inadequacies in current CPU core design (inability to tolerate latency) Lets not forget about latency!







Auto-tuning

Problem: want to compare best potential performance of diverse architectures, avoiding

- Non-portable code
- Labor-intensive user optimizations for each specific architecture
- Our Solution: Auto-tuning
 - Automate search across a complex optimization space
 - Achieve performance far beyond current compilers
 - achieve performance portability for diverse architectures!







Conclusion

- Huge opportunities for energy-efficiency
 improvement simply by optimizing code performance
- Compilers cannot achieve this because of insufficient information
 - Assume flat machine model (which is wrong)
 - Cannot exploit domain-specific knowledge
- Auto-tuners
 - Can exploit domain-specific abstraction (motifs)
 - Can automate search of design space for performance portability
- Languages:
 - Need to expose correct machine model (flat model is wrong)
 - Need to express locality







How Can We Achieve our Goals Cost Effectively?







Intel HPC Market Overview



HPC is built with of pyramid investment model



U.S. DEPARTMENT OF ENERGY

ASC/OASCR Collabs

Dec 11, 2008



U.S. DEP

Processor Technology Trend

- 1990s R&D computing hardware dominated by desktop/COTS
 - -Had to learn how to use COTS technology for HPC
- 2010 R&D investments moving rapidly to consumer electronics/ embedded processing
 - Must learn how to leverage embedded processor technology for future HPC systems
 Market in Japan(B\$)



ERSOnsumer Electronics has Replaced PCs as NATIONAL ENERGY RESEARE Dominant Market Force in CPU Design!!









Embedded Design Automation (Example from Existing Tensilica Design Flow)







Technology Continuity for A Sustainable Hardware Ecosystem



ERSC

NATIONAL ENERGY RESEARCH SCIENTIFIC COMPUTING CENTER

U.S. DEPARTMENT OF ENERGY

Need building blocks for a compelling





If this is such a great idea, then why don't you do it?

Eating our own dogfood







Green Flash Overview

- Research effort: study feasibility and share insight w/community
- Elements of the approach
 - Choose the science target first (climate for example)
 - **Design systems for applications** (rather than the reverse)
 - Design hardware, software, scientific algorithms together using hardware emulation and auto-tuning

What is NEW about this approach

- Leverage commodity processes used to design power efficient embedded devices (redirect the tools to benefit scientific computing!)
- Auto-tuning to automate mapping of algorithm to complex hardware
- RAMP: Fast hardware-accelerated emulation of new chip designs

Applicable to broad range of scientific computing applications







Identify Target First! (Global Cloud Resolving Climate Model)

Surface Altitude (feet)



200km Typical resolution of IPCC AR4 models



25km Upper limit of climate models with cloud parameterizations 1km Cloud system resolving models are a transformational change





Requirements for 1km Climate Computer

Must maintain 1000x faster than real time for practical climate simulation

- ~2 million horizontal subdomains
- 100 Terabytes of Memory
 - 5MB memory per subdomain
- ~20 million total subdomains
 - 20 PF sustained (200PF peak)
 - Nearest-neighbor communication
- New discretization for climate model
 - CSU Icosahedral Code







Embedded Design Automation (Example from Existing Tensilica Design Flow)



ERSC Climate System Design Concept

Strawman Design Study

NATIONAL ENERGY RESEARCH SCIENTIFIC COMPUTING CENTER





Green Flash Strawman System Design In 2008

We examined three different approaches:

- AMD Opteron: Commodity approach, lower efficiency for scientific applications offset by cost efficiencies of mass market
- BlueGene: Generic embedded processor core and customize system-on-chip (SoC) services to improve power efficiency for scientific applications
- Tensilica XTensa: Customized embedded CPU w/SoC provides further power efficiency benefits but maintains programmability

Processor	Clock	Peak/ Core (Gflops)	Cores/ Socket	Sockets	Cores	Power	Cost 2008
AMD Opteron	2.8GHz	5.6	2	890K	1.7M	179 MW	\$1B+
IBM BG/P	850MHz	3.4	4	740K	3.0M	20 MW	\$1B+
Green Flash / Tensilica XTensa	650MHz	2.7	32	120K	4.0M	3 MW	\$75M

rrrr





What we have learned from our more detailed design study

Mark Horowitz 2007: "Years of research in lowpower embedded computing have shown only one design technique to reduce power: <u>reduce waste</u>."

