FLOPS to BYTES: Towards Accelerating Beyond Moore's Law

TSUBAME 2.0

Satoshi Matsuoka Professor, GSIC, Tokyo Institute of Technology / Director, AIST-Tokyo Tech. Big Data Open Innovation Lab / Fellow, Artificial Intelligence Research Center, AIST, Japan / Vis. Researcher, Advanced Institute for Computational Science, Riken

PPAM2017 2017/9/11 Lublin, Poland Overview of TSUBAME3.0 BYTES-centric Architecture, Scalaibility to all 2160 GPUs, all nodes, the entire memory hiearchy



Hot Pluggable ICE-

Smaller than server no cables or pipes

The Martin



Figuid Cooled NVMe

9-25/11 83071

1165

Xeon x 2

20 PeraFloos

256 Ci Byte Memory

D

PCIe NVMe

rive Bay x

15 Compute Racks DDN Storage Racks 3 Perioheral & SW racks

Total 22 Racks

Integrated 100/200Gbps Fabric Backplane

Toxyo Institute of Tec

CPUS/Fack

CP

Running Benchmarks (HPL, HPCG) Lokyo Lech HPE/Sci Japan Jean Het Performance

Peta

Green

0

0

0

SOQ CERTIFICATE

Tsubame 3.0, a modified HPE ICE XA System at the GSIC Center, Tokyo Institute of Technology

is ranked

--- No. 1 in the Green500

among the World's TOP500 Supercomputers with 14.1 GFlops/Watt Linpack Power-Efficiency on the Green500 List published at ISC High Performance, June 19, 2017

Congratulations from the Green500 Editors

Kirk Cameron Virginia Tech

2016 #1 = 9.46

erecce.

Announcement SC170 pening, 9 AM Monday, 19th Stree World #1 !!!

Benchmark team



Tremendous Recent Rise in Interest by the Japanese Government on Big Data, DL, AI, and IoT

- Three national centers on Big Data and AI launched by three competing Ministries for FY 2016 (Apr 2015-)
 - METI AIRC (Artificial Intelligence Research Center): AIST (AIST internal budget + > \$200 million FY 2017), April 2015
 - Broad AI/BD/IoT, industry focus
 - MEXT AIP (Artificial Intelligence Platform): Riken and other institutions (\$~50 mil), April 2016
 - A separate Post-K related AI funding as well.
 - Narrowly focused on DNN
 - MOST Universal Communication Lab: NICT (\$50~55 mil)
 - Brain –related AI
 - \$1 billion commitment on inter-ministry AI research over
 10 years => Supplanting HPC activities?



Vice Minsiter Tsuchiya@MEXT Annoucing AIP estabishment

7



2015- AI Research Center (AIRC), AIST





Director: Jun-ichi Tsujii

appointment as

since July 2017





Other Big Data / AI research organizations and proposals JST BigData CREST JST AI CREST





METI AIST-AIRC ABCI



as the worlds first large-scale OPEN AI Infrastructure

- ABCI: <u>AI</u> Bridging <u>Cloud</u> Infrastructure
 - Top-Level SC compute & data capability for DNN (>130 AI-Petaflops)
 - <u>Open Public & Dedicated</u> infrastructure for AI & Big Data Algorithms, Software and Applications
 - Platform to accelerate joint academic-industry R&D for AI in Japan



- < 3MW Power
- < 1.1 Avg. PUE
- Operational 2017Q4 ~2018Q1







The "Real" ABCI – 2018Q1

• Extreme computing power

- w/ >130 AI-PFlops for AI/ML especially DNN
- x1 million speedup over high-end PC: 1 Day training for 3000-Year DNN training job
- TSUBAME-KFC (1.4 AI-Pflops) x 90 users (T2 avg)

• Big Data and HPC converged modern design

- For advanced data analytics (Big Data) and scientific simulation (HPC), etc.
- Leverage Tokyo Tech's "TSUBAME3" design, <u>but differences/enhancements</u> <u>being AI/BD centric</u>
- Ultra high BW & Low latency memory, network, and storage
 - For accelerating various AI/BD workloads
 - Data-centric architecture, optimizes data movement
- Big Data/AI and HPC SW Stack Convergence
 - Incl. results from JST-CREST EBD
 - Wide contributions from the PC Cluster community desirable.
- Ultra-Green (PUE<1.1), High Thermal (60KW) Rack
 - Custom, warehouse-like IDC building and internal pods
 - Final "commoditization" of HPC technologies into Clouds









