

BLUE WATERS

SUSTAINED PETASCALE COMPUTING



CUG.2017.CAFFEINATED COMPUTING

Redmond, Washington May 7-11, 2017

HPCG and HPGMG benchmark tests on MPMD mode
on Blue Waters – a Cray XE6/XK7 hybrid system

JaeHyuk Kwack and Gregory H Bauer
National Center for Supercomputing Applications



GREAT LAKES CONSORTIUM
FOR PETASCALE COMPUTATION

CRAY

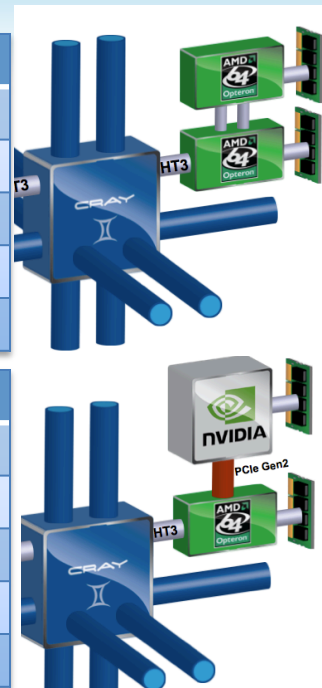
This paper has been submitted as an article in a special issue of Concurrency and Computation Practice and Experience on the Cray User Group 2017.

The Blue Waters system



22,640 XE6 compute nodes	
Number of Core Modules	32
Peak Performance	313 Gflops/sec
Memory Size	64 GB per node
Memory Bandwidth (Peak)	102 GB/sec
Interconnect Injection Bandwidth (Peak)	9.6 GB/sec per direction

4,228 XK7 compute nodes with NVIDIA Kepler (GK110) GPUs	
Host Processor	AMD Series 6200 (Interlagos)
Host Processor Performance	156.8 Gflops
Kepler Peak (DP floating point)	1.32 Tflops
Host Memory	32GB, 51 GB/sec
Kepler Memory	6GB GDDR5 capacity, > 180 GB/sec

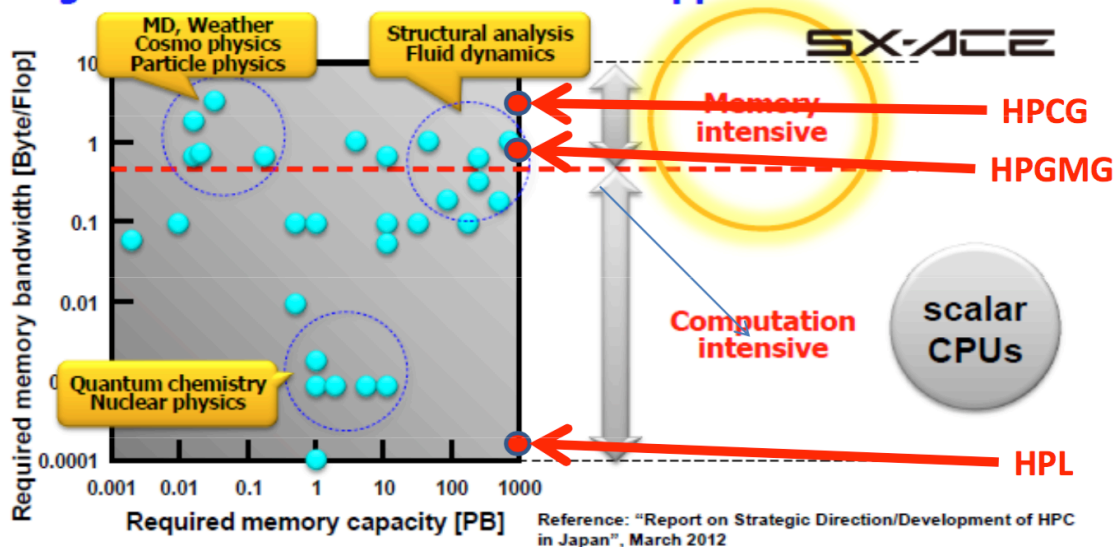


Byte/Flop in Real Applications and HPC Benchmarks

According to Japanese Government (MEXT) working group report for a wide variety of strategic segment applications, diverse characteristics are observed.

MEXT: Ministry of Education, Culture, Sports, Science & Technology

**B/F requirement from each application differs greatly.
Any single architecture cannot cover all application areas.**

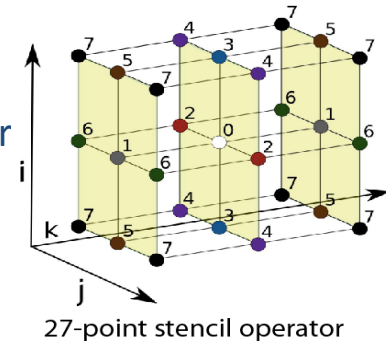


Source: Marjanovic, Gracia, Glass (HLRS, Germany), HPC-Benchmarking: Problem Size Matters, PMBS-16, 2016