Seymour Cray 1977: "Don't put anything in to a supercomputer that isn't necessary."





Peel Back the Historical Growth of Instruction Sets (accretion of cruft)



System Memory Interface Coherent MP Split Block Bridges Base 16b GP DSP Debua Interrupts Timers Memory MMU Computation Processor Instruction Control Set Time per variant: days Area = silicon cost and power BERKELEY LA

Configurable Processor Family



A Short List of x86 Opcodes that Science Applications Don't Need!

mnemonic	opl	<u>op2</u>	<u>op3</u>	<u>op4</u>	iext	pf	<u>OF</u> p	0 50	<u>0</u> <u>p</u> <u></u>	<u>oc</u> <u>s</u>	<u>t</u> m	<u>r1</u>	<u>x</u> <u>t</u>	ested f	modif f	<u>def f</u>	<u>undef f</u>	<u>f values</u>	description, notes
AAA	AL	ЪN.					3	7					1.	a	oszapc	a.c	05%.p.		ASCII Adjust After Addition
AAD	AL	AN					D	5 0 A					\square		oszapc	5z.p.	0 a .c		ASCII Adjust AX Before Division
AAM	AL	AN					D	4 0 A					H		oszapc	5z.p.	0 a .c		ASCII Adjust AX After Multiply
AAS	AL	AN					3	F					Π.	a	oszapc	a.c	052.p.		ASCII Adjust AL After Subtraction
ADC	r/m8	1 8					1	.0	т		+		L .	c	oszapc	oszapc			Add with Carry
ADC	r/m16/32/64	r15/32/54					1	1	т				L .	c	oszapc	oszapc			Add with Carry
ADC	rô	r/m8					1	.2	т				Π.	c	oszapc	oszapc			Add with Carry
ADC	r16/32/64	r/m15/32/54				\square	1	.3	т		+		Η.	c	oszapc	oszapc			Add with Carry
ADC	AL	imm8					1	.4			\top		Π.	c	oszapc	oszapc			Add with Carry
ADC	rAX	imm16/32					1	.5					1.	c	oszapc	oszapc			Add with Carry
ADC	r/m8	imm8					8	0	2		\top		L .	c	oszapc	ossapc			Add with Carry
ADC	r/m16/32/64	imm16/32					8	l	2				L .	c	oszapc	oszapc			Add with Carry
ADC	r/m8	imm8					8	2	2		\top		L .	c	ossapc	ossapc			Add with Carry
ADC	r/m16/32/64	imm8					8	3	2				L .	c	ossapc	ossapc			Add with Carry
ADD	r/m8	1 8					0	0	т				L		ossapc	ossapc			Add
ADD	r/m16/32/64	r15/32/54					0	1	т				L		ossapc	oszapc			Add
ADD	r 8	r/m8					0	2	т				\square		ossapc	ossapc			Add
ADD	r16/32/64	r/m16/32/64					0	3	т				\square		oszapc	oszapc			ådd .
ADD	AL	imm8					0	4					\square		ossapc	ossapc			Add
ADD	rAX	imm15/32					0	5					\square		ossapc	ossapc			Add