ABCI Cloud Infrastructure

Ultra-dense IDC design from ground-up

- Custom inexpensive lightweight "warehouse" building w/ substantial ABCI AI-IDC CG Image earthquake tolerance
- x20 thermal density of standard IDC
- Extreme green
 - Ambient warm liquid cooling, large Li-ion battery storage, and highefficiency power supplies, etc.
 - Commoditizing supercomputer cooling technologies Clouds (60KW/rack)
- Cloud ecosystem
 - Wide-ranging Big Data and HPC standard software stacks
- Advanced cloud-based operation
 - Incl. dynamic deployment, container-based virtualized provisioning, multitenant partitioning, and automatic failure recovery, etc.
 - Joining HPC and Cloud Software stack for real
- Final piece in the commoditization of HPC (into IDC)
- Open Sourcing of Next-Gen IDC Architecture for AI







ABCI Datacenter ~\$10 million (Just broke ground, to be completed late 2017)



Cutting Edge Research AI Infrastructures in Japan Accelerating **BD/AI** with HPC

Being

(and my effort to design & build them)

In Production



Oct. 2015 **TSUBAME-KFC/DL** (Tokyo Tech.) 1.4 AI-PF(Petaflops)



Manufactured x2.8~4.2 Aug. 2017 TSUBAME3.0 (Tokyo Tech.) 47.2 AI-PF (65.8 AI-PF w/Tsubame2.5)

x5.0~7.7 Mar. 2018 **ABCI (AIST-AIRC)** 130-200 AI-PF

Draft RFC out

construction

IDC under

~x1000 in 3.5 years

Mar. 2019? "ExaAI" ~1 AI-ExaFlop Undergoing Engineering Study

凤

Tokyo Tech



Mar. 2017 Al Supercomputer **Riken AIP** 4.1 AI-PF

R&D Investments into world leading AI/BD HW & SW & Algorithms and their co-design for cutting edge Infrastructure absolutely necessary (just as is with Japan Post-K and US ECP in HPC)

ORNL Summit

• ~200 Petaflops FP64, ~3 Exaflop FP16 by 1H2018



Peak FP64 Exaflop TSUBAME in 2020 – Just getting Flops/W is within reach

- 7nm+ post Volta GPU (Pascal P100 16nm)
- ~10,000 CUDA Cores (P100 3840), 12.5 Teraflops/Chip (P100 5.3TF) w/matrix engine
- 80,000 chips => 80 million small cores
- 4 GPUs/node => 20,000 nodes (x40 TSUBAME3, 500~600 racks)
- Scalable High-Dimensional Torus or Hypercube topology (Tsubame3 : Full Fattree)
- x3 power efficiency 50GF/W (x1.9 via process, x1.6 via arch) (TSUBAME3.0 14.1GF/W)
- 1 Exa DFP Peak, ~600 PF Linpack, 12MW Facility Power
- So the DARPA Exascale report projection turned out to be fairly accurate
- But is just getting FLOPS all that valuable?



ENERGY EFFICIENCY

Ē

The GREEN



50GFlops/W in late 2020?



Many core was a good step but we already used it once, and cannot use it again for boosting



Measured for the 2011 Gordon Bell Award Dendritic Solidification App Flop/s/W = Total #Flops / J = energy to solution given same problem What is worse: Moore's Law will end in the 2020's

- Much of underlying IT performance growth due to Moore's law
 - "LSI: x2 transistors in 1~1.5 years"
 - Causing qualitative "leaps" in IT and societal innovations
 - The main reason we have supercomputers and Google...
- •But this is slowing down & ending, by mid 2020s...!!!
 - End of Lithography shrinks
 - End of Dennard scaling
 - End of Fab Economics

The curse of <u>constant</u> <u>transistor power</u> shall



Gordon Moore

- How do we sustain "performance growth" beyond the "end of Moore"?
 - Not just one-time speed bumps
 - Will affect all aspects of IT, including BD/AI/ML/IoT, not just HPC
 - End of IT as we know it

20 year Eras towards of End of Moore's Law



flat performance

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

Need to realize the next 20-year era of supercomputing

The "curse of constant transistor power"

- Ignorance of this is like ignoring global warming -
- Systems people have been telling the algorithm people that "FLOPS will be free, bandwidth is important, so devise algorithms under that assumption"
- This will certainly be true until exascale in 2020...
- But when Moore's Law ends in 2025-2030, constant transistor power (esp. for logic) = FLOPS will no longer be free!
- So algorithms that simply increase arithmetic intensity will no longer scale beyond that point
- Like countering global warming need disruptive change in computing in HW-SW-Alg-Apps etc. for the next 20 year era

Performance growth via <u>data-centric computing:</u> <u>"From FLOPS to BYTES"</u>