HIGH PERFORMANCE CONJUGATE GRADIENTS (HPCG) BENCHMARK

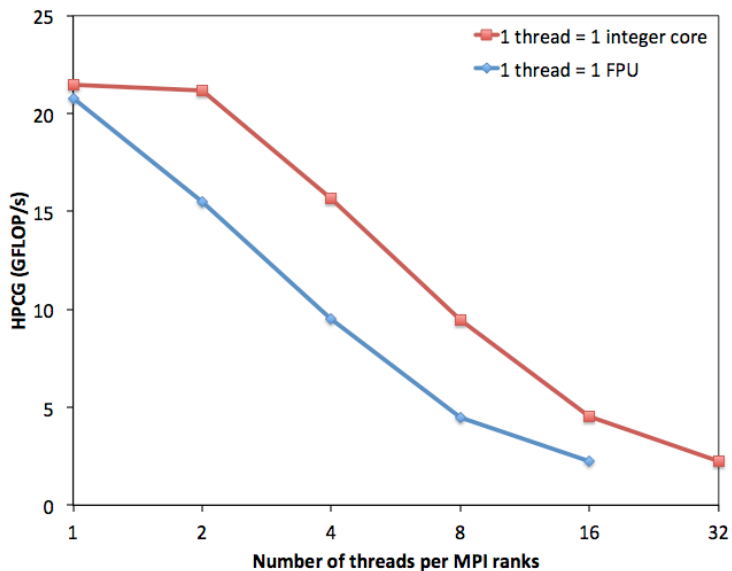
Introduction of HPCG benchmarking

- Essential computational/communicational patterns in a variety of methods for PDEs
 - Dense/sparse computations
 - Dense/sparser collectives
 - Multi-scale execution of kernels via Multi-Grid V cycle
 - Data-driven parallelism (unstructured sparse triangular solves)
- Model problem description
 - An elliptic PDE in 3D with $n_x \times n_y \times n_z$ local domain on $n_p \times n_p \times n_p$ process
 - Zero Dirichlet BCs with synthetic RHS
 - A symmetric positive definite sparse matrix from 27-point stencil operator
- Rewards investment in
 - High-performance collective ops.
 - Local memory system performance
 - Low latency cooperative threading



Source: Dongarra, Heroux, Luszczek (hpcg-benchmark.org), HPCG Update, ISC-16, 2016

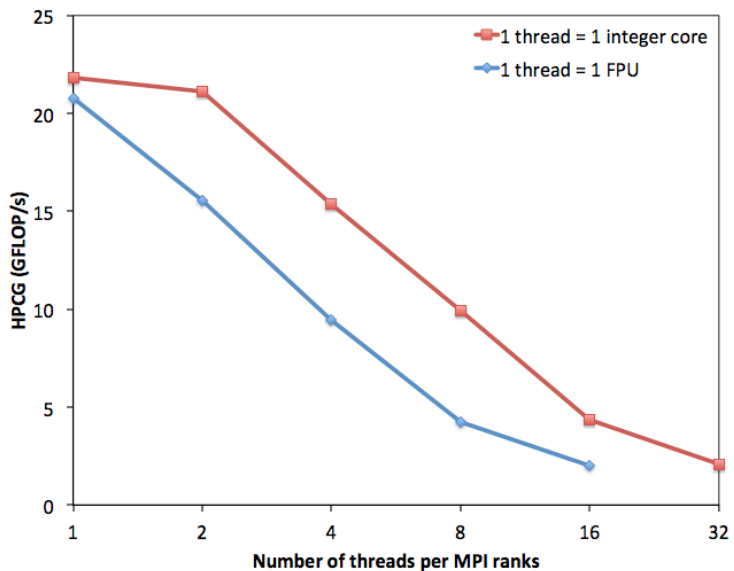
Weak scaling for MPI rank on 4 XE nodes: HPCG-CPU-GitHub-GNU⁽¹⁾



Core type per thread	MPI rank	Number of threads	HPCG (GFLOP/s)	Number of Equations in total	Memory (GB)
1 integer core per thread	128	1	21.4798	143,982,592	102.99
	64	2	21.1808	71,991,296	51.50
	32	4	15.6547	35,995,648	25.75
	16	8	9.47234	17,997,824	12.87
	8	16	4.53633	8,998,912	6.44
	4	32	2.2442	4,499,456	3.22
1 FPU per thread	64	1	20.7478	71,991,296	51.50
	32	2	15.5168	35,995,648	25.75
	16	4	9.50418	17,997,824	12.87
	8	8	4.45687	8,998,912	6.44
	4	16	2.24646	4,499,456	3.22

(1): an executable built with gcc/4.9.3 and cray-mpich/7.3.0 (HPCG 3.0 reference code)

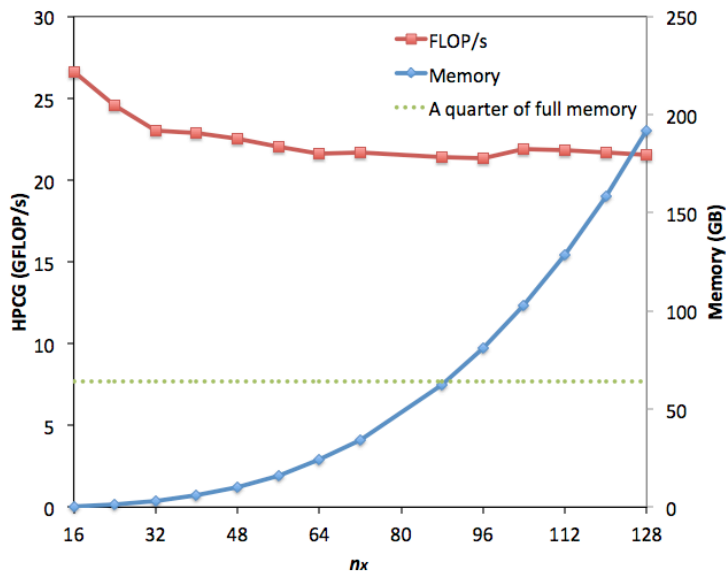
Strong scaling for MPI rank on 4 XE nodes: HPCG-CPU-GitHub-GNU⁽¹⁾



