





Outline

- Power and energy consumption basics
- Power consumption in processors
- Multicore: power and utilization walls
- Energy advantages of hardware accelerators
- Playing with accuracy for reducing energy
- Towards heterogeneous manycores
 - Can 3D stacking help?
 - Optical interconnect



_				
Гуре	Model	Mhash/s	Mhash/J	Power (W)
GPP	Intel Xeon X5355 (dual)	22.76	0.09	120
GPP	ARMCortex-A9	0.57	1.14	1.5
GPP	Intel Core i7 3930k	66.6	0.51	130
GPU	AMD 7970x3	2050	2.41	850
GPU	Nvidia GTX460	158	0.66	240
ASIC	AntMiner S1	180.000	500	360
ASIC	AntMiner S5	1.155.000	1957	590
FPGA	Bitcoin Dominator X5000	100	14.7	6.8
FPGA	Butterflylabs Mini Rig	25.200	20.16	1250

Time has Come for Specialization

- Microsoft Unveils Catapult to Accelerate Bing!
 - One FPGA per blade
 - 6 × 8 2-D torus topology
 - High-end Stratix V FPGAs
- Running Bing Kernels for feature extraction and machine learning
- Increase ranking throughput by 95% at comparable latency to software-only
- Increase power consumption by 10%
- Increase total cost of ownership by less than 30%





Outline

- Power and energy consumption basics
- Power consumption in processors
- Multicore: power and utilization walls
- Energy advantages of hardware accelerators
- Playing with accuracy for reducing energy
- Towards heterogeneous manycores
 - Can 3D stacking help?
 - Optical interconnect





- Relies on the ability of many applications to tolerate some loss of quality
- Allows substantially improved energy efficiency by relaxing the need for fully precise operations
 - Trade quality against performance/energy





Intrinsic Application Resilience

• e.g. Image Segmentation





Algebraic approximations

- Real digital operators are not equal to mathematical operators
 - E.g. 1/sqrt(x) can be approximated by various techniques or polynomial orders
- Applications often accept even more aggressive approximations



Approximate Operators • Pruning logic cells - Han-Carlson and Kogge-Stone adders Ranking metrics on logic gates ⊖ - Area ⊟ - Delay 13 Energy Activi 12 Energ -Delay Produc Normalized Gains (Conventional/Proposed) 1/22 11 10 1/2 g 64-bit Kogge Stone Adde $1/2^{4}$ 1/29. 15:014:013:012:011:010:0 9:0 8:0 7:0 6:0 5:0 4:0 3:0 2:0 1:0 0:0 29 28 22 215 214 213 212 211 210 2 2 2 23 21 2 -Output Significance (a) Ranking Scheme for Minimizing Rel. Error and Average Error 10⁻³ 10⁻² 10⁻¹ 10¹ 10 10 Relative Error Magnitude % A. Lingamneni, C. Enz, J.-L. Nagel, K. Palem and C. Piguet, Energy Parsimonious 48 Circuit Design through Probabilistic Pruning, DATE, 2011 (CF12, ACM TECS 2011)

Approximate Adders (CSEM/EPFL/RICE)

• Results after an FFT operation





Hard and Soft SIMD

Hard SIMD

Needs HW modifications in the operators (e.g., manipulate the carry propagation of an adder)

Very limited number of vector configurations (e.g., 1x64b, 2x32b, 4x16b, 8x8b)



Soft SIMD

The operator is unaware of the vector configuration

Large number of vector configurations (e.g., 5x12b, 6x10b, 1x32b+2x16b, etc.)



Outline

- Power and energy consumption basics
- Power consumption in processors
- Multicore: power and utilization walls
- Energy advantages of hardware accelerators
- Playing with accuracy for reducing energy
- Towards heterogeneous manycores
 - Can 3D stacking help?
 - Optical interconnect

Chips go 3D!

- 3D Integrated Circuits
 Stack Multiple Dies
- Wire Length Reduction
 - Replace long, high capacitance wires by Through Silicon Vias (TSVs)
 - Low latency, low energy, high bandwidth
- Heterogeneous Integration
 - Image Sensors, Sensor Network Nodes
 - Processor + Memory



TSVs







Wide I/O, Hybrid Memory Cube (HMC)

- Wide I/O memory interface
- 3D die stacks with TSVs
- 2.5D interposer
- 50GB/sec

- Micron/Intel's HMC couples a logic layer with 3D-stacked DRAM
- 160GB/sec







On-chip Optical Interconnects





Towards Heterogeneous Multicores

- Different cores on a single chip
 - GPPs, HW accelerators, memory, network-on-chip
- Self-adapting micro-architectures
 - Dynamically adapt the hardware to the application
 - Provide support for static and dynamic compilation



