

Data Centers

Current trends and impacts

Xavier VIGOUROUX HOSCAR Workshop Bordeaux Sept. 2, 2013







A VERY QUICK SUMMARY



Article suggestion: http://herbsutter.com/welcome-to-the-jungle/





"THE SHIFT IS ALSO AN INDICATION THAT THE INTEL PENTIUM 4 LINE HAD REACHED ITS PERFORMANCE LIMIT FOR THE AMOUNT OF POWER IT CONSUMED."

N. BROOKWOOD - ANALYST



2004 – EARLY ADOPTERS FOR GPGPU



2001 – programmable shader (NVidia) 2002 – name "GPGPU" by Mark Harris 2003 – cloud simulation on GPGPU by Mark Harris

Source: http://gpgpu.org/s2004, http://gpgpu.org/oldsite/data/history.shtml, http://gpgpu.org/about,



2005 – TANSTAAFL* ... BUT IT'S OVER



Ref: "The Free Lunch Is Over - A Fundamental Turn Toward Concurrency in Software" By Herb Sutter in Dr. Dobb's Journal * "There ain't no such thing as a free lunch"





And now ?



Land & Ocean Temperature Percentiles Jul 2013

NOAA's National Climatic Data Center

Data Source: MLOST version 3.5.4



Source: http://www.ncdc.noaa.gov/sotc/global/





New constraints On most consuming parts









12kW/rack



6°C – 30kW



30°C – 80kW/Rack Bull

DIRECT LIQUID COOLING

- □ Water has closed as possible to the heat source (CPU at 70 C)
- \Box Water can be hotter (as ΔT is key)
- Room can be hotter (remove CRAC)
- And without any change in maintenance process
 - CPU can be changed,
 - DIMM can be changed
 - Blades can be removed

PUE < 1.08



DIRECT LIQUID COOLING



Most efficient pure x86 supercomputer in Green 500 june' 13 (météo France)
 New version with GPGPU is coming



ADIABATIC COOLING





Cooling capacity	108 kW
Power consumption	2,8 kW
Water used (estim. h/year)	472
Water used (estim. m ³)	53







New constraints On most consuming parts





FREQUENCY AND POWER CONSUMPTION

IVB Power vs. Voltage -1.6 GHz and 2.4 GHz, 90°C



Power vs. Frequency, 1.26 V, variable temperature



IVB Voltage-Frequency Shmoo plot







CPU FREQUENCY







[fr q] x [#flop cycle] x [#cc es]













	SEVENTH FR.	AMEWOI	RK	
	Proposal acronym: Proposal full title:	font-Blanc,	Mont-Bland European scalab C platform based	c ¹ ole and power on low-power
Alex Ramin Guadalupe List of parti	ez (Technical Manager) Moreno (Project Manager) cipants:			
 Alex Ramin Guadalupe List of participant no. 	ez (Technical Manager) Moreno (Project Manager) cipants: Participant organisation n Barrolona Supercomputing Center	ame	Part. short name	Country
 Alex Ramin Guadalupe List of participant no. 1 	ez (Technical Manager) Moreno (Project Manager) ipants: Participant organisation n Barcelona Supercomputing Center Bull SAS	ате	Part. short name BSC Bull	Country Spain France
 Alex Ramin Guadalupe List of parti Participant 1 2 3 	ez (Technical Manager) Moreno (Project Manager) cipants: Participant organisation n Barcelona Supercomputing Center Bull SAS ARM Limited	ame	Part. short name BSC Bull ARM	Country Spain France UK
Alex Ramin Guadalupe List of parti Participant no. 1 2 3 4	ez (Technical Manager) Moreno (Project Manager) :ipants: Participant organisation n Barcelona Supercomputing Center Bull SAS ARM Limited Gnodal Ltd.	ame	Part. short name BSC Bull ARM Gnodal	Country Spain France UK UK
 Alex Ramin Guadalupe List of participant no. 1 2 3 4 5 	ez (Technical Manager) Moreno (Project Manager) ipants: Participant organisation n Barcelona Supercomputing Center Bull SAS AM Limited Gnodal Ltd. Forschungszentrum Jülich GmbH	ame	Part. short name BSC Bull ARM Gnodal JUELICH	Country Spain France UK UK Germany
 Alex Ramin Guadalupe List of parti Participant no. 1 2 3 4 5 6 	ez (Technical Manager) Moreno (Project Manager) ipants: Participant organisation n Barcelona Supercomputing Center Bull SAS ARM Limited Gnodal Ltd. Forschungszentrum Jülich GmbH Leibbiz-Rechenzentrum der Bayeri Akademie der Wissenschäfter	ame	Part. short name BSC Bull ARM Gnodal JUELICH BADW-LRZ	Country Spain France UK UK Germany Germany
Alex Ramin Guadalupe List of participant no. 1 2 3 4 5 6 7 7	ez (Technical Manager) Moreno (Project Manager) cipants: Participant organisation n Barcelona Supercomputing Center Bull SAS ARM Limited Gnodal Ltd. Forschungszentrum Jülich GmbH Leibniz-Recherzentrum der Bayeris Akademie der Wissenschaften Grand Equipement National de Cal	ame schen cul Intensif	Part. short name BSC Bull ARM Goodal JUELICH BADW-LRZ GENCI	Country Spain France UK UK Germany Germany France
Alex Ramin Guadalupe List of parti Participant no. 1 2 3 4 5 6 7 8	ez (Technical Manager) Moreno (Project Manager) :jpants: Participant organisation n Barcelona Supercomputing Center Bull SAS ARM Limited Gnodal Ltd. Forschungszentrum Jülich GmbH Leibniz-Rachenzentrum der Bayeri Akademie der Wissenschäften Grand Equipement National de Cal Consorzio Interuniversitario CINECO	ame schen cul Intensif	Part. short name BSC Bull ARM Gondal JUELICH BADW-LRZ GENC1 CINECA	Country Spain France UK UK Germany Germany France Italy