ADD	r/m8	imm8					8	0	0				L		ossapc	ossapc			Add
ADD	r/m16/32/64	imm15/32					8	1	0				L		ossapc	ossapc			Add
ADD	r/m8	imm8					8	2	0				L		ossapc	oszapc			Add
ADD	r/m16/32/64	imm8					8	3	0				L		oszapc	oszapc			Add
ADDPD	xanan.	xmm/m128			55e2	66	0F 5	8	r P4	+			\square						Add Packed Double-FP Values
ADDPS	xmm.	xmm/ml28			ssel		0F 5	8	r P3	+			\square						Add Packed Single-FP Values
ADD SD	xmm.	xmm/m64			55e2	F2	0F 5	8	r P4	+			\square						Add Scalar Double-FP Values
ADDSS	xmm.	xmm/m32			ssel	F3	0F 5	8	r P3	+			\square						Add Scalar Single-FP Values
ADDSUBPD	жтт.	xmm/m128			sse3	66	OF D	0	r P4	++									Packed Double-FP Add/Subtract
ADDSUBPS	xmm.	xmm/ml28			sse3	F 2	OF D	0	r P4	++			Π						Packed Single-FP Add/Subtract
ADX	AL	AN	imm8				D	5							oszapc	5z.p.	0 a .c		Adjust AX Before Division
ALTER						64			P4	+ v	1		Π						Alternating branch prefix (used only with Jcc instructions)
AVIX	AL	AN	imm8				D	4			+		H		oszapc	5z.p.	0 a .c		Adjust AX After Multiply
AND	r/m8	1 8				\square	2	0	т		+		L		oszapc	05z.pc	a	oc	Logical AND
AND	r/m16/32/64	r15/32/54					2	1	т		+		L		oszapc	05z.pc	a	oc	Logical AMD
AND	r 8	r/m8				\square	2	2	т		+		H		oszapc	05z.pc	a	oc	Logical AND
AND	r16/32/64	r/m15/32/54					2	3	т		+		H		oszapc	05%.pc	a	oc	Logical AMD
AND	AL	imm8				\square	2	4					\square		oszapc	05z.pc	a	oc	Logical AMD
AND	rAX	imm16/32					2	5			+		H		oszapc	052.pc	a	oc	Logical AMD
AND	r/m8	imm8				\square	8	0	4				L		oszapc	05z.pc	a	oc	Logical AMD
AND	r/m16/32/64	imm16/32					8	1	4				L		oszapc	052.pc	a	oc	Logical AMD
AND	r/m8	imm8				\square	8	2	4		+		L		oszapc	05z.pc	a	oc	Logical AMD
AND	r/m16/32/64	imm8				\square	8	3	4 03	+			L		oszapc	052.pc	a	oc	Logical AMD
ANDNPD	xmm.	xmm/m128			sse2	66	0F 5	5	r P4	+									Bitwise Logical AND NOT of Packed Double-FP Values
AND NP S	xanan.	xmm/m128			ssel		0F 5	5	r P3	+									Bitwise Logical AND NOT of Packed Single-FP Values
ANDPD	xanan.	xmm/m128			sse2	66	0F 5	4	r P4	+									Bitwise Logical AND of Packed Double-FP Values
ANDP 3	xanan.	xmm/m128			ssel		0F 5	4	r P3	+									Bitwise Logical AND of Packed Single-FP Values