Towards Heterogeneous Multicores

• Embedded and High-Performance Computing



Embedded heterogeneous multicore



Heterogeneous platforms



FPGA accelerators for HPC

• C to hardware high-level synthesis boosts hardware designer productivity

Conclusions

- Energy consumption is a major issue
 - True in embedded systems since 20 years...
 - But also true now in HPC, mobile clouds, data centers, etc.
- What did we learn?
 - Hardware specialization
 - Inexact computation
 - Emerging technologies
- Dark Silicon is an opportunity
 - Heterogeneous manycore architectures
 - Bring a new demand for genuinely high level synthesis tools and (JIT) compilers that map programs to accelerators

63

References

- Hadi Esmaeilzadeh, Emily Blem, Renee St. Amant, Karthikeyan Sankaralingam, and Doug Burger. 2011. Dark silicon and the end of multicore scaling. SIGARCH Comput. Archit. News 39, 3 (June 2011), 365-376.
- Venkatesh, Ganesh, Sampson, Jack, Goulding, Nathan, Garcia, Saturnino, Bryksin, Vladyslav, Lugo-Martinez, Jose, Swanson, Steven, and Taylor, Michael Bedford, Conservation cores: reducing the energy of mature computations, Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2010.
- William J. Dally, James Balfour, David Black-Shaffer, James Chen, R. Curtis Harting, Vishal Parikh, Jongsoo Park, David Sheffield "Efficient Embedded Computing" IEEE Computer, July 2008.
- 4. William J. Dally, keynote at IPDPS2011.
- Nathan Goulding, Jack Sampson, Ganesh Venkatesh, Saturnino Garcia, Joe Auricchio, Jonathan Babb, Michael Bedford Taylor, and Steven Swanson, GreenDroid: A Mobile Application Processor for a Future of Dark Silicon, Hot Chips 22, Stanford, CA, Aug. 2010.
- Paul Chow, Why Put FPGAs in Your CPU Socket?, keynote at FPT 2013. http://www.fpt2013.org/Day3_keynote.pdf
- Jongpil Jung, Kyungsu Kang and Chong-Min Kyung, Design and management of 3d-stacked NUCA cache for chip multiprocessors. In Proc. of the 21st edition of the great lakes symposium on Great lakes symposium on VLSI, GLSVLSI '11, pages 91–96, 2011. ISBN 978-1-4503-0667-6.
- Denis Dutoit, Eric Guthmuller and Ivan Miro-Panades, 3d integration for power- efficient computing. In Design, Automation Test in Europe Conference Exhibition (DATE), 2013, mars 2013.
- Taeho Kgil, Shaun D'Souza, Ali Saidi, Nathan Binkert, Ronald Dreslinski, Trevor Mudge, Steven Reinhardt and Krisztian Flautner, PicoServer: using 3D stacking technology to enable a compact energy efficient chip multiprocessor. SIGOPS Oper Syst. Rev., 40(5):117–128, octobre 2006. ISSN 0163-5980.

References

- D.Menard, R.Rocher and O.Sentieys. Analytical Fixed-Point Accuracy Evaluation in Linear Time-Invariant Systems. IEEE Transactions on Circuits and Systems I: Regular Papers, 55(10):3197–3208, November 2008.
- K. Parashar, R. Rocher, D. Menard, O. Sentieys, D. Novo, and F. Catthoor. Fast performance evaluation of fixed-point systems with un-smooth operators. In Proc. of the IEEE/ACM International Conference on Computer-Aided Design (ICCAD), pages 9–16, San Jose, CA, November 2010.
- H.-N. Nguyen, D. Menard, and O. Sentieys. Dynamic precision scaling for low power wcdma receiver. In Proc. of the IEEE International Symposium on Circuits and Systems (ISCAS), pages 205–208, Taipei, Taiwan, May 2009.
- 4. <u>http://blogs.msdn.com/b/satnam_singh/archive/2011/01/18/reconfigurable-data-processing-forclouds.aspx</u>
- 5. http://research.microsoft.com/en-us/projects/kiwi/default.aspx
- David Greaves and Satnam Singh, Designing Application Specific Grouits with Concurrent C# Programs, in ACM/IEEE International Conference on Formal Methods and Models for Codesign, IEEE, 26 July 2010
- 7. <u>http://researcher.watson.ibm.com/researcher/view_group.php?id=122</u>
- 8. http://queue.acm.org/detail.cfm?id=2000516
- 9.
 Satnam Singh, Computing without Processors. Queue 9, 6, Pages 50 (June 2011), 14

 pages.
 http://doi.acm.org/10.1145/1989748.2000516
- 10. Shekhar Borkar, Andrew A. Chien, The Future of Microprocessors, Communications of the ACM, Vol. 54 No. 5, Pages 67-77, 10.1145/1941487.1941507