- Identify the new parameter(s) for scaling over time
- Because data-related parameters (e.g. capacity and bandwidth) will still likely continue to grow towards 2040s
- Can grow transistor# for compute, but CANNOT use them AT THE SAME TIME(Dark Silicon) => <u>multiple computing units specialized to type of data</u>
- <u>Continued capacity growth</u>: 3D stacking (esp. direct silicon layering) and low power NVM (e.g. ReRAM)
- <u>Continued BW growth</u>: Data movement energy will be <u>capped constant</u> by dense 3D design and advanced optics from silicon photonics technologies
- Almost back to the old "vector" days(?), but no free lunch latency still problem, locality still important, need <u>general algorithmic acceleration</u> <u>thru data capacity and bandwidth</u>, not FLOPS

Leadership Systems Ratios (Tier1 memory)(slide courtesy Nic Dube@HPE)



Enterprise

Leadership Systems Ratios (sum)) (slide courtesy Nic Dube@HPE)



Leadership Systems Ratios (wAvg) (slide courtesy Nic Dube@HPE)



Hewlett Packard Enterprise

Y	#	Τ	HPCG Top 10 ranking June 201	Manufacturer	Computer	Country	HPCG [Pflop/s]	Rmax [Pflop/s]	HPCG/ Peak	HPCG/ HPL
	1	8	RIKEN Advanced Institute for Computational Science	Fujitsu	K Computer SPARC64 VIIIfx 2.0GHz, Tofu Interconnect	Japan	0.6027	10.5	5.3%	5.7%
	2	2	National University of Defense Technology	NUDT	Tianhe-2 NUDT TH-IVB-FEP, Xeon 12C 2.2GHz, IntelXeon Phi	China	0.5801	33.9	1.1%	1.7%
	3	3	Swiss National Supercomputing Centre (CSCS)	Cray	Piz Daint Cray XC50, Xeon E5 12C 2.6GHz, Aries, NVIDIA Tesla P100	Switzerland	0.4700	19.6	1.9%	2.4%
	4	7	JCAHPC Joint Center for Advanced HPC	Fujitsu	Oakforest-PACS PRIMERGY CX1640 M1, Intel Xeons Phi 7250 68C 1.4 GHz, OmniPath	Japan	0.3855	13.6	1.5%	2.8%
	5	1	National Supercomputing Center in Wuxi	NRCPC	Sunway TaihuLight NRCPC Sunway SW26010, 260C 1.45GHz	China	0.3712	93.0	0.3%	0.4%
	6	6	Lawrence Berkeley National Laboratory	Cray	Cori Cray XC40, Intel Xeons Phi 7250 68C 1.4 GHz, Aries	USA	0.3554	14.0	1.3%	2.5%
	7	5	Lawrence Livermore National Laboratory	IBM	Sequoia BlueGene/Q, Power BQC 16C 1.6GHz, Custom	USA	0.3304	17.2	1.6%	1.9%
	8	4	Oak Ridge National Laboratory	Cray	Titan Cray XK7, Opteron 16C 2.2GHz, Gemini, NVIDIA K20x	USA	0.3223	17.6	1.2%	1.8%
	9	10	Los Alamos NL / Sandia NL	Cray	Trinity Cray XC40, Xeon E5 16C 2.3GHz, Aries	USA	0.1826	8.10	1.6%	2.3%
-	10	15	NASA/ Ames Research Center/NAS	HPE	Pleiades SGI ICE X,	USA	0.1750	5.95	2.5%	2.9%

Sparse BYTES: The Graph500 – 2015~2016 – world #1 x 4 K Computer #1 Tokyo Tech[Matsuoka EBD CREST] Univ. Kyushu [Fujisawa Graph CREST], Riken AICS, Fujitsu





Characteristics of Big Data and Al Computing As BD / Al Dense LA: DNN

Graph Analytics e.g. Social Networks

Sort, Hash, e.g. DB, log analysis

Symbolic Processing: Traditional AI



As HPC Task Integer Ops & Sparse Matrices Data Movement, Large Memory Sparse and Random Data, Low Locality

Acceleration, Scaling

Opposite ends of HPC computing spectrum, but HPC simulation apps can also be categorized likewise



Acceleration via Supercomputers adapted to AI/BD

Inference, Training, Generation



As HPC Task Dense Matrices, Reduced Precision Dense and well organized neworks and Data



Acceleration, Scaling

(Big Data) BYTES capabilities, in bandwidth and capacity, unilaterally important but often missing from modern HPC machines in their pursuit of FLOPS...

0.8

0.6

0.2

음 – Nodes

Exec<u>ution tim</u>e [s]

- Need <u>BOTH bandwidth and capacity</u> (BYTES) in a HPC-BD/AI machine:
 - Obvious for lefthand sparse ,bandwidthdominated apps
 - But also for righthand DNN: Strong scaling, large networks and datasets, in particular for future 3D dataset analysis such as CTscans, seismic simu. vs. analysis...)