UCICIICC

U.S. DEPARTMENT OF ENERGY





ARPL BOUN BSF BSR BSWA BT вт BTC BTC BTR BTR BTS BTS CALL CALL CALL CALL CALL CALL

CBW св₩ CWDE CDQE CDQ CLC CLD CLFL

CL I CLTS

CMC

CMOUNBE

CMOVA CMOMBEL

CMOVGE

M0176

CMOUNLE

Moro Wastad Oncodes

						CVTSD233	xmm.
						CVTS I2SD	xmm.
ARPL	r/m16		r16		смо	CVTS1233	xmm
DOOND	T10/32		mib/ 32610/ 32	eriags	СМО	CUTSS2SD	жттт.
B31	r16/32/	64	I/ M10/ 32/ 04		СМО	00733231	x32/64
BSR	r16/32/	64	r/m15/32/64			CUTTED 2D0	
BSWAP	r16/32/	64			CMP		Annun
вт	r/m16/3	82/54	r15/32/64		CMP	COTTPDZPI	mm.
BT	r/m16/3	32/54	imm8		CMDP	CUTTPS2DQ	XIIIII
BTC	r/m16/3	32/64	imm8		CMIP	CUTTP32PI	πm
BTC	r/m16/3	32/64	r15/32/54		CMDP	CVTTSD2SI	r32/64
BTR	r/m16/3	32/64	r16/32/64		CMDP	CUTTSS2SI	r32/64
BTR	r/m16/3	32/64	imm8		CMDP	CWD	DX
BTS	r/m16/3	32/64	r16/32/64		CMDP	CIMD	77
BTS	r/m16/3	32/64	imm8		CMDP	CTD	FDV
CALL	rel 16/3	32			CMDP		
CALL	re132				CMIP	CQU	RDX
CALL	r/m16/3	32			CMP	CWDE	EAX
CALL	r/m54				CMP	DAA	AL
CALLF	ptr16:1	L 5/ 32			CMIP	DAS	AL
CALLF	m15:15/						
CBW	AN	•V	Va on	ly no	A ha	20 Aut	t nt
CBW	AN			iy iic	eu c		
CWDE	EAX	-	4	4 •			
CDQE	RAX	In	struc	tion	set!		
CDQ	EDX				••••		
CLC							
CLD							
CLFLUSH	m8	~					\sim 7
CL I		•5	itill ha	ve al	l of t	ne 80	18
CLTS	CRO						• ·
смс		•\/	Vida S		Doo	en't N	1al
CMOVB	r16/32/	۷			DUE	511 L IV	/lar
CMOUNAE	r16/32/				0		<u> </u>
CMOVC	r16/32/	•	leithei	r doe	ຣ () ຊ	nche (Cor
CMOVBE	r16/32/						501
CMOUNA	r16/32,	- N	laitha				da
CMOVL	r16/32,	•1\	ieime	uue	5 AV		iue
CMOUNGE	r16/32,				_		_
CMOVLE	r16/32/		•Cr	eate	s nin	eline	hu
CMOUNG	r16/32/			cale	y hih	Cinic	bu
CMOUNB	r16/32/					···· 11 . 11	-
CMOUAE	r16/32/		•Be	etter t	o un	roll It	ac
CMOUNC	r16/32/						~ ~

	CVTPS2PD	xmm.	xmm/m128				U	160		
	CVTPS2PI	mm.	xmm/mδ4			-				
	CVTSD2SI	r32/64	xmm/mδ4							
	CUTSD2SS	жттт.	xmm/m54					F3(2)(A	SP	STi
	CUTSI2SD	xmm	r/m32/54		×15/32/64	r/m15/22/54				
101	CVTSI2SS	xmm.	r/m32/54		r16/32/64	r/m16/32/64		FXCH4	ST	STi
101	CVTSS2SD	xmm	xmm/m32	1	r16/32/64	r/m15/32/54		FXCH7	SI	STi
P	CVTSS2SI	r32/64	xmm/m32		r/m8	r 8		FXCH7	SI	STi
P	CVTTPD2DQ	xmm	xmm/m128		r/m16/32/64	r16/32/54		FXRSTOR	SI	ST1
P	CVTTPD2P I	ராக	xmm/m128		1 8	r/m8		FYRSTOR	50	501
P	CUTTPS2DQ	xmm	xmm/m128		r15/32/54	r/m15/32/54		TYPANT	- 510	CT .
P	CUTTPS2PI	π. π .	xmm/m54	l	AL	imm8		TASKOL	m512	31
P	CUTTSD2SI	r32/64	xmm/m54		TAX	imm15/32		FXSAUE	m512	57
P	CUTTSS2SI	r32/64	xmm/m32		r/m8	imm8		FXTRACT	SI	
P	CWD	DX	AX		r/m16/32/64	1700110/32	<u> </u>	FYL2X	ST1	ST
P	CWD	DX	AX		r/m15/22/54	immo	<u> </u>	FYL2XP1	ST1	ST
P	CDQ	EDX	BAX		Strate, State, State	xmm/m128	imm8	GS	GS	
PI PI	CQO	RDX	RAX		xmm	xmm/m128	imm8	HADDPD	XTTUTI.	xmm/m128
P	CWDE	ERX	AX		m8	m8		HADDES	NUMB	xmm/m128
P	DAA	AL.			70 B	m 8		HI.T		
=			<u> </u>		m15	m16				

