Life After Dennard and How I Learned to Love the Picojoule

STEVE KECKLER DIRECTOR OF ARCHITECTURE RESEARCH



Performance-oriented Computing is Everywhere



Providing Rich Immersive Experiences





Addressing Societal Challenges





Product Design and Simulation



Alternative Fuels



Medical Diagnosis



Climate Change



Materials Design



Understanding Viruses

Insatiable Demand for Performance













> 1 TFlop ~1KW





Which of These are Power Constrained?



NVIDIA



Silicon Scaling – the end of the easy ride

Where Does Processor Energy Go?

Strategies for Energy Efficiency

Challenges

Moore's Law is Only Part of the Story



Classic Dennard Scaling 2.8x chip capability in same power





Post Dennard Scaling 2x chip capability at 1.4x power 1.4x chip capability at same power





Implications



- End of clock-rate scaling => need more parallelism
 - Capability drops down to 1.4 per generation.
 - Optimistic because wires scale worse than transistors
 - 32x gap in a decade from what continued Dennard scaling would give
- Scaling up conventional core count is not feasible
 - Best case 40% more cores per generation
 - Far cry from 2x performance per generation
- Need better materials, devices, circuits, architectures, and SW © NVIDIA 2011

How is power spent in a CPU?





Fundamental Energy Costs





Today's high perf CPUs 1-2 nJ/instruction

20-40x energy overhead

- Instruction control
- Data movement

Energy Shopping List



Processor Technology	40 nm	10nm	
Vdd (nominal)	0.9 V	0.7 V	
DFMA energy	50 pJ	7.6 pJ	
64b 8 KB SRAM Rd	14 pJ	2.1 pJ	
Wire energy (256 bits, 10mm)	310 pJ	174 pJ	
	100		
Memory Technology	45 nm	16nm	
DRAM interface pin bandwidth	4 Gbps	50 Gbps	
DRAM interface energy	20-30 pJ/bit	2 pJ/bit	
DRAM access energy	8-15 pJ/bit	2.5 pJ/bit	
	Vogelsang [Micro 2010]		

FP Op lower bound = 8 pJ



Strategies for Energy Reduction

- Simpler processor architectures
- Reduce waste
- Improve physical locality of data
- Exploit heterogeneity and specialization
- Signaling and packaging
- Inch down voltage

Simpler/Slower Cores = Energy Efficiency











Throughput Processors at Petascale





#1 : K Computer 68K Fujitsu Sparc CPUs 10.5 PFLOPS, 12.6MW



#2 : Tianhe-1A 7K Tesla GPUs 2.6 PFLOPS, 4MW

#4 : Nebulae 4K Tesla GPUs 1.3 PFLOPS, 2.5MW



#3 : Jaguar 36K AMD Opteron CPUs 1.8 PFLOPS, 6.9MW © NVIDIA 2011

Titan 18K Tesla GPUs >20 PFLOPS, 8.6MW



#5 : Tsubame 2.0 4K Tesla GPUs 1.2 PFLOPS, 1.3MW (most efficient PF system)



Throughput Processor Optimizations



Hierarchical hardware thread scheduling

Register File caching

SIMT and Temporal SIMT execution

Streaming Multiprocessor (SM)



Multithreading with: 32 threads per Warp 48 Warps/SM

Large register files

- 1500+ threads/SM
- Fermi: 128KB of RF
- Heavily banked to provide high bandwidth
- Large operand read/write energy
 - ~15% of power of SM



Optimization Opportunities



- Large number of threads hide two types of latency
 - Long: global memory access (~400 cycles)
 - Short: ALU and shared memory access (8-20 cycles)



Two-Level Scheduling

- Active warps used to tolerate short latency events
- Simplified instruction scheduler only considers active warps each cycle
- 8 active warps (of 48) is enough to sustain full SM throughput



Gebhart, et. al [ISCA 2011, Micro 2011]

Register File Caching (RFC)



- 4-6 entries per thread
- Operand routing when results are needed by shared units
- 48x smaller than MRF
 - When combined with 2-level scheduling
- Reduces up to 80% of RF reads/writes



SM Cluster - replicated 8 times to form SM

© NVIDIA 2011

NVIDIA

SIMD versus MIMD versus SIMT?







Temporal SIMT

Spatial SIMT (current GPUs)

32-wide datapath

🔿 🖿> АААААААААААААААААААААААААААААААА

Pure Temporal SIMT





Temporal SIMT Optimizations



Control divergence – hybrid MIMD/SIMT



Strategies for Energy Reduction



- Simpler processor architectures
- Reduce waste
- Improve physical locality of data
- Exploit heterogeneity and specialization
- Signaling and packaging
- Inch down voltage

Heterogeneity and Specialization

Programmable processor heterogeneity
CPU + GPU in same system/chip

Functional specialization
Fixed/limited programmability accelerators

Power specialization

Same function, different power/performance profile

NVIDIA Tegra3 Mobile SOC



5 ARM Cortex A9 Cores
Range of accelerators
Some programmable

- Some fixed function
- **1-5W power envelope**



Variable SMP



- Same CPU core design
- Optimized for different voltage/frequency range
 - **Different transistors**
 - Different synthesis

Low Power Process: inefficient at high performance ranges

Fast Process: higher leakage in active standby



Typical Power Breakdown





Packaging/Signaling/Architecture





Micron Hybrid Memory Cube



Silicon Interposer Packaging (Xilinx)

Voltage Scaling – is there any left?



- Answer #1: Hard scale V_{dd} down toward 2*V_t
- Answer #2: Harder scale down to near V_t
 - Principle challenges
 - Slow transistors
 - Reliability
 - Process variation

Aggressive Voltage Scaling





Energy Efficiency Potential at 10nm



Natural process scaling	6x	Assuming 0.9=>0.7V	
Energy efficient chip architecture	3-5x		
Circuits/Packaging	2-2.5x	Includes on-chip ckts	
Slower clock frequency	2x		
Aggressive voltage scaling	1.5x		
Total:	20-40x	~3x area penalty	

Architecture specialization	~10x	On compute limited applications
Near-threshold Voltage	4x (theoretical)	~30x area penalty



"Teraflops"





ASCI Red @ Sandia Labs





2019

5 Watts

1997

© NVIDIA 2011

40

"Tens" of Teraflops

100 Watts





"Hundreds" of Teraflops

1000 Watts

NVIDIA





Challenges for Research Community

- Fundamentally energy efficient architectures
 - Core, memory system, interconnect, system
- Low voltage
 - Clocking, RAM, signaling, variation tolerant
- Specialization
 - In general purpose system what are the "accelerators"?
- Emerging and important workloads are throughput



'Super' Computing From Super Computers to Super Phones