Core type per thread	MPI rank	Number of threads	HPCG (GFLOP/s)	Number of Equations in total	Memory (GB)
1 integer core per thread	128	1	21.8334	143,982,592	102.99
	64	2	21.1425	143,982,592	102.98
	32	4	15.3969	143,982,592	102.96
	16	8	9.95385	143,982,592	102.95
	8	16	4.34153	143,982,592	102.94
	4	32	2.09014	143,982,592	102.92
1 FPU per thread	64	1	20.74	143,982,592	102.98
	32	2	15.5413	143,982,592	102.96
	16	4	9.45932	143,982,592	102.95
	8	8	4.23567	143,982,592	102.94
	4	16	2.00099	143,982,592	102.92

(1): an executable built with gcc/4.9.3 and cray-mpich/7.3.0 (HPCG 3.0 reference code)

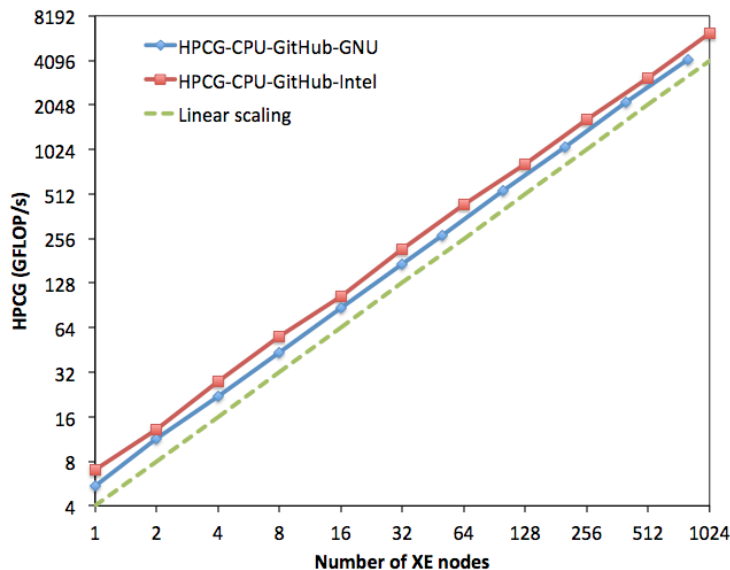
HPCG-CPU-GitHub-GNU⁽¹⁾ performance for various local blocks



n_x	n_y	n_z	HPCG (GFLOP/s)	Number of equations in total	Memory (GB)
16	16	16	26.5972	524,288	0.38
24	24	24	24.5402	1,769,472	1.27
32	32	32	23.0211	4,194,304	3.01
40	40	40	22.9008	8,192,000	5.87
48	48	48	22.5263	14,155,776	10.14
56	56	56	22.0341	22,478,848	16.09
64	64	64	21.6028	33,554,432	24.02
72	72	72	21.7075	47,775,744	34.19
88	88	88	21.3858	87,228,416	62.41
96	96	96	21.3235	113,246,208	81.01
104	104	104	21.8639	143,982,592	102.99
112	112	112	21.8343	179,830,784	128.63
120	120	120	21.6664	221,184,000	158.20
128	128	128	21.5274	268,435,456	191.98

(1): an executable built with gcc/4.9.3 and cray-mpich/7.3.0 (HPCG 3.0 reference code)

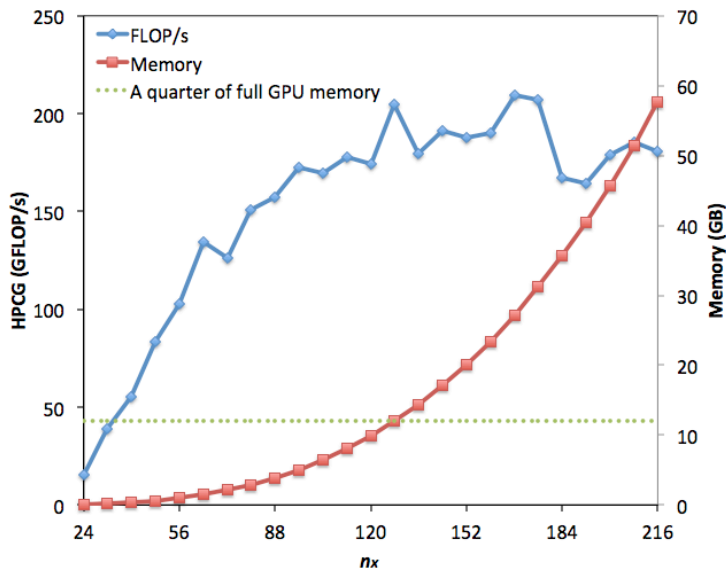
Scaling test for HPCG-CPU-GitHub-Intel⁽¹⁾ up to 4.6% Blue Waters XEs



