

Energy-Aware Dense and Sparse Linear Algebra

Enrique S. Quintana-Ortí





(or... Doing Nothing to Save Energy in Matrix Computations)

Enrique S. Quintana-Ortí





20.1 PFLOPS (10¹⁵ flops/sec.)

2012 DOE/NNSA/LLNL Sequoia

- 10⁹ core level (Power BQC, 1.60 GHz \rightarrow 12.8 GFLOPS)
- 10¹ node level

(16 cores/node)

10⁵ cluster level

(98.304 nodes)







2020 EFLOPS (10¹⁸ flops/sec.)

- 10^{9.5} core level
- 10³ node level!
- 10^{5.5} cluster level



Green500 vs Top500 (November 2012)

Rank	Site	#Cores	MFLOPS/W
Green/Top			
1/253	National Institute for Computational	9.216 (Intel Xeon E5)	2,499.44
	Sciences/University of Tennessee	8.640 (Intel Xeon Phi)	
3/1	DOE/SC/Oak Ridge National	560,640 (AMD Opteron)	2,142.77
		261,632 (NVIDIA K20)	



Green500 vs Top500 (November 2012)

Rank	Site	#Cores	MFLOPS/W	MW to
Green/Top				
1/253	National Institute for	9.216 (Intel Xeon E5)	2,499.44	400.09
	Sciences/University of Tennessee	8.640 (Intel Xeon Phi)		
3/1	DOE/SC/Oak Ridge National	560,640 (AMD Opteron)	2,142.77	466.69
		261,632 (NVIDIA K20)		



Most powerful reactor under construction in France Flamanville (EDF, 2017 for US \$9 billion): 1,630 MWe



Green500 vs Top500 (November 2012)

Rank Green/Top	Site	#Cores	MFLOPS/W	MW to EXAFLOPS?
1/253	National Comput Science Tenness (visit xe.co	315 MEuro/year m for ₤ ;-)	2,499.44	400.09
3/1	DOE/SC/Oak Ridge National Laboratory	560,640 (AMD Opteron) 261,632 (NVIDIA K20)	2,142.77	466.69



Most powerful reactor under construction in France Flamanville (EDF, 2017 for US \$9 billion): 1,630 MWe



Systems ranked #1 in Green500





Systems ranked #1 in Green500





- Reduce energy consumption!
 - Costs over lifetime of an HPC facility often exceed acquisition costs
 - Carbon dioxide is a hazard for health and environment
 - Heat reduces hw. reliability
- Personal view
 - Hardware features energy saving mechanisms
 - Scientific apps. are in general energy-oblivious



- Trading stability for performance
 - Old idea, revisited many times
- Mixed precision
 - Also an old idea (Von Neumann), but interesting
- Probabilistic (or unreliable) systems-on-a-chip
- Energy-aware numerical algorithms

. . .





- Energy-aware numerical algorithms
 - ILUPACK for multicore processors

• The CG method for hybrid CPU-GPU platforms

Basics...



• Remember:

$$E = \int_{0,T} P dt = P_{avg} \cdot T$$



ILUPACK on multicore



Incomplete LU Package (<u>http://ilupack.tu-bs.de</u>)

- Iterative Krylov subspace methods
- Multilevel ILU preconditioners for general/symmetric/Hermitian positive definite systems
- Based on inverse ILUs with control over growth of inverse triangular factors
- Specially competitive for linear systems from 3D PDEs

ILUPACK on multicore Task parallelism



- Multi-threaded parallelism (real s.p.d. systems)
 - Leverage task parallelism



ILUPACK on multicore Task parallelism



- Multi-threaded parallelism (real s.p.d. systems)
 - Out-of-order, dependency-aware execution of tasks





ILUPACK on multicore Task parallelism



- Multi-threaded parallelism (real s.p.d. systems)
 - Out-of-order, dependency-aware execution of tasks





- Dynamic scheduling via (OpenMP) runtime



tasks (no dependencies)



- Sparse linear system benchmark
 - Laplacian equation $-\Delta u = f$ in a 3D unit cube $\Omega = [0,1]^3$
 - Linear system Au = b with A → n x n, n = 252³ ≈ 16 million unknowns and 111 millions of nonzero entries



CPU (performance) P-states:

AMD

- 2 AMD Opteron 6128, 48GB
- DVFS per core

P-state	$V CC_i$	f_i
P_0	1.23	2.00
P_1	1.17	1.50
P_2	1.12	1.20
P_3	1.09	1.00
P_4	1.06	0.80

Intel

- 2 Intel Xeon E5504, 32GB
- DVFS per socket

P-state	VCC_i	f_i
P_0	1.04	2.00
P_1	0.98	1.73
P_2	0.95	1.60
P_3	1.01	1.87





- CPU (power) C-states:
 - C0: normal operation mode
 - C1, C2,...: disable core components (L1/L2 caches), clock signal, memory controller,...
 - Trade off power for wakeup time



