Green Energy Prediction for Distributed Datacenters

Prof. Tajana Šimunić Rosing Dept. of Computer Science

System Energy Efficiency Lab seelab.ucsd.edu









see

Semiconductor Research Corporatior

Energy efficiency at global scale See

NSF GreenLight & IRNC: TransLight/StarLight





Key issues for distributed renewable-powered datacenters



- Green energy availability varies dramatically
 - Instantaneous use leads to significant energy efficiency losses
 - Prediction is needed
- Datacenter computing requires consistent performance
 - Infrastructure that monitors and manages computation in datacenters has to be aware of performance costs
 - Service response times are around 100ms, Max 10% batch job throughput hit
- Energy costs of datacenters are typically higher than green energy availability
 - Brown energy needs to be present to both supplement green and as "insurance" to meet performance constraints
 - Improvements in computation & networking infrastructure energy efficiency are necessary (power, thermal and cooling management)



Energy efficiency of the infrastructure



NSF Project GreenLight

- Green cyber-infrastructure in energy-efficient mobile facilities
- Closed-loop power and thermal management

Dynamic power management (DPM)

- HW level: adaptive power gating gives 40% energy savings with no perf. impact
- SW level: 92% reduction in performance variability with DVFS
- Optimal DPM for a class of workloads
- Machine learning to adapt
 - Select among specialized policies
 - Measured energy savings of 70%

Dynamic thermal management (DTM)

- Workload scheduling:
 - Machine learning for dynamic adaptation
- Proactive thermal management
 - Reduces thermal hot spots by average 80% with no performance overhead
- Cooling aware management
 - Savings of 70% in cooling subsystem













Texas Instruments Inc.

NSF GreenLight:Image: Optimized sectorDashboard & History plotsSector

GLIMPSE

GreenLight Project

Water Temperature Temperature Temperature Total Power Measured: 22011 W Used: 28660 W PUE: 1.30 Note: * this is a rough estimation