(Source: http://www.dgi.com/images/cvmain_overview/CV4DOverview_Model_001.jpg)





important

Our measurement on

Number of GPUs = 8 per node

TSUBAME3: A Massively BYTES Centric Architecture for Converged BD/AI and HPC



~4 Terabytes/node Hierarchical Memory for Big Data / AI (c.f. K-compuer 16GB/node)

➔ Over 2 Petabytes in TSUBAME3, Can be moved at 54 Terabyte/s or 1.7 Zetabytes / year

TSUBAME3: A Massively BYTES Centric Architecture for Converged BD/AI and HPC



~4 Terabytes/node Hierarchical Memory for Big Data / AI (c.f. K-compuer 16GB/node)

➔ Over 2 Petabytes in TSUBAME3, Can be moved at 54 Terabyte/s or 1.7 Zetabytes / year



Fujitsu SPARC64[™] Xifx (2015)



x10 ports

Japan Flagship 2020 "Post K" Supercomputer

✓CPU

- <u>ARM v8 + 512 bit SVE extensions</u>
- Multi-hundred petaflops peak total
- Power Knob feature for saving power

✓Memory

- ✓ 3-D stacked DRAM, Terabyte/s BW
- ✓Interconnect
 - TOFU3 CPU-integrated 6-D torus network
- I/O acceleration
- <u>30MW+ Power</u>
- Being designed and will be manufactured by Fujitsu
- Development Leaders: Yutaka Ishikawa, Mitsuhisa Sato (Riken)



Prime Minister Abe visiting K Computer 2013




Post-K will be/have

- Continuum of K
- Pretty good DFP FLOPS
- Very Good low precision FLOPS
- Awesome Memory Bandwidth
- (but) relatively low memory capacity
- Awesome Network Injection Bandwidth
- Very Good Network Bisection Bandwidth
- Relatively low I/O speed



Many Core Era



Post Moore Era



Flops-Centric Algorithms and Apps

Flops-Centric System Software



Hardware/Software System APIs Flops-Centric Massively Parallel Architecture



Transistor Lithography Scaling (CMOS Logic Circuits, DRAM/SRAM) ~2025 M-P Extinction Event Bytes-Centric Algorithms and Apps

Bytes-Centric System Software

Hardware/Software System APIs Data-Centric Heterogeneous Architecture



Novel Devices + CMOS (Dark Silicon) (Nanophotonics, Non-Volatile Devices etc.) Post-Moore is NOT a More-Moore device as a panacea

Device & arch. advances improving data-related parameters over time

Runtime "Rebooting Computing" in terms of devices, architectures, software.New memory Devices PC-RAM Algorithms, and ReRAM applications necessary STT-MRAM => Co-Design even 3D architecture more important fabrication c.f. Exascale



Problem Specific Architectures to exploit dark silicone "What are they good for?" – c.f. Berkeley Dwarfs

- Deep Neural Network Accelerator (Many, incl. Google)
- Spiking Neuromorphic Architecture (Manchester SpiNNaker, IBM TrueNorth, Heidelberg BrainScaleS)
- Ising Model optimization architecture (Hitachi)
- Automata Processor (Micron)
- Advanced FPGAs (Alltera, Xilinx)
- Network & I/O accelerator (Mellanox)
- .
- And of course Quantum Annealing and Computing (D-Wave)

Fujitsu Deep Learning Processor (DLUTM) Fujitsu





Supercomputer K technologies

- Architecture designed for Deep Learning
- High performance HBM2 memory
- Low power design

DLU[™] features

→ Goal: 10x Performance/Watt compared to others



Massively parallel : Apply supercomputer interconnect technology

- → Ability to handle large scale neural networks
- → TOFU Network derivative for massive scaling

Designed for Scalable Learning, technically superior to Google TPU2

"Exascale" Al possible in 1H2019

Neuromorphic Architectures (Not to be confused with DNN Accelerators)



IBM TrueNorth







Asynchronous Spiking Neural Computation Hardware for low-power real-time operation in Closed-Loop Systems

... to simulate 1 Billion Spiking Neurons in real-time



- Multi-channel spiking input and output
- Stand-alone spiking computing system
- Simulates ~20.000 neurons in real time
- Small (~20x20mm); low power (~600mW)
- Flexibly configurable, extendable, stackable

Manchester SpiNNaker (ARM Based)

Computing with Ising model

- Ising model: expressing behavior of magnetic spins
- Using Ising model as natural phenomenon to map problems

CMOS Ising computing

- Mimicking Ising model with CMOS circuits
- Easy to manufacture, easy to use, good scalability



- Hitachi@ISSCC2015 "An 1800-Times-Higher Power-Efficient 20k-spin Ising Chip for Combinational Optimization Problem with CMOS Annealing"
- Competitive to Quantum Annealing, room temperature, easy to scale
- Could be applicable to deep learning?