XEs	MPI ranks	HPCG (GFLOP/s)	Number of equations in total	Memory (GB)	Efficiency (%)
1	32	7.02902	35,995,648	25.75	100
2	64	13.0688	71,991,296	51.50	93
4	128	27.9278	143,982,592	102.99	99
8	256	55.3573	287,965,184	205.99	98
16	512	104.433	575,930,368	411.97	93
32	1024	218.144	1,151,860,736	823.94	97
64	2048	435.228	2,303,721,472	1,647.88	97
128	4096	817.343	4,607,442,944	3,295.75	91
256	8192	1646.64	9,214,885,888	6591.51	92
512	16384	3103.43	18,429,771,776	13183	86
1024	32768	6218.57	36,859,543,552	26366	86

(1): an executable built with intel/16.0.3.210 and cray-mpich/7.3.0 (HPCG 3.0 reference code)

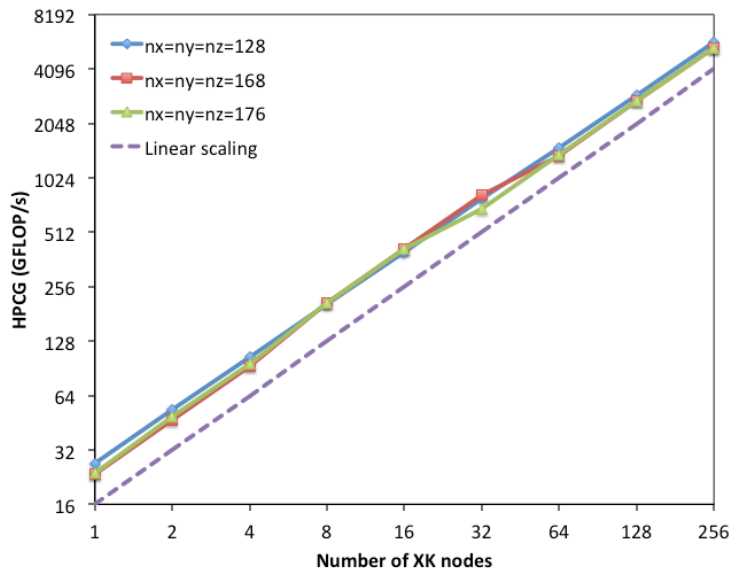
HPCG-GPU-NVIDIA⁽¹⁾ for various local blocks on 8 XK nodes



n_x	n_y	n_z	HPCG (GFLOP/s)	Number of equations in total	Memory (GB)
24	24	24	15.5542	110,592	0.079
32	32	32	38.661	262,144	0.188
40	40	40	55.1535	512,000	0.367
48	48	48	83.3617	884,736	0.634
56	56	56	102.637	1,404,928	1.006
64	64	64	134.636	2,097,152	1.501
72	72	72	126.136	2,985,984	2.137
80	80	80	150.928	4,096,000	2.931
88	88	88	157.044	5,451,776	3.900
96	96	96	172.29	7,077,888	5.063
104	104	104	169.535	8,998,912	6.437
112	112	112	177.864	11,239,424	8.039
120	120	120	174.563	13,824,000	9.887
128	128	128	204.912	16,777,216	11.999
136	136	136	179.439	20,123,648	14.392
144	144	144	191.406	23,887,872	17.083
152	152	152	187.779	28,094,464	20.091
160	160	160	190.252	32,768,000	23.432
168	168	168	209.338	37,933,056	27.125
176	176	176	207.152	43,614,208	31.187
184	184	184	167.078	49,836,032	35.635
192	192	192	164.078	56,623,104	40.487
200	200	200	178.81	64,000,000	45.761
208	208	208	185.213	71,991,296	51.474
216	216	216	180.76	80,621,568	57.644

(1): an executable compatible with cray-mpich on Blue Waters provided by NVIDIA

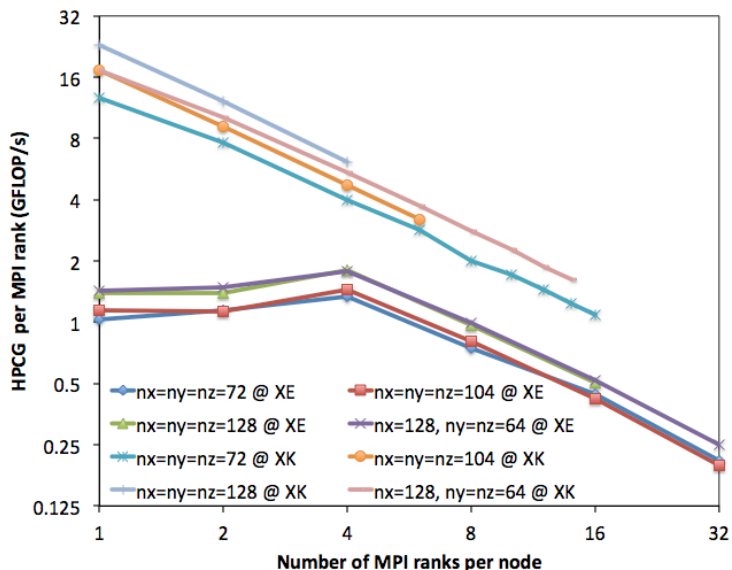
HPCG-GPU-NVIDIA⁽¹⁾ up to 6.1% Blue Waters XKs