 Multiple sensor data: temperature, fan speed, liquid flow rate & temp, power

Heat Xcha 4

Heat Xchg 5

Heat Xcho 6

 Use measurements to develop models needed for energy management

7	Heat	Xchg 8		Ma Island I.					
4	bow i	Arciage 1074		and a standard	alalah dalah d	al al ada da da da da da		a talaha baha baha baha baha baha baha baha	de la
а. З	50W 🗖	Eab 26 16:00	Eab 27.0-00	Eab 27 8:00	Eab 27 16:00	Eab 28 0:00	Eab 28 8-00	Eab 28 16:00	
1		reb 26 16.00	160 27 0.00	Feb 27 0.00	160 27 10:00	reb 20 0.00	PED 20 0.00	reb 20 10.00	
4	888 🖿	مر مع المناطق المراجع		- Aller aler and a second	and have	And the second s			
		Feb 26 16:00	Feb 27 0:00	Feb 27 8:00	Feb 27 16:00	Feb 28 0:00	Feb 28 8:00	Feb 28 16:00	
2 - 5	75W	thumper-1-ps1							
4	50W	Average = 426	569W						
4	25W	Average - 420.	5071	molence	mundalle	mandra	Manufa	marken	~~
4	oow 🕒				40.00		24.00		
a - 5		17:00	18	5:00	19:00	20:00	21:00	22:00	
4	50W	and a start of the Martha	- Marine and a second	Martin Martine		A CARLENS	lasta ta ta ta da	al a	باساد
2		Feb 26 16:00	Feb 27 0:00	Feb 27 8:00	Feb 27 16:00	Feb 28 0:00	Feb 28 8:00	Feb 28 16:00	
5	20W	thumpor 1 pc2							
2.4	BOW	Average = 447	71.014						
1 4	60W	Average - 447.	710 99	~					
4	40W 🗖								
				Sunday E	obruory 28	2010 0.50.17	DM value in	440.00144	
5	2 50W	1:50 21	1:55	Sunday, F	ebruary 28,	2010 9:59:17	'PM value is	449.00W 22:	:20
	2	1:50 21	1:55	Sunday, F	ebruary 28,	2010 9:59:17	PM value is	5 449.00W 222:	20
u unitere		Feb 26 16:00	1:55 Feb 27 0:00	Sunday, F	ebruary 28, Feb 27 16:00	2010 9:59:17 Feb 28 0:00	PM value is Feb 28 8:00	Feb 28 16:00 22:	20
a contract	2	11:50 21 Feb 26 16:00	1:55 Feb 27 0:00	Sunday, F	February 28, 5	2010 9:59:17 Feb 28 0:00	PM value is Feb 28 8:00	5 449.00W 22: Feb 28 16:00	:20
n - 6	2 88 -	Feb 26 16:00	1:55 Feb 27 0:00	Sunday, F	February 28, 1	2010 9:59:17 Feb 28 0:00	PM value is	5 449.00W 22: Feb 28 16:00	20
- 6i	2 88 -	Feb 26 16:00	Feb 27 0:00	Sunday, F	Feb 27 16:00	2010 9:59:17 Feb 28 0:00	PM value is	Feb 28 16:00	20
n 60	2 88 50 00 00 00 00 00	Feb 26 16:00 Feb 26 16:00 thumper-2-ps1 Average = 536.	Feb 27 0:00	Sunday, F	Feb 27 16:00	2010 9:59:17 Feb 28 0:00	PM value is	Feb 28 16:00	20
- 64 51 51 42	2 5000 5000 5000	Feb 26 16:00 thumper-2-ps1 Average = 536.4	Feb 27 0:00	Sunday, F	Feb 27 16:00	2010 9:59:17 Feb 28 0:00	PM value is	Feb 28 16:00	20
- 61 51 51 42	2 58% = 50W = 50W = 50W =	Feb 26 16:00 Feb 26 16:00 thumper-2-ps1 Average = 536. 7:00	Feb 27 0:00 450W 8:00	: Sunday, F Feb 27 8:00 9:00	Feb 27 16:00	2010 9:59:17 Feb 28 0:00	7 PM value is Feb 28 8:00	Feb 28 16:00	20
6 5 5 4 GOLDA		Feb 26 16:00 thumper-2-ps1 Average = 536. 7:00	Feb 27 0:00 450W 8:00	2 Sunday, F Feb 27 8:00 9:00	Feb 27 16:00	2010 9:59:17 Feb 28 0:00	Feb 28 8:00	Feb 28 16:00	00
6 5 5 4 Grund		Feb 26 16:00 thumper-2-ps1 Average = 536 7:00 Feb 26 16:00	Feb 27 0:00 450W 8:00 Feb 27 0:00	: Sunday, F Feb 27 8:00 9:00 Feb 27 8:00	Feb 27 16:00	Eeb 28 0:00	Feb 28 8:00	Feb 28 16:00	20
		Feb 26 16:00 thumper-2-ps1 Average = 536. 7:00 Feb 26 16:00	Feb 27 0:00 450W 8:00 Feb 27 0:00	: Sunday, F Feb 27 8:00 9:00 Feb 27 8:00	Feb 27 16:00 10:00 Feb 27 16:00	2010 9:59:17 Feb 28 0:00 11:00 Feb 28 0:00	Feb 28 8:00	Feb 28 16:00 0 13: Feb 28 16:00	20
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1		Feb 26 16:00 thumper-2-ps1 Average = 536. 7:00 Feb 26 16:00	Feb 27 0:00 450W 8:00 Feb 27 0:00	2 Sunday, F Feb 27 8:00 9:00 Feb 27 8:00	Feb 27 16:00 10:00 Feb 27 16:00	Eeb 28 0:00	PM value is Feb 28 8:00 12:0 Feb 28 8:00	Feb 28 16:00 Peb 28 16:00 Feb 28 16:00	00
10000000000000000000000000000000000000		11:50 21 Feb 26 16:00 1 thumper-2-ps1 Average = 536. 7:00 1 Feb 26 16:00 1 thumper-2-ps2 1	Feb 27 0:00 450W 8:00 Feb 27 0:00	2 Sunday, F Feb 27 8:00 9:00 Feb 27 8:00	Feb 27 16:00 10:00 Feb 27 16:00	2010 9:59:17 Feb 28 0:00 11:00 Feb 28 0:00	7 PM value is Feb 28 8:00 12:0 Feb 28 8:00	Feb 28 16:00	00
6 5 5 4 6 5 5 6 4 5 5 6 4		11:50 21 Feb 26 16:00 1 thumper-2-ps1 Average = 536,. 7:00 1 Feb 26 16:00 1 thumper-2-ps2 Average = N/A	Feb 27 0:00 450W 8:00 Feb 27 0:00	2 Sunday, F Feb 27 8:00 9:00 Feb 27 8:00	Feb 27 16:00 10:00 Feb 27 16:00	2010 9:59:17 Feb 28 0:00 11:00 Feb 28 0:00	PM value is Feb 28 8:00	Feb 28 16:00	00
10044 6 6 5 5 5 6 4 4 6 5 7 5 5 6 4 4 4 4		11:50 21 Feb 26 16:00 1 thumper-2-ps1 Average = 536, 7:00 1 Feb 26 16:00 1 thumper-2-ps2 Average = N/A	Feb 27 0:00 450W 8:00 Feb 27 0:00	: Sunday, F Feb 27 8:00 9:00 Feb 27 8:00	Feb 27 16:00 10:00 Feb 27 16:00	Eeb 28 0:00	PM value is Feb 28 8:00 12:0 Feb 28 8:00	Feb 28 16:00	00
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5		11:50 21 Feb 26 16:00 1 thumper-2-ps1 Average = 536, 7:00 1 Feb 26 16:00 1 thumper-2-ps2 Average = N/A Feb 26 16:00 1	Feb 27 0:00 450W 8:00 Feb 27 0:00 Feb 27 0:00	: Sunday, F Feb 27 8:00 9:00 Feb 27 8:00 Feb 27 8:00	Feb 27 16:00 Feb 27 16:00 Feb 27 16:00 Feb 27 16:00	Eeb 28 0:00 Feb 28 0:00 Feb 28 0:00 Feb 28 0:00	Feb 28 8:00 Feb 28 8:00 12:0 Feb 28 8:00 Feb 28 8:00	Feb 28 16:00 Feb 28 16:00 Feb 28 16:00 Feb 28 16:00	20
	2 500	11:50 21 Feb 26 16:00 1 thumper-2-ps1 Average = 536. 7:00 5 Feb 26 16:00 1 thumper-2-ps2 Average = N/A Feb 26 16:00 1	Feb 27 0:00 450W 8:00 Feb 27 0:00 Feb 27 0:00 Feb 27 0:00	: Sunday, F Feb 27 8:00 9:00 Feb 27 8:00 Feb 27 8:00 Feb 27 8:00	Feb 27 16:00 Feb 27 16:00 Feb 27 16:00 Feb 27 16:00	Eeb 28 0:00 Feb 28 0:00 Feb 28 0:00 Feb 28 0:00	PM value is Feb 28 8:00 12:0 Feb 28 8:00 Feb 28 8:00	Feb 28 16:00 Feb 28 16:00 Feb 28 16:00 Feb 28 16:00	00