	8 <u>-</u>
1k-spin sub-array ∕ 780 x 380 um2	Numb P Cl
	Are
1k-spin	Nu SR
SDAM VE	Me
	Intera
← 4 mm	Opera of co

© 2015 IEEE International Solid-State Circuits Conference

Items	Value			
Number of spins	20k (80 x 256)			
Process	65 nm			
Chip area	4x3=12 mm ²			
Area of spin	11.27 x 23.94 =270 μm ²			
Number of SRAM cells	260k bits Spin value: 1 bit Interaction factor: 2 bit x 5=10 bits External magnetic coefficient: 2 bits			
Memory IF	100 MHz			
Interaction speed	100 MHz			
Operating current of core circuits (1.1 V) Do not include IO	Write: 2.0 mA Read: 6.0 mA Interaction: 44.6 mA			

 1,800 times higher energy efficiency than conventional approximation algorithm on CPU



Randomly generated problems, energy for same preciseness solution Ising chip: VDD=1.1 V, 100-MHz interaction, best solution among 10-times trial is selected. Approximation algorithm: SG3(*) is operated on Core i5, 1.87 GHz, 10 W/core.

> 24.3: An 1900-Times-Higher Power-Efficient 20k-apin Jaing Chip for Combinational Optimization Problem with CMOS Annealing

(*): Sera Kahruman et al., "On Greedy Construction Heuristics for the Max-Cut Problem," International Journal on Computational Science and Engineering, Volume 3, Number 3/2007, pp. 211-21



21 o © 2015 IEEE International Solid-State Circuits Conference

Conditions:

Towards Understanding the Performance of FPGAs using OpenCL Benchmarks [HiPEAC Reconfigurable Computing WS 2015 Extended version to appear SC16]

> Hamid Zohouri (Tokyo Tech), Naoya Maruyama (Riken AICS), Satoshi Matsuoka (Tokyo Tech), Motohiko Matsuda (RIKEN AICS)





In collaboration with: Aaron Smith (Microsoft Research),

Supported by Altera

Microsoft[®] Research

Parallelism in Altera OpenCL

Explicit: Thread/SIMD parallelism

Implicit: Pipeline parallelism

Inter pipelines

 Configurable number of duplicated pipelines

Intra pipeline

SIMD

Instantiate SIMD units base on user direction (attribute num_simd_work_items)

Optimization Effects

Туре	Optimizat ion	F _{max} (MHz)	Run Time (ms)	Power Dissipation (Watt)	Power Usage (J)
MT	None	277.2	16574	12.01	199.1
Pipeline	None	243.4	117523	10.59	1245.2
MT	Basic	194.7	2445	16.94	41.4
Pipeline	Basic	249.1	116457	9.93	1156.7
Pipeline	Advanced	148.0	251	15.44	3.8

Sliding window is 66x faster than baseline

Stratix V Speed-up vs. Baseline

- Up to 133x speed-up
- CFD speed-up is minimal due to lack of area for optimization



FPGA Speed-up vs. Baseline

Stratix V vs. E5-2670 & K20c

- CFD only has basic optimization
 - Performance is not good
- Stratix V:
 - Beats CPU performance in all benchmarks except Hotspot
 - Beats CPU power efficiency in all benchmarks
 - Cannot beat GPU performance in any benchmark
 - Beats GPU power efficiency in all benchmarks up to 3.4x



Portfolio of Accelerators – What are they good for – (1) general purpose

- Vector (SIMD) HPC accelerators => General purpose HPC esp. SIMD
 - GPUs, Xeon Phi, Shenwei SW26010 , Post-K ARM-SVE
- Macro Dataflow Processors => asynch threads, functional programming etc.
 - ETL EM-4, Wave computing, ...
- High memory bandwidth accelerators (FLOPS increase = Moore's law ending) => memory capacity and BW increase, =>
 - NEC SX-Aurora Vector processor BYTES/FLOPS ~= 1
 - Future 3-D die stacked architectures, NVM including NV-DIMM
- FPGAs programmable HLS languages e.g. OpenCL
 - Intel Stratix 10, Xilinx Vertex 9
- Superconducting Accelerators
 - Massive single thread performance to accelerate serial bottleneck