XKs	MPI ranks	$n_x=n_y=n_z=128$		$n_x=n_y=n_z=168$		$n_x=n_y=n_z=176$	
		HPCG (GFLOP/s)	Efficiency (%)	HPCG (GFLOP/s)	Efficiency (%)	HPCG (GFLOP/s)	Efficiency (%)
1	1	27.029	100	23.544	100	23.734	100
2	2	53.353	99	46.811	99	48.960	103
4	4	103.84	96	93.085	99	95.674	101
8	8	204.87	95	209.33	111	207.19	109
16	16	397.87	92	415.67	110	412.30	109
32	32	790.06	91	827.76	110	692.05	91
64	64	1493.0	86	1347.1	89	1365.9	90
128	128	2941.0	85	2741.2	91	2724.6	90
256	256	5756.2	83	5365.2	89	5359.8	88

(1): an executable compatible with cray-mpich on Blue Waters provided by NVIDIA

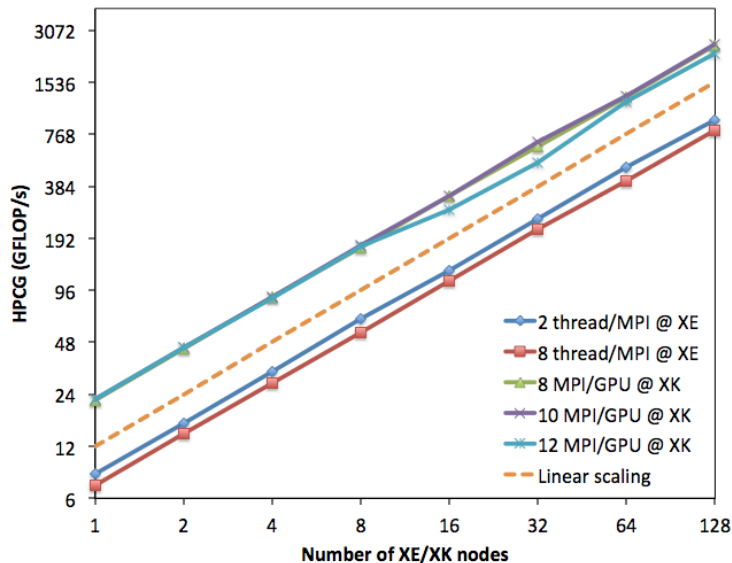
HPCG-MPMD-CPU/GPU⁽¹⁾ on a single XE or XK node



Node type	MPI ranks	Threads	nx=ny=nz=72	nx=ny=nz=104	nx=ny=nz=128	nz=128, ny=nz=64
XE	1	32	1.0335	1.13907	1.39787	1.43499
	2	16	2.28724	2.2585	2.76631	2.98869
	4	8	5.3681	5.78568	7.21688	7.14705
	8	4	5.94444	6.4733	7.73771	7.88999
	16	2	7.04848	6.74828	8.02983	8.32372
	32	1	6.69432	6.31722	OOM ⁽²⁾	7.97788
XK	1	1	12.7187	17.2333	22.9948	17.2444
	2	1	15.2734	18.4187	24.5339	20.31
	4	1	15.9005	18.9973	24.5516	21.899
	6	1	17.0141	19.2494	OOM ⁽²⁾	22.3401
	8	1	16.0738	OOM ⁽²⁾	OOM ⁽²⁾	22.3905
	10	1	17.1923	OOM ⁽²⁾	OOM ⁽²⁾	22.7505
	12	1	17.3912	OOM ⁽²⁾	OOM ⁽²⁾	22.5492
	14	1	17.2224	OOM ⁽²⁾	OOM ⁽²⁾	22.7725
	16	1	17.4242	OOM ⁽²⁾	OOM ⁽²⁾	OOM ⁽²⁾

(1): MPMD executable binaries compatible with cray-mpich on Blue Waters provided by NVIDIA, (2) OOM: Out-of-Memory error

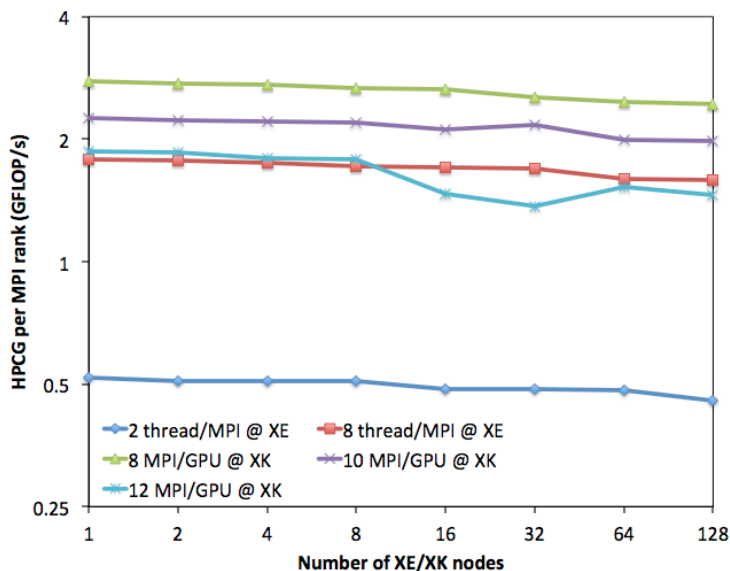
Scalability tests of HPCG-MPMD-CPU/GPU⁽¹⁾ with up to 128 XE/XKs