GMVQ VM Power Cost Prediction See

- Goal: Estimate how much a VM consumes and predict what the cost would be if it migrates to another machine
- Approach: Gaussian Mixture Vector Quantization (GMVQ) to fit a GMM to the training data



GMVQ is 3x better than regression

Energy management with virtualization Sec

vGreen



- Scheduling
 - Co-locate guests with orthogonal characteristics
- Management policies
 - · Based on the metrics maintained per guest



VMs: Energy and QoS



Tajana Simunic Rosing

Managing multi-tier applications \mathcal{O}_{SCC}

- What happens when we combine:
 - Latency sensitive jobs (e.g. RUBIS)
 - Throughput sensitive (e.g. batch jobs)
- Preliminary results:
 - More than 10x improvement in SLA with background jobs relative to the default scheduler



RUBiS: auction website





VM Scheduling Policies Throughput

• State of the art:



- Baseline: running services and batch jobs on separate servers
- Selective Consolidation (tChar): vGreen
- Capping (tCap): Cap the CPU allotment to bVM to mitigate interference effects (Padala@EuroSys'09, Nathuji@EuroSys'10)
- Controller:
 - Dynamically control the vCPU allocation of the service VM to maximize the batch job throughput while meeting service response times



Batch job throughput

Our controller is within 7% of baseline; CPU capping has 25% lower average throughput

Energy Efficiency Improvements



Maximize CPUs of various batch jobs while meeting Rubis SLAs



- 70% more efficient than running service & batch VMs separately while within 7% of maximum batch job throughput
- 35% more efficient than the ideal version of state of the art

Green Energy Prediction



- Predict green energy availability for the next 30min window
 - Schedule additional MapReduce jobs accordingly; they take max 30mins



- Data from solar panels at UCSD
- State of the art: exponential weighted average
 - EWMA: 32.6 % error
 - Extended eEWMA: 23.4% error
- Our algorithm:
 - WCMA: less than 9.6% error



- Data gathered from a wind farm in Lake Benton, available by NREL
- State-of-the-art:
 - Integrated predictor: 48.2% error
- Our algorithm: 21.2% error
 - Combination of a weighted nearestneighbor (NN) tables and wind power curve models

Methodology



- Use green energy to schedule "extra" batch jobs.
 - MapReduce (MR) type jobs for this purpose.
 - Initiate more subtasks with the available green energy.
 - Increased throughput
 - Reduced completion time limit reduction to 10%
 - Kill a subtask if the green energy supply level drops



Experimental setup & validation

- Globally distributed datacenters connected
 - 5 datacenters, 12 routers, modeled after ESnet
 - Solar traces from UCSD
 - Wind traces from NREL