f the nearly 300 ASM instructions in the x86

- and 8088 instructions!
- ke Sense with Small Cores
- nerence
- or Sqrt for loops
 - bbles
- ross the loops (like IBM MASS libraries) •Move TLB to memory interface because its still too huge (but still get

precise exceptions from segmented protection on each core)

Science U.S. DEPARTMENT OF ENERGY

r16/32 r16/32,

r16/32

r16/32 r16/32/

r16/32

INTO	eFlags	
INVD		
INVLPG	m	





-0000

Xtensa

-000-

æ

Xtensa

	Intel Core2 (Penryn)	Intel Atom core	Tensilica core w/ 64-bit FP
Die area (mm²)	53.5	25	5.32
Process	45 nm	45 nm	65 nm
Power	18W	0.625W	0.091W
Freq	2930 MHz	800MHz	370MHz
Flops / Watt	162	1280	4065







JUGILE

U.S. DEPARTMENT OF ENERGY

Global address space

Architectural Support for PModels Make hardware easier to program!



- Logical topology is a full crossbar
- Each local store mapped to global address space
- To initiate a DMA transfer between processors:
 - Processors exchange starting addresses through TIE Queue interface
 - Optimized for small transfers
 - When ready, copy done directly from LS to LS
 - Copy will bypass cache hierarchy



CMP Architecture - Physical View



ERSC

- Concentrated torus
 - Direct connect
 between 4
 processors on a
 tile
 - Packet switched network connecting tiles
- Between 64 and 128
 processors per die





Fault Tolerance/Resilience

- Our Design does not expose unique risks
 - Faults proportional to # sockets (not # cores) and silicon surface area
 - We expose less surface area and fewer sockets with our approach
- Hard Errors
 - Spare cores in design (Cisco Metro)
 - SoC design (fewer components and fewer sockets)
 - Use solder (not sockets)
- Soft Errors
 - ECC for memory and caches
 - On-board NVRAM controller for localized checkpoint
 - Checkpoint to neighbor for rollback





Memory: Perhaps we don't need Byte/FLOP (Scripted Memory Movement)

Trace analysis key to memory requirements

R

NATIONAL ENERGY RESEARCH

- Actually running the code gives realistic values for memory footprint, temporal reuse, DRAM bandwidth requirements
- Memory footprint: unique addresses accessed → size of local store needed
- Temporal reuse: maximum number of addresses which will be reused at any time → size of cache needed
- DRAM bandwidth
 - (instruction throughput) X (memory footprint)/(instruction count)







Bandwidth Requirements (MB/s) (Instructions/Cycle=1, 500 MHz)



Auto-Tuning Can Change Hardware Design Requirements





LH2 (small domain, reordered)



- Memory footprint: 160 KB
- Cache size requirement: 160 KB
- < 50% instructions are floating-point</p>
 - Huge overhead for address generation
- Although code streams through data, loop ordering was bad → cachelines reused although addresses were not
 Office of Science

- Memory footprint: 160 KB
- Cache size requirement: 1 KB
- > 85% instructions are floating-point
 - Good ordering → simpler addressing

160x reduction in cache size!2x savings in execution time





Generalized Stencil Auto-Tuning Framework

Ability to tune many stencil-like kernels

- No need to write kernel-specific perl scripts
- Uses semantic information from existing Fortran