Portfolio of Accelerators - – What are they good for – (2) ML / Al

- DNN accelerators accelerating tensor operations => many ML apps
 - Small Matrix-Tensor Engine (NVIDIA Volta TensorCore, Intel LakeCrest?)
 - Systoric Array (Google TPU2)
 - Small Vector processor arrays (Fujitsu DLU)
- Neuromorphic accelerators –Spiking Neural Networks => brain simulation, more power efficient ML?
 - IBM TrueNorth
 - Manchester U Spinnaker
 - Heidelberg U BrainScales
 - Many others (DoE, U-Tokyo NEC, SIngapore, ...)
- Symbolic Computing accelerators => string searches. E.g. Genomics?
 - Micron Automata Processor

Portfolio of Accelerators – – What are they good for – (3) Quantum and Pseudoquantum

- Quantum Annealers (theory invented at Tokyo Tech.) => ML?
 - D-Wave
 - Others in the lab
- Pseudoquantum (CMOS) Annealers => ML?
 - Hitachi Ising chip (ISSCC 2015)
 - Fujitsu (Pseudo)quantum annealing chip
- Quantum Gate processors => Quantum simulation, Cryptography
 - Many ongoing work

I/O and Data accelerators – What are they good for (4) I/O

- Cray MTA graph operations
- Burst Buffers I/O intensive ops e.g. checkpoints
- Database accelerators classic

Non-Volatile Memory and 3-D Stacking

- Many devices
- Various stacking technologies

- Results: Massive capacity, extreme bandwidth, low power
- Exploits Z-direction locality
- New breed of "in memory computing"
- Could persist as a trajectory for the next 20 years



When does data movement dominate?

(Original Slide Courtecy John Shalf@LBNL)



0.5mm 1.2 mm

Core Energy/Area est.

 Area:
 12.25 mm²

 Power:
 2.5W

 Clock:
 2.4 GHz

 E/op:
 651 pj

Area: 0.6 mm² Power: 0.3W (<0.2W) Clock: 1.3 GHz E/op: 150 (75) pj

Area: 0.046 mm² Power: 0.025W Clock: 1.0 GHz E/op: 22 pj Data Movement Cost

Compute Op == data movement Energy @ **108mm** Energy Ratio for 20mm **0.2x**

Compute Op == data movement Energy @ **12mm** Energy Ratio for 20mm **1.6x**

Compute Op == data movement Energy @ **3.6mm** Energy Ratio for 20mm **5.5x** Could be reduced by orders of magnitude by 3D, as Z-direction movement is under 1mm

Capacity by dense NVM w/DRAM cache







tensilica

Example Innovation: Tungsten TSV at 2um ultra fine pitch with die thinning by Tezzaron Semiconductor

- Suppose 4TF SFP @ 7nm, 16TB/s internal chip BW vs. 200GB/s external chip mem BW => 80 times speedup!
- High-density, high-signaling TSV challenge
 - Wide I/O 2 1024 bits 1 Ghz -> 2~3 Ghz
 - We need 128,000 bits @ 1Ghz !
 - 10 micron TSV estimation
 - 400 x 400 TSVs on 20mx20m chip -> 50 micron spacing
 - With tungsten TSVs the chip area is negligible

Many-layer stacking via aggressive wafer thinning and selfdiagnostics

14.2 mm



0.5mm

4x nm node

40 nm node

1/0 layer

Finished thickness:

65 nm node

Source: Tezzaron website http://www.tezzaron.com

DiRAM4 Stack Overview

(Tezzaron slides taken from http://www.tezzaron.com/media/Tezzaron-Presentation-EPS-100814-dist-.pptx)

- **64 Gb** of Memory in 175 mm²
- 256 fully independent RAMs
- 16 Banks per RAM
- 64 bit Sep I/O Data per RAM
- 7ns Access Time (Closed page to data)
- **12ns tRC** (Page Open to Page Open in a Bank)
- 16 Tb/s Data Bandwidth
- Competitive Manufacturing Cost





Tezaron semiconductor

Tezzaron Semiconductor

Super Building Block Archiecture (Amano, Keio U)



System View on "Post-Moore" Architecture

Not just a new device, but focus on how they are interconnected, and integrated as a system controlling their power

A Hub architecture that employs Inductive (3D) TCI and programmable FPGA+Switch



ドーターチップ接続のイメージ



Microphotograph of stacked test chips.

