DVIDIA.

Challenges for Future Computing Systems

Bill Dally | Chief Scientist and SVP, Research NVIDIA | Professor (Research), EE&CS, Stanford

Exascale Computing Will Enable Transformational Science



Climate



Comprehensive Earth System Model at 1KM scale, enabling modeling of cloud convection and ocean eddies.





Combustion



First-principles simulation of combustion for new high-efficiency, low-emision 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.





Exascale Computing Will Enable Transformational Science

High-Performance Computers are Scientific Instruments



TITAN

18,688 NVIDIA Tesla K20X GPUs27 Petaflops Peak: 90% of Performance from GPUs17.59 Petaflops Sustained Performance on Linpack



Tsubame KFC 4.5GFLOPS/W #1 on Green500 List





US to Build Two Flagship Supercomputers







SIERRA

150-300 PFLOPS Peak Performance IBM POWER9 CPU + NVIDIA Volta GPU NVLink High Speed Interconnect 40 TFLOPS per Node, >3,400 Nodes

2017

Major Step Forward on the Path to Exascale



Scientific Discovery and Business Analytics

Driving an Insatiable Demand for More Computing Performance









2 Trends

The End of Dennard Scaling

Pervasive Parallelism



2 Trends

The End of Dennard Scaling

Pervasive Parallelism



Moore's Law doubles processor performance every few years, right?



Moore's Law doubles processor performance every few years, right?

Wrong!



Moore's Law gives us transistors Which we used to turn into scalar performance



Moore, Electronics 38(8) April 19, 1965







Classic Dennard Scaling 2.8x chip capability in same power





But L³ energy scaling ended in 2005



Moore, ISSCC Keynote, 2003











"Moore's Law gives us more transistors... Dennard scaling made them useful."



Bob Colwell, DAC 2013, June 4, 2013



Also, ILP was 'mined out' in 2000



ISAT LCC: 24

Dally et al. "The Last Classical Computer", ISAT Study, 2001



Result: The End of Historic Scaling



Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten Dotted line extrapolations by C. Moore

C Moore, Data Processing in ExaScale-ClassComputer Systems, Salishan, April 2011

<mark> NVIDIA</mark>.

The End of Dennard Scaling

- Processors aren't getting faster, just wider
 - Future gains in performance are from parallelism
- Future systems are energy limited
 - Efficiency /S Performance
- Process matters less
 - One generation is 1.2x, not 2.8x

Its not about the FLOPs

DFMA 0.01mm² 10pJ/OP – 2GFLOPs

A chip with 10⁴ FPUs: 100mm² 200W 20TFLOPS

Pack 50,000 of these in racks 1EFLOPS 10MW

16nm chip, 10mm on a side, 200W



Overhead

Locality





Optimized for Latency Caches



Westmere 32 nm

GPU 140 pJ/flop

Optimized for Throughput Explicit Management of On-chip Memory



Kepler 28 nm



Fixed-Function Logic is Even More Efficient

	Energy/Op
CPU	1.7nJ
GPU	140pJ
Fixed-Function	10pJ



How is Power Spent in a CPU?



Dally [2008] (Embedded in-order CPU)

Natarajan [2003] (Alpha 21264)









4/11/

NVIDIA.



📀 NVIDIA.

Simpler Cores = Energy Efficiency





Payload Arithmetic 20pJ

Overhead 20pJ





15% of SM Energy



Hierarchical Register File



<mark> NVIDIA</mark>.

Register File Caching (RFC)





Energy Savings from RF Hierarchy 54% Energy Reduction







Temporal SIMT

Spatial SIMT (current GPUs)

32-wide datapath







d

ld

ld

C

ld ld

ld

ld

ld

ld ld

0

2

3

5 6

7

8

9

10

(threads)



1 warp instruction = 32 threads

Temporal SIMT Optimizations

Control divergence – hybrid MIMD/SIMT





Scalar Instructions in SIMT Lanes



Scalarization Eliminates Work



Communication Dominates Arithmetic





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

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

FP Op lower bound = 4 pJ







DRAM Wordline Segmentation

Local Wordline Segmentation (LWS)

- Activate a portion of each mat
- Activate the same segment in all mats
- Master Wordline Segmentation (MWS)
 - Activate subset of mats
 - New data layout such that a DRAM atom resides in the activated mats
 - Multiple narrow subchannels in each channel





Activation Energy Reduction



Reductions in DRAM Row Energy

- LWS: 26% MWS: 49%
- LWS+MWS: up to 73%, average 68%



GRS Test Chips







Eye Diagram from Probe



Test Chip #2 fabricated on production GPU

Poulton et al. ISSCC 2013, JSSCC Dec 2013



Circuits: Ground Referenced Signaling

Goal: reduce Energy/bit 200fJ/bit-mm → 20fJ/bit-mm



Repeater

On-chip Signaling

- Testsite on GM2xx GPU
- <45fJ/bit-mm

Measurement Results



With GPU noise

Without GPU noise



The Locality Challenge Data Movement Energy 77pJ/F





Optimized Architecture & Circuits 77pJ/F -> 18pJ/F



<mark> (</mark>NVIDIA





2 Trends

The End of Dennard Scaling

Pervasive Parallelism



Processors aren't getting faster just wider



Thread Count (CPU+GPU)

	2013 (28nm)	2020 (7nm)
Cell Phone	4+1,000	16+4,000
Server Node	12+10,000	32+40,000
Supercomputer	10 ⁵ +10 ⁸	10 ⁶ +10 ⁹



Parallel programming is not inherently any more difficult than serial programming

However, we can make it a lot more difficult



A simple parallel program

```
forall molecule in set { // launch a thread array
   forall neighbor in molecule.neighbors { // nested
      forall force in forces { // doubly nested
      molecule.force =
      reduce_sum(force(molecule, neighbor))
   }
```



Why is this easy?

```
forall molecule in set { // launch a thread array
   forall neighbor in molecule.neighbors { // nested
      forall force in forces { // doubly nested
      molecule.force =
      reduce_sum(force(molecule, neighbor))
   }
```

No machine details All parallelism is expressed Synchronization is semantic (in reduction)



```
pid = fork() ; // explicitly managing threads
```

```
lock(struct.lock) ; // complicated, error-prone synchronization
// manipulate struct
unlock(struct.lock) ;
```

code = send(pid, tag, &msg) ; // partition across nodes

We could make it hard



Programmers, tools, and architecture Need to play their positions



Map foralls in time and space Map molecules across memories Stage data up/down hierarchy Select mechanisms Exposed storage hierarchy Fast comm/sync/thread mechanisms







Optimized Circuits 77pJ/F -> 18pJ/F



© NVIDIA.

Autotuned Software 18pJ/F -> 9pJ/F





Conclusion

Power-limited: from data centers to cell phones

Perf/W is Perf

Throughput cores

Reduce overhead



Data movement

۲

- **Circuits:** 200 -> 20
- Optimized software

Parallel programming is simple - we can make it hard

Target-independent programming - mapping via tools