Target multiple architectures

- Search over many optimizations for each architecture
- Currently supports multi/manycore, GPUs

Better performance = Better energy efficiency







Multi-Targeted Auto-Tuning

NATIONAL ENERGY RESEARCH SCIENTIFIC COMPUTING CENTER

ERSC



0

1

2 4 8 Threads

16

8

n 8 reads Science

U.S. DEPARTMENT OF ENERGY

0

1

16

2 4 8 16 Threads

12

10

8

4

2

GFlop/s



BERKELEY LAB



How many processors per chip



Design Trade-offs

pack fewer cores in socket to minimize memory bandwidth
maximize cores in socket

to minimize surface-to-

volume ratio

•Little's Law latency hiding

Logical View of Processor Network











BERKELEY LAB



Inserting Scientific Apps into the Hardware Development Process

- Research Accelerator for Multi-Processors
 (RAMP)
 - Simulate hardware before it is built!
 - Break slow feedback loop for system designs
 - Enables tightly coupled hardware/software/science co-design (not possible using conventional approach)







HW/SW Co-Tuning for Energy Efficiency

The approach: Use auto-tuned code when evaluating architecture design points



Co-Tuning can improve powerefficiency and area-efficiency by ~4x









Green Flash Hardware Demo

- Demonstrated during SC '08
- Proof of concept
 - CSU atmospheric model ported to Tensilica Architecture
 - Single Tensilica processor running atmospheric model at 50MHz
- Emulation performance advantage
 - Processor running at 50MHz vs. Functional model at 100 kHz
 - 500x Speedup
- Actual climate code not representative benchmark











Summary

- Power is leading design constraint for future HPC
 - Future technology driven by handheld space
 - Notion of "commodity" moving on-chip
- Approach for Power Efficient HPC
 - Choose the science target first (climate in this case)
 - **Design systems for applications** (rather than the reverse)
 - Design hardware, software, scientific algorithms together using hardware emulation and auto-tuning
 - This is the right way to design efficient HPC systems!






More Info

- Green Flash
 - <u>http://www.lbl.gov/CS/html/greenflash.html</u>
 - <u>http://www.lbl.gov/CS/html/greenmeetings.html</u>
- NERSC Science Driven System
 Architecture Group
 - http://www.nersc.gov/projects/SDSA











Ability to Verify

2004

2000

Ofer Sacham Stanford

20

10

1988

1992

1996





Source: SIA

Roadmap, 2001

Parallel Computing Everywhere Cisco CRS-1 Terabit Router



ERS

NATIONAL ENERGY RESEARCH

Replaces ASIC using 188 GP cores! Emulates ASIC at competitive power/performance Better power/performance than FPGA! New Definition for "Custom" in SoC



- 188+4 Xtensa general purpose processor cores per Silicon Packet Processor
- Up to 400,000 processors per system
 - (this is not just about HPC!!!)





Growth in Power Consumption (Top50) Excluding Cooling







Part III

A Short diversion on Metrics







Metrics: Can't improve what you don't measure

- Collecting Metrics for HPC Power consumption (Green500, Top500, SpecHPC)
 - Raise Community Awareness of HPC System
 Power Efficiency
 - Push vendors toward more power efficient solutions (shine a light on inefficiency)
- Choice of measurement has a dramatic effect on the outcome (Law of unintended consequences)
 - Suddenly everything is "green"
 - But is anything *really* getting better? (everything looks better on an exponential curve)







Anatomy of a "Value" Metric

Good Stuff

Bad Stuff







Anatomy of a "Value" Metric

Bogus!!! FLOP/s **Natts** Potentially Bogus!!







Anatomy of a "Value" Metric

Choose your own metric for performance! (doesn't need to be HPL, or FLOPS)



Formal process for collecting this data emerging (Green500, Top500, and eventually SpecPowerHPC)







Performance/measured_watt is much more useful than FLOPs/peak_watt But, are we getting the desired response?