Strawman BYTES-Oriented Post-Moore Architecture

Low voltage & power CPU for direct stacking and large silicon area

Domain-specific hetero- and customizable processor configurations, including PIM

Extreme multi-layer DRAM & NVRAM stacking via high density tungsten TSV

Direct WDM optics onto Interposer



Direct Chip-Chip Interconnect with DWDM optics Low Power Processor allows Direct 3D Stacking Configurable Low-power CPU



Fujitsu SPARC64[™] Xifx (2015)



x10 ports

Interconnect Shortcomings

- Current technology:
 - 10\$ / Gbps and 50 pJ per bit, per link
 - 1 exaflops -> 10 PB/s injection bw
 - O (1B\$) and O(5MW) (node link only)
- First stepping stone: mid-board optics vcsels
 - Advanced Development program at HPE
 - Cheaper, more efficient, can be water-cooled
- Exascale technology target: silicon photonics ring resonators
 - 10 cents per Gbps, 5 pJ per bit





Separate 1xN_Bottom Emitting VCSEL_Arrays w/ integrated lens







12dB

1298

Optical Network Technology for Future Supercomputers & IDC



Luxtera 2.5D Photonic Data Pump



- 2.5pJ/bit power
- Bare metal protocol
 - Ultra low latency
 - Protocol agnostic
- 8 core Fiber
- 25Gb SERDES or 3.125Gb interface
- Self-calibrating self-tuning
- >1.6Tb/s payload

(Tezzaron slides taken from
http://www.tezzaron.com/media/TezzaronPresentation-EPS-100814-dist-.pptx)

Тегатоп

Tezzaron Semiconductor

6A560

09648

32 x 32 Optical Circuit Switch (Courtesy NTTAIST)

After FC bonding

- Ceramic LGA interposer with 0.5-mm pitch
- Flip-chip bonding with Au bumps and non-conductive paste
- LGA socket to contact PCB

LGA socket





Problem: heavy optical loss

VICTORIES

http://www.aist.go.jp

Fast Optical Crossbar Swtch (EECS, UCB)

Seok et. al. "Large-scale broadband digital silicon photonic switches with vertical adiabatic couplers" Optica, 3-1, 2016

 $\mathbf{\nabla}$



Fig. 1. Schematics of silicon photonic MEMS switches. (a) Matrix architecture of silicon photonic MEMS switch, (b) close-up view of a MEMS-actuated adiabatic coupler, (c) switch unit cell in the OFF state, and (d) switch unit cell in the ON state. The adiabatic coupler is precisely positioned at the optimum distance to the bus waveguide.

- Array of 64x64 MEMS optical crossbar switch
- 3.7db on-chip insertion loss
- 0.91microsecond switching time
- At 100,000 ports 9 hop network
 - 33db+ loss
 - 8.2 microsecond switching time => 1Tb 800Kbyte BW x Delay



Solution: Hybrid EO Network

- Idea1: use low (latency/diameter, bandwidth, power) electrical network for low latency messages, and use optical circuits for high bandwidth and fixed topology messages
- Idea2: merge the electronic switch and optical MEMS switch, and use the latter as the control plane of the optical MEMs circuit
 - Thus the electronic switches become the optical speculative "buffer"

Hybrid Electro-Optical Network w/shortcuts [Takizawa&Matsuoka LSPP07]

"Locality Aware MPI Communication on a Commodity Opto-Electronic Hybrid Network"



NICT Optical Packet Switch Node (Slides courtesy NICT)



- 4 x 4 OPS node with optical packet (OP) transponder
- 100Gb/s OPS port, 10GbE x 10 Client ports
- Stability: Tolerance for environmental disturbance (Polarization, Power fluctuation)
- Total throughput : 800 Gb/s
- Total power consumption: 141 W (w/o Transponder)
- 10-node hopping, 450 km fiber transmission



vitching 2015. H. Furukawa, et.al, no.P.4.16, ECOC2015.

70

Applications & Algorithms

Slides by Kengo Nakajima Information Technology Center The University of Tokyo

New Frontiers of Computer & Computational Science towards Post Moore Era December 22, 2015, Takeda Auditorium, The University of Tokyo

Assumptions & Expectations towards Post-Moore Era

- Higher Bandwidth, Larger & Heterogeneous Latency
 - Memory: 3D Stacked Memory
 - Network: Optical Communication
 - Both of Memory & Network will be more hierarchical
- Larger Size of Memory & Cache
- Transaction/Transactional Memory
- Application-Customized Hardware, FPGA
- Large Number of Nodes/Number of Cores per Node
 - under certain constraints (e.g. power, space ...)
Applications & Algorithms in Post-Moore Era (1/2)

- Compute Intensity ⇒ Data Movement Intensity
 - Non-Blocking Method, Out-of-Core Algorithm
- Implicit scheme strikes back !
 - I believe it was never defeated
 - Improvement of performance on sparse matrix computations
 - Big change and advancement are expected in all research areas related to algorithms for sparse matrices including preconditioning
 - Everything might be easier... but don't relax too much!
 - Other Compute to Data Algorithms: H-Matrices