Number of nodes	HPCG-HPMD-CPU		HPCG-MPMD-GPU		
	<i>2thr/mpi</i> ⁽²⁾	<i>8thr/mpi</i>	<i>8mpi/gpu</i> ⁽³⁾	<i>10mpi/gpu</i>	<i>12mpi/gpu</i>
1	8.31327	7.14167	22.2208	22.5616	22.3774
2	16.2669	14.1543	43.985	44.5186	44.4957
4	32.4941	28.0463	87.3184	88.2059	86.0493
8	65.1978	54.7415	170.352	175.311	170.902
16	124.731	108.783	339.408	338.977	282.692
32	248.387	216.612	648.375	692.816	525.783
64	493.589	409.541	1266.32	1273.18	1173.86
128	930.88	812.448	2498.76	2529.21	2245.37

(1): MPMD executable binaries compatible with cray-mpich on Blue Waters provided by NVIDIA,
 (2) thr/mpi: number of threads per MPI rank, (3) mpi/gpu: number of MPI ranks per GPU

Per-MPI performance of HPCG-MPMD-CPU/GPU⁽¹⁾



Number of nodes	HPCG-HPMD-CPU		HPCG-MPMD-GPU		
	<i>2thr/mpi</i> ⁽²⁾	<i>8thr/mpi</i>	<i>8mpi/gpu</i> ⁽³⁾	<i>10mpi/gpu</i>	<i>12mpi/gpu</i>
1	0.5196	1.7854	2.7776	2.2562	1.8648
2	0.5083	1.7693	2.7491	2.2259	1.8540
4	0.5077	1.7529	2.7287	2.2051	1.7927
8	0.5094	1.7107	2.6618	2.1914	1.7802
16	0.4872	1.6997	2.6516	2.1186	1.4724
32	0.4851	1.6923	2.5327	2.1651	1.3692
64	0.4820	1.5998	2.4733	1.9893	1.5285
128	0.4545	1.5868	2.4402	1.9759	1.4618

(1): MPMD executable binaries compatible with cray-mpich on Blue Waters provided by NVIDIA,
(2) thr/mpi: number of threads per MPI rank, (3) mpi/gpu: number of MPI ranks per GPU

HPCG performance of MPMD runs with up to 128 XEs and 128 XKs

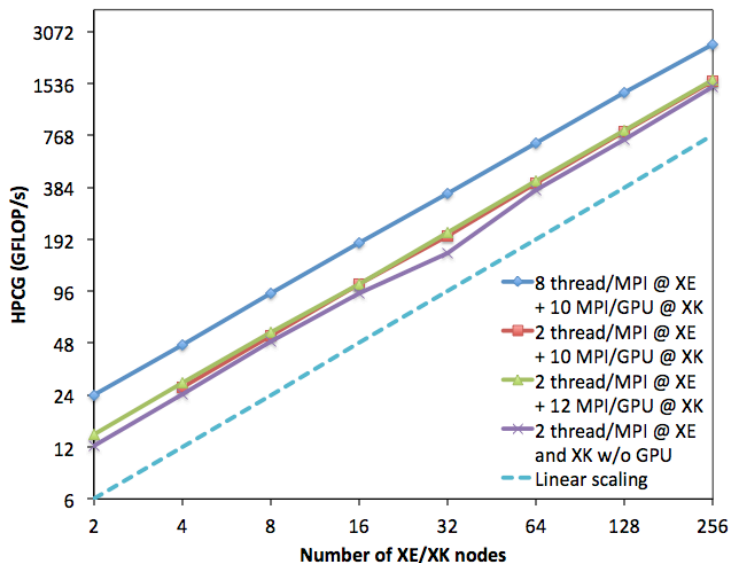
HPCG performance of MPMD runs (unit: GFLOP/s)

XE nodes	XK nodes	Case I ⁽¹⁾	Case II ⁽²⁾	Case III ⁽³⁾	Case IV ⁽⁴⁾
1	1	24.0831		14.2045	12.2138
2	2	46.7258	26.6157	28.2871	24.3066
4	4	92.937	53.0355	55.8481	48.7609
8	8	182.796	105.505	106.092	93.1538
16	16	350.558	201.75	208.912	159.893
32	32	693.168	402.937	417.709	371.003
64	64	1366.53	803.817	815.194	725.243
128	128	2599.56	1571.16	1609.91	1461.7

HPCG performance per MPI rank of MPMD runs (unit: GFLOP/s)

XE nodes	XK nodes	Case I ⁽¹⁾	Case II ⁽²⁾	Case III ⁽³⁾	Case IV ⁽⁴⁾
1	1	1.7202		0.5073	0.5089
2	2	1.6688	0.5118	0.5051	0.5064
4	4	1.6596	0.5100	0.4986	0.5079
8	8	1.6321	0.5072	0.4736	0.4852
16	16	1.5650	0.4850	0.4663	0.4164
32	32	1.5473	0.4843	0.4662	0.4831
64	64	1.5251	0.4831	0.4549	0.4722
128	128	1.4506	0.4721	0.4492	0.4758