UNIVERSITAT JAUME

- National Instruments NI9205+NIcDAQ-9178
- 1,000 samples/sec. per channel





 $P = P^{(S)Y(stem)} + P^{C(PU)} = P^{Y} + P^{S(tatic)} + P^{D(ynamic)}$

- P^{C} is the power dissipated by CPU (socket): $P^{S} + P^{D}$
- P^S is the static power
- *P^D* is dynamic power
- **P**^Y is the power of remaining components (e.g., RAM)

Considerations:

- *P^Y* and *P^S* are constants (though *P^S* grows with temperature)
- Hot system

System power:

 $P = P^Y + P^S + P^D$

Estimated as *idle* power Due to off-chip components: e.g., RAM (only mainboard)

$P^{Y} \approx P^{I} = 80.15 \text{ W}$



• Static power:

 $P = P^{Y} + P^{S} + P^{D}$







JAUME ·

 $P = P^{Y} + P^{S} + P^{D}$





P-state P _i	Vcc _i	f _i	α_i	β_i	ΔP_i^S	ΔP_i^D
P ₀	1.23	2.00	168.59	9.12	_	_
P_1	1.17	1.50	161.10	5.77	-9.52	-32.14
P_2	1.12	1.20	155.90	4.23	-17.09	-50.25
P ₃	1.09	1.00	152.94	3.15	-21.47	-60.73
P ₄	1.06	0.80	150.61	2.44	-25.73	-70.30

• P_i^{S} depends on Vcc_i^2

 $E.g., P_0 (1.23V) \rightarrow P_3 (1.09V): -21.47\%$

- P_i^D depends on $Vcc_i^2 \cdot f_i$ *E.g.*, P_0 (1.23V, 2.00GHz) $\rightarrow P_3$ (1.09V, 1.00GHz): -60.73%
- These values agree within 2.5% with the experimental linear regression models!



P-state P _i	Vcc _i	f _i	α_i	β_i	ΔP_i^S	ΔP_i^D
P_0	1.23	2.00	168.59	9.12	_	—
P_1	1.17	1.50	161.10	5.77	-9.52	-32.14
P_2	1.12	1.20	155.90	4.23	-17.09	-50.25
P_3	1.09	1.00	152.94	3.15	-21.47	-60.73
P ₄	1.06	0.80	150.61	2.44	-25.73	-70.30

- Moving to a higher P-state results in ↓power
- \downarrow Power = \downarrow energy?
 - For a compute-bounded operation, f_i is linear to time⁻¹
 - In principle, for a memory-bounded operation (ILUPACK), reducing f_i should have a minor impact on performance



1st attempt: Dynamic Static voltage-frequency scaling

P-state P _i	T_i	\bar{P}_i^T	Ei	ΔT_i	$\Delta \bar{P}_i^T$	ΔE_i
P_0	34.06	282.87	9,634.78	····	_	-
P_1	43.57	235.64	10,267.72	21.88	-16.69	6.53
P_2	54.48	210.86	11.478.79	59.91	-25.45	19.20
P ₃	61.58	197.01	12.132.79	80.73	-30.35	25.87
P ₄	76.50	186.86	14,295.18	124.47	-33.94	48.28

Why? Is really ILUPACK a memory-bounded operation?



1st attempt: Dynamic Static voltage-frequency scaling

P-state P _i	Vcc _i	f _i	T_i	ΔT_i	BWi	ΔBW_i
P_0	1.23	2.00	34.06	_	30.29	
P_1	1.17	1.50	43.57	21.88	24.63	-18.67
P_2	1.12	1.20	54.48	59.91	20.46	-32.44
P_3	1.09	1.00	61.58	80.73	17.48	-42.30
P_4	1.06	0.80	76.50	124.47	14.00	-53.77

- Combined effect of linear decrease of CPU performance and memory bandwidth!
- Decrease of $P_i^s(P_0 \rightarrow P_3: -21.47\%)$, decrease of $P_i^{D_i}(P_0 \rightarrow P_3: -60.73\%)$ but P^{Y} does not change!



- 2nd attempt: DVFS during idle periods





Calculation of preconditioner



2nd attempt: DVFS during idle periods





2nd attempt: DVFS during idle periods



tasks (no dependencies)



Active polling for work...