Highly-Scalable Atmospheric Simulation Framework (ACM Gordon Bell Prize 2016)



Slide courtesy Haohuan Fu



Weak-scaling results



国家超级计算无锡中则 National Supercomputing Center in Wuxi

Ę

The 488-m res run: 0.07 SYPD, 10.6M cores, dt=240s, 89.5X speedup over explicit

GeoFEM Benchmark: ICCG for FEM Performance of a Node: Flat MPI

	SR11K/J2 Power5+	T2K AMD	FX10	K	Earth Sim 1
Core #/Node	16	16	16	8	8
Peak Performance (GFLOPS)	147.2	147.2	236.5	128.0	64.0
STREAM Triad (GB/s)	101.0	20.0	64.7	43.3	256.0
B/F	0.686	0.136	0.274	0.338	4.00
GeoFEM (GFLOPS)	19.0	4.69	16.0	11.0	25.6
% to Peak	12.9	3.18	6.77	8.59	40.0
LLC/core (MB)	18.0	2.00	0.75	0.75	-

Sparse Solver: Memory-Bound

Improvement of performance on sparse matrix computations due to higher memory bandwidth



Assumptions & Expectations towards Post-K/Post-Moore Era

- Post-K (-2020)
 - Memory Wall
 - Hierarchical Memory (e.g. KNL: MCDRAM-DDR)
- Post-Moore (-2025? -2029?)
 - Larger Size of Memory & Cache
 - Higher Bandwidth, Larger & Heterogeneous Latency
 - 3D Stacked Memory, Optical Network
 - Both of Memory & Network will be more hierarchical
 - Application-Customized Hardware, FPGA
- Common Issues
 - Hierarchy, Latency (Memory, Network etc.)
 - Large Number of Nodes/Number of Cores per Node
 - under certain constraints (e.g. power, space ...)

Parallel-in-Space/Time (PiST)

- MG is scalable, but improvement of performance is limited by parallelization only in space direction
- Time-Dependent Problems: Concurrency in Time Dir.
- Multigrid in (Space+Time) Direction
 - ✓ Traditional time-dependent method: Point-Wise Gauss Seidel
 - ✓ XBraid : Lawrence Livermore National Laboratory
 □Application to nonlinear problems (Transient Navier-Stokes Eqn's)
- MS with 3 sessions in SIAM PP16 (April 2016)
- PiST approach is suitable for the Post-Moore Systems with a complex and deeply hierarchical network



Co-Designing Post-Moore HPC System Architecture



FLOPS-Oriented => BYTES-Oriented Numerical Applications and Algorithms



Programming Models and Abstraction?ming & Comm MW for

Exabit Optical Inteconnect

Perf Modeling of **Super Building Block Architecture** System SW and & Programming of **Deep and High** Bandwidth Memory Hiearchy

Next Gen Exabit-class **Optically Switched** Interconnect

2.5D-CPU Card Wavelength Single-source Bank Vavelengths MUX Fiber Waveleng**Suppty** (Optical comb) ilicon Photonics Integratio

Specialized/Integrated/Re configurable Super **Building Block**

000000

Architecture

Advanced 3-D stacked Non-Volatile Memory >Tbytes, >10TByte/s

00000

Optica

SV Interpose

 \bigcirc \bigcirc \bigcirc

Post Moore Era Supercomputing Workshop @ SC17

- https://sites.google.com/site/2017pmes/
- Jeff Vetter (ORNL), Satoshi Matsuoka (Tokyo Tech) et. al.



Search this sit

2016 Post-Moore's Era Supercomputing (PMES) Workshop Home

News

Call For Position Papers - Submission Deadline - June 17 Invited Speakers Photos Program Resources Workshop Venue Sitemap

2016 Post-Moore's Era Supercomputing (PMES) Workshop Home

Co-located with <u>SC16</u> in Salt Lake City Monday, 14 November 2016

Workshop URL: <u>http://j.mp/pmes2016</u> CFP URL: <u>http://j.mp/pmes2016cfp</u> Submission URL (EasyChair): <u>http://j.mp/pmes2016submissions</u> Submission questions: <u>pmes16@easychair.org</u>

tp://j.mp/pmes2016 mp/pmes2016cfp ttp://i.mp/pmes2016submissions

Submission questions. <u>pmesro@easychair.org</u>

This interdisciplinary workshop is organized to explore the scientific issues,

challenges, and opportunities for supercomputing beyond the scaling limits of

News PMES Submission Site Now Open! PMES Workshop Confirmed for SC16! Submissions open for PMES Position Papers on April 17

Important Dates

Submission Site Opens: 17 April 2016

Output a final Decalling of 47 June 2040