(1): 8 thread/MPI @ XE + 10 MPI/GPU @ XE, (2): 2 thread/MPI @ XE + 10 MPI/GPU @ XK,
(3): 2 thread/MPI @ XE + 12 MPI/GPU @ XK, (4): 2 thread/MPI @ XE and XK w/o GPU



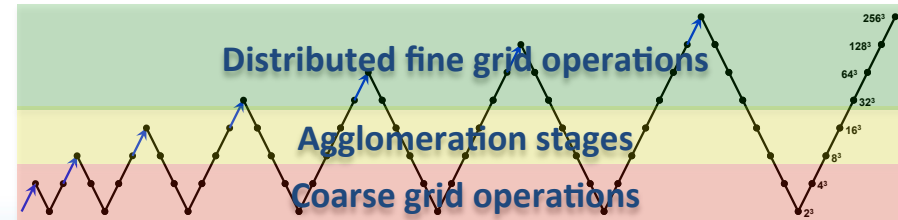
Concluding remarks - HPCG

- HPCG benchmarking on CPU-based XEs
 - Consistent performance for various numbers of equations per MPI rank
 - More than or equal to 86% parallel efficiency with up to 1024 XE nodes
- HPCG benchmarking on GPU-enabled XKs
 - Performance increases exponentially w/ number of equations per MPI rank at the beginning, and then converges to a certain level
 - More than or equal to 88% parallel efficiency with up to 256 XK nodes
 - HPCG on XKs is around 2.7 times faster than HPCG on the same number of XEs
- HPCG on MPMD mode
 - Synchronization of per-MPI performance
 - Increasing number of OpenMP threads on XEs
 - Increasing number of MPI ranks per GPU on XKs
 - The synchronization results in 60% to 100% improvement in performance

HIGH PERFORMANCE GEOMETRIC MULTI- GRID (HPGMG) BENCHMARK

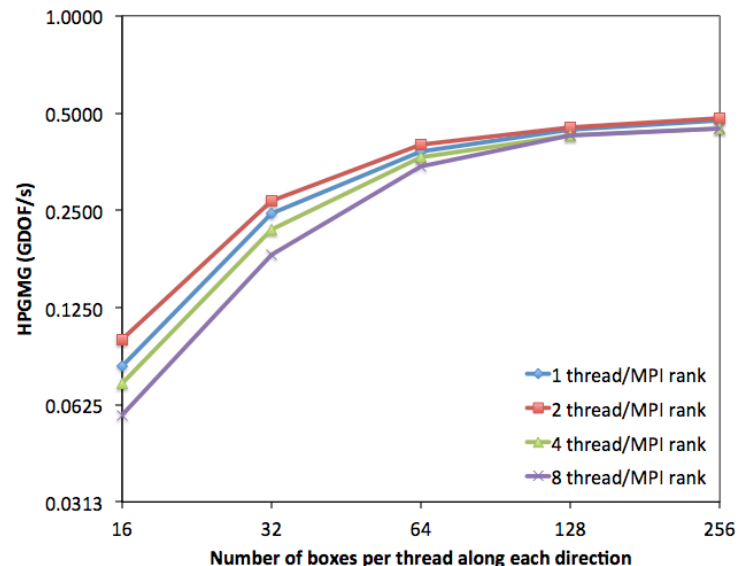
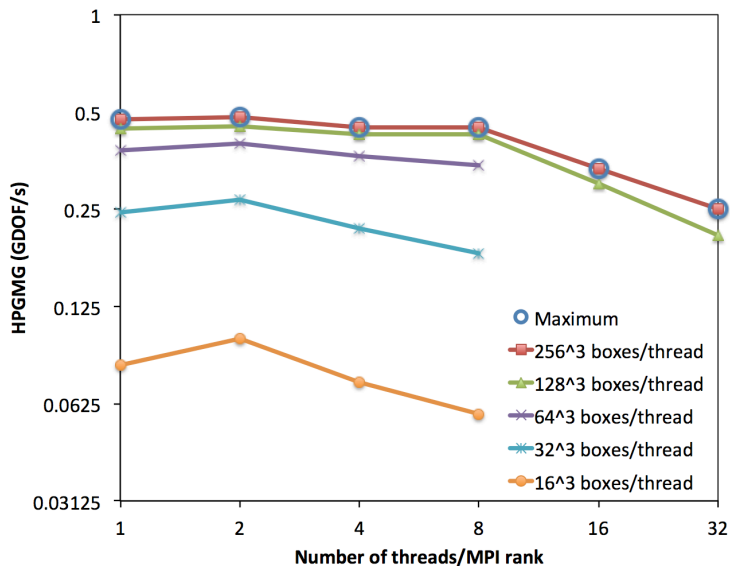
Introduction of HPGMG benchmarking