- 3rd attempt: DVFS and idle-wait





3rd attempt: DVFS and idle-wait





- 3rd attempt: DVFS and idle-wait:
 - Savings of 6.92% of total energy
 - Negligible impact on execution time
- ...but take into account that
 - Idle time: 23.70%
 - Dynamic power: 39.32%
 - Upper bound of savings: 39.32 · 0.2370 ≈ 9.32%



P-state	$V CC_i$	f_i	$ BW_i $	ΔBW_i
P_0	1.04	2.00	12.72	
P_1	0.98	1.73	12.61	-0.87
P_2	0.95	1.60	12.55	-1.34
P_3	1.01	1.87	12.58	-1.10

DVFS

	P_i	T_i	\bar{P}_i^T	E_i	ΔT_i	$\Delta \bar{P}_i^T$	ΔE_i
	P_0	56.43	135.17	7,627.97	_	_	_
	P_1	59.06	127.96	$7,\!557.87$	4.67	-5.33	-0.92
	P_2	62.93	121.99	$7,\!676.98$	11.52	-9.75	0.64
	P_3	67.05	116.22	7,792.77	18.82	-18.82	2.16
	P_0	148.94	155.27	$23,\!123.99$	_	_	_
Solvo	P_1	148.52	151.07	$22,\!434.73$	-0.28	-2.70	-2.98
Solve	P_2	154.86	145.11	$22,\!469.38$	3.97	-6.55	-2.83
	P_3	159.08	138.50	$22,\!033.14$	6.81	-10.80	-4.72

DVFS per socket, not per core!

ILUPACK on multicore Leveraging P- and C-states (Intel)





The CG method on CPU-GPU





- GPU (or other hardware accelerator, e.g., Intel Xeon Phi):
 - High throughput
 - Reasonable power
 - Reasonable(?) cost
 - (3,000 Euro NVIDIA K20)
 - MFLOPS/W!
- Systems #1 in both green500 and top500!

The CG method on CPU-GPU



- Leveraging P-states on CPU-GPU platforms?
 - DVFS in the CPU while computation proceeds on the GPU?
- Leveraging C-states on CPU-GPU platforms?
 - What is the idle CPU doing?

The CG method on CPU-GPU **Experimental setup**

Platform:

- Intel i7-3770K (Ivy Bridge), 16GB
- NVIDIA GeForce GTX480 (Fermi)

Cases from two matrix collections

Source	Matrix	#nonzeros (n_z)	Size (n)	n_z/n
	audikw_1	77,651,847	943,645	82.28
	BMWCRA1	10,641,602	148,770	71.53
UFMC	CRANKSEG_2	14,148,858	63,838	221.63
OFMC	F1	26,837,113	343,791	78.06
	INLINE_1	38,816,170	503,712	77.06
	LDOOR	42,493,817	952,203	44.62
	A100	6,940,000	1,000,000	6.94
	A126	13,907,370	2,000,376	6.94
Laplace	A159	27,986,067	4,019,679	6.94
	A200	55,760,000	8,000,000	6.94
	A252	111,640,032	16,003,001	6.94



















CG: Sparse matrix-vector (SpMV) + CUBLAS

```
while((k < maxiter) && (res > epsilon)){
   SSpMV <<<Gs,Bs>>> (n, rowA, colA, valA, d, z);
   tmp = cublasSdot (n, d, 1, z, 1);
   rho = beta / tmp;
   gamma = beta;
   cublasSaxpy (n, rho, d, 1, x, 1);
   cublasSaxpy (n, -rho, z, 1, r, 1);
   beta = cublasSdot(n, r, 1, r, 1);
   alpha = beta / gamma;
   cublasSscal (n, alpha, d, 1);
   res = sqrt(beta);
   k++;
} // end-while
```



CG: Sparse matrix-vector (SpMV) + CUBLAS

```
while( ( k < maxiter ) && ( res > epsilon ) ){
   SSpMV <<<Gs,Bs>>> ( n, rowA, colA, valA, d, z );
   tmp = cublasSdot ( n, d, 1, z, 1 );
   rho = beta / tmp:
```

Leveraging P-states:

- Basically all computation performed on the GPU
- Apply static VFS to reduce power of CPU!

```
aipna - Deta / gamma;
cublasSscal (n, alpha, d, 1 );
cublasSaxpy (n, one, r, 1, d, 1 );
res = sqrt( beta );
k++;
} // end-while
```



CG: Sparse matrix-vector (SpMV) + CUBLAS

```
while( ( k < maxiter ) && ( res > epsilon ) ){
   SSpMV <<<Gs,Bs>>> ( n, rowA, colA, valA, d, z );
   tmp = cublasSdot ( n, d, 1, z, 1 );
   rho = beta / tmp:
```

Leveraging C-states:

- What is the CPU doing while idle?
- CUDA offers polling (active-wait) vs blocking (idle-wait) operation modes

```
cublasSaxpy (n, one, r, 1, d, 1 );
res = sqrt( beta );
k++;
} // end-while
```



Trading off power for time:
 CUDA blocking mode w.r.t. CUDA polling mode





Total power polling

Trading off power for time:
 CUDA blocking mode w.r.t. CUDA polling mode





Trading off power for time:
 CUDA blocking mode w.r.t. CUDA polling mode



The CG method on CPU-GPU Merged implementation



- Can we attain polling performance and blocking energy advantage?
- Requires a reformulation of CG (merge kernels)



The CG method on CPU-GPU Merged implementation

Time vs. CPU energy

Maintain performance of polling...