- o ARM sells IP not a physical core
- o ARM is present in
 - o Mobile (phone, tablet, laptop)
 - o Embedded (electronics)
 - o Enterprise (switch, disk, printers, servers...)
 - o Home (TV, camera, games consoles)
- o Shipped chips
 - o 2012 27 billions
 - o 2017 41 billions

Why ARM cores ?

- o Price (\$1 \$25)
- o Power consumption is key for mobile device (100mW)
- o Very versatile

HOW TO SIZE A MACHINE

[Money] => [CAPEX] + [OPEX]

[CAPEX]

Flops/\$ is decreasing Cheaper More flops More Watt

cpu.h



More Watt (from CAPEX and PUE) \$/Watt is increasing Much More Expensive

W.h







Soft dev & Users	 Write code taking into account "execution policy" Execute your job with a choice of policy (time to result, Joule to result, Price to result) Pay less by a cossing the policy (time to result, Joule to result)
Midleware providers	 You will have to bridge the gap between Hardware and Users Write code taking into account "execution policy" (upward and downward) You will face heterogeneity
Cloud And Data Center Director	 You will charge in Watt.hour You can adjust your batch scheduler according to policy (green energy, power capping,) Efficient hardware will be valuable (DLC,) You will adjust your electricity contract



WHAT IS TO IMPROVE





Pluggins: RAPL, IPMI (OS) and RRD
Per job (global value & time slice)
Per node

- O Per user
- New srun parameter to allow CPU frequency scaling for job execution







High Definition Energy Efficiency Monitoring









POLICIES COMPARISON





SESAMES FRAMEWORK: A MULTI-CRITERIA GREEN SCHEDULER FOR CONSUMING LESS AND BETTER



M.E.M. Diouri, O. Glück and L. Lefèvre,

SESAMES: A Smart Grid based framework for consuming less and better in extreme scale infrastructures, GreenCom 2013, 20-23 August, Beijing 2013.



COMPUTE PERFORMANCE (flops/s)

- Applications are written for <u>compute</u> <u>performance</u>
- Tools are seeking performance issues (host spot, idle time)
- Cost of a run is expressed in core.hour (because it's easy to measure)

POWER CONSUMPTION (W.h)

- Becomes an increasing budget
- Will be charged on end users
- Is a strict constraint for large systems
- flops/Watt is becoming a metric (see green500).

HDEEM MERGES THE BOTH WORLDS

Precise metrics in performance AND power consumption



TU DRESDEN IN A NUTSHELL



The technische universität dresden is the largest institute of higher education in the city of dresden, the largest university in saxony and one of the 10 largest universities in germany.

#students	34'993
#foreign students	3′442 (9,8%)
#employees	6′ 123
#faculties	14
Total Budget	€ 500M (circa)





INTERESTED IN HDEEM



Developping tools

Ŷ



SOLUTION OVERVIEW





DON'T FORGET ABOUT THE JUNGLE







center for excellence in parallel programming









Architect of an Open World™



Production: 1908 – 1927 50% of the world market Low cost (\$850 – 4 months) Easy Maintenance 93 min to build one



IF ALL YOU HAVE IS A HAMMER...