Simulator validation	Measured Value	Simulated Value	Ave. Error
Avg. Power Consumption	246 W	251 W	3%
Rubis QoS ratio	0.08	0.085	6%
Avg. MapReduce Completion Time	112 min	121 min	8%

- Jobs run within vGreen VMs on Nehalem server
 - Rubis used for services with 100ms 90th%ile response time constraint
 - MapReduce used for batch jobs with 10% max job completion time reduction (max 5 cores on Nehalem server)
 - Use VM migration (with quantified performance impact) Tajana Simunic Rosing

Benefits of Green Energy Prediction



 Compare our green energy predictive jobs scheduler with instantaneous usage of green energy



Prediction has 93% GE Efficiency

GE Efficiency: ratio of green energy consumed for useful work vs. the total green energy available

Benefits of Green Energy Prediction





Prediction has 22% faster batch job completion time vs. instantaneous

On average, **5x fewer batch tasks need to be terminated** when using GE prediction vs. instantaneous usage

38% more jobs complete with prediction vs. instantaneous

GE Job %: ratio of batch jobs completed with GE over all jobs.

Next steps: Green energy powered global routing



(jointly with Inder Monga, Esnet)



Focus Center Research Program



<u>Sponsors</u>



Intel MICRON Freescale Texas Instruments IBM Xilinx GLOBALFOUNDRIES AMD



Applied Materials Novellus Cadence

DOD & DARPA



Raytheon United Technologies



Multi-Scale Systems Research Center [10 Universities] <u>Director: Prof. Jan Rabaey (UC-Berkeley)</u>. High-level systems design addressing distributed sense and control systems, largescale and small-scale information technologies systems.



Gigascale Systems Research Center [15 Universities] <u>Director: Prof. Sharid Malik (Princeton)</u>. Platform architectures; concurrent systems programming; platform viability; resilient systems and alternative computation models.

Center for Circuit & Systems Solutions [13 Universities] <u>Director: Prof. Larry Pileggi (CMU)</u>. Circuit/module infrastructure; enterprise systems; portable electronics; functional diversity and emerging circuits for post CMOS.

Interconnect Focus Center [13 Universities] <u>Director: Prof. Paul Kohl (Georgia Tech)</u>. Nanoscale electrical & optical interconnects; energy delivery and thermal management; wireless connectivity; modeling, analysis and assessment of new connectivity solutions.

Materials, Structures and Devices [15 Universities] <u>Director: Prof. Dimitri Antoniadis (MIT</u>). Integration of new materials enabling CMOS extension; carbon-based devices; novel embedded memory; functional diversification; theory, modeling and simulation of new devices.



Functional Engineered Nano-Architectonics [14 Univ.] <u>Director: Prof. Kang Wang (UCLA)</u>. Novel materials and processes which enable fabrication of nanoscale devices and interconnects.

MuSyC in a Nutshell

GRAND CHALLENGE : "Energy-smart" distributed systems, that
Are deeply aware of balance between energy availability and demand
Adjust behavior through dynamic and adaptive optimization through all scales of design hierarchy.



19 Faculty Distributed over 9 US Universities

see

- 60 Students (many only partly funded)
- 109 publications
- Monthy e-seminars and bi-annual e-workshops



Multiscale Systems Center: Energy Balanced Datacenters Theme leader: Tajana Simunic Rosing



Realize closed-loop energy management strategies that enable "energy-intensive" large-scale systems to be orders of magnitude more energy-efficient, while ensuring that mission-critical goals are met.



The Multiscale Systems Center

Combined Energy Thermal & Cooling, CETC Management Results



- □ CETC performance overhead < 0.2%
- CETC page migration rate < 5 pages/sec -> negligible overhead & high stability

Average cooling savings of 70% relative to state-of-the-art

Intel Xeon Dual Socket Quad core Server; with state-of-the-art PI fan controller



• **CETC:** Our policy

DLB: Dynamic load balancing (baseline)

- NFMO: Only page migrations allowed
- NMM: No memory clustering
- NCM: No CPU scheduling optimizations

Workload	Socket A	Socket B	Ambient ^o C
W1	3eon	eon + mcf + gcc	42
W2	2eon + mcf	eon + bzip2 + mcf	42
W3	2bzip2 + 2mcf	2bzip2 + 2mcf	42
W4	2perl + 2eon	2gcc + 2mcf	42
W5	2perl + 2bzip2	2gcc + 2mcf	39
W6	2perl + bzip2	gcc+2mcf	42
W7	perl + 3gcc	perl + 2gcc	40
W8	perl + 3gcc	perl + 2gcc	42
W9	3eon	3eon	42
W10	2eon+mcf	2eon+mcf	42
W11	2bzip2	2bzip2	42