The CG method on CPU-GPU Merged implementation



Time vs. CPU energy

Maintain performance of polling...



...while leveraging power-efficiency of C-states via CUDA blocking mode





- A battle to be fought in the core arena
 - More concurrency
 - Heterogeneous designs





- A related battle to be fought in the energy arena
 - P-states and C-states:

"Doing nothing to save energy in matrix computations"

- "Do nothing, efficiently..." (V. Pallipadi, A. Belay) or "Doing nothing well" (D. E. Culler)
- Better (numerical) algorithms will reduce dynamic power but system power & static power are for the architects!





- Being energy-aware scientists?
 - Spring in Castellón: 21° C degrees (71° F) today!

Date	Fri 22			Sat 23				Sun 24		Mon 25		Mar 26	Wed 27	Thu: 20
	6-12	12-18	18-24	0-6	6-12	12-18	18-24	0-12	12-24	0-12	12-24	Mar 26	wed 27	110 28
Sky conditions	÷ģ		C	Ċ	5	<u>&</u> -	b	Ś	Ś	-;	Ś	¢2	Č-	-Ś
Prob Precip.	0%	0%	0%	15%	25%	15%	5%	25%	30%	5%	5%	25%	25%	15%
Snow bound prov. (M)					1800	1800		1700	1600				1700	1800
Temp. Min. / max. (° C)	10 / 21			9 / 19				9 / 18		11 / 21		12 / 21	14 / 21	14 / 21
Wind	K	K	-	*	N	K	-	1	>	>	>	7	1	
(Km / h)	10	10	0	5	10	10	0	10	15	10	15	20	10	0
Maximum UV Index	6			4				4		5		5		
Notices	No Risk			No Risk										



70N

65N

60N

55N

50N

45N

40N



NCEP GFS 2-meter TEMPERATURE [*F] Init: 12Z11AUG2012 -- [0] hr --> Valid Sat 12Z11AUG2012

110 105 100 95 90 85 80 75 70 65 60 55 50 45 40 35 20 15 10 -15 -20 -25 -30 -35 -40 -45 -55 -60 -65 -70 -75 -85 -90 -95



- Being energy-aware scientists?
 - ...but come back to Manchester for the summer

Summary



NCEP GFS 1760x880 sflux Forecast Grid

Acknowledgments





A. F. Martín







J. I. Aliaga, M. Barreda, M. F. Dolz, R. Mayo, J. Pérez, E. S. Quintana-Ortí



P. Alonso





 EU FP7 318793 Project "EXA2GREEN. Energy-Aware Sustainable Computing on Future Technology - Paving the Road to Exascale Computing"



Project CICYT TIN2011-23283
 "PA-HPC. Power-Aware High Performance Computing"



More information



- "Leveraging task-parallelism in energy-efficient ILU preconditioners", J. I. Aliaga, M. F. Dolz, A. F. Martín, E. S. Quintana-Ortí. ICT-GLOW 2012. Vienna (Austria)
- "Modeling power and energy of the task-parallel Cholesky factorization on multicore processors", P. Alonso, M. F. Dolz, R. Mayo, E. S. Quintana-Ortí. EnaHPC 2012. Hamburg (Germany)
- "Exploiting thread-Level parallelism in the iterative solution of sparse linear systems". J. I. Aliaga, M. Bollhöfer, A. F. Martín, E. S. Quintana-Ortí. Parallel Computing, 2011
- Reformulated Conjugate Gradient for the Energy-Aware Solution of Linear Systems on GPUs. J. I. Aliaga, J. Pérez, E. S. Quintana-Ortí, ICPP 2013. Lyon (France). Submitted

More information



- "Leveraging task-parallelism in energy-efficient ILU preconditioners", J. I. Aliaga, M. F. Dolz, A. F. Martín, E. S. Quintana-Ortí. ICT-GLOW 2012. Vienna (Austria)
- "Modeling power and energy of the task-parallel Cholesky factorization on multicore processors", P. Alonso, M. F. Dolz, R. Mayo, E. S. Quintana-Ortí. EnaHPC 2012 ARKSrg210 Gany) QUESTIONS?
- "Exploiting thread-Level parallelism in the iterative solution of sparse linear systems". J. I. Aliaga, M. Bollhöfer, A. F. Martín, E. S. Quintana-Ortí. Parallel Computing, 2011
- Reformulated Conjugate Gradient for the Energy-Aware Solution of Linear Systems on GPUs. J. I. Aliaga, J. Pérez, E. S. Quintana-Ortí, ICPP 2013. Lyon (France). Submitted