- An effort for HPC performance benchmarking on geometric multi-grid methods
 - Community-driven development process
 - Long-term durability
 - Scale-free specification and scale-free communication
- HPGMG-FE(Finite Element): compute-intensive and cache-intensive
- HPGMG-FV(Finite Volume): memory bandwidth-intensive
 - Used for the official list
 - Solving an elliptic problem on isotropic Cartesian grids with 4th order accuracy
 - 4× FP ops, 3× MPI messages, 2× MPI message size w/o DRAM data movement compared to 2th order HPGMG-FV
 - Employing the Full Multi-grid (FMG) F-cycle
 - A series of progressively deeper geometric multi-grid V-cycles



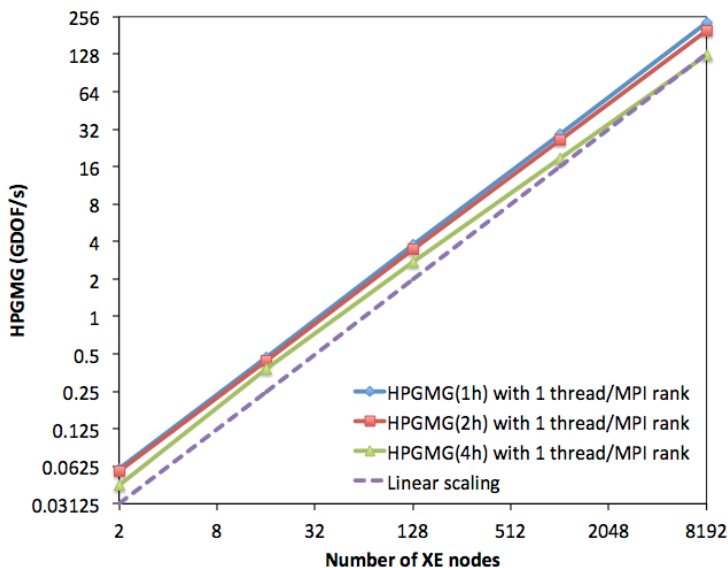
Source: Williams (hpgmg.org), HPGMG BoF, SC-16, 2016

HPGMG performance on 16 XE nodes: HPGMG-CPU-8a2f0e1⁽¹⁾



(1): a HPGMG executable binary built with gcc/4.9.3 and cray-mpich/7.3.3 (commit: 8a2f0e1)

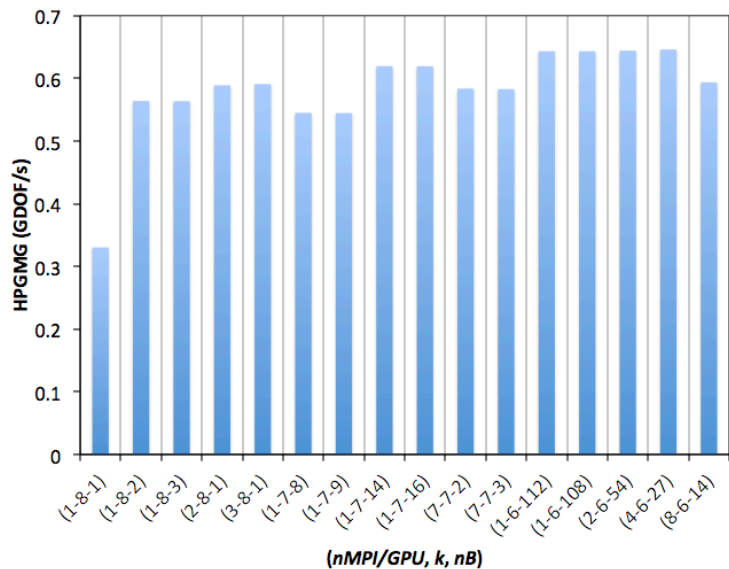
HPGMG performance⁽¹⁾ with up to 36% Blue Waters XEs



<i>nthr</i>	XE nodes	MPI ranks	Cells/thread	HPGMG (GDOF/s)			Efficiency (%)	HPGMG /MPI
				1h	2h	4h		
1	2	64	256 ³	0.06016	0.05753	0.04387	100	0.000940
1	16	512	256 ³	0.474	0.4441	0.3767	98	0.000926
1	128	4096	256 ³	3.762	3.466	2.752	98	0.000918
1	1024	32768	256 ³	29.35	26.18	18.62	95	0.000896
1	8192	262144	256 ³	229.6	198	126.3	93	0.000876
2	2	32	256 ³	0.06054	0.05777	0.05312	100	0.001892
2	16	256	256 ³	0.4708	0.4493	0.3952	97	0.001839
2	128	2048	256 ³	3.784	3.496	2.884	98	0.001848
2	1024	16384	256 ³	28.85	25.39	19.03	93	0.001761
2	8192	131072	256 ³	223.7	189.5	125.1	90	0.001707
4	2	16	256 ³	0.05672	0.05449	0.04811	100	0.003545
4	16	128	256 ³	0.4473	0.4245	0.3591	99	0.003495
4	128	1024	256 ³	3.526	3.286	2.654	97	0.003443
4	1024	8192	256 ³	26.63	23.3	17.08	92	0.003251
4	8192	65536	256 ³	205.9	172.8	112.5	89	0.003142
8	2	8	256 ³	0.05662	0.05462	0.04338	100	0.007078
8	16	64	256 ³	0.4485	0.4244	0.3367	99	0.007008
8	128	512	256 ³	3.503	2.312	2.463	97	0.006842
8	1024	4096	256 ³	26.3	22.86	15.59	91	0.006421
8	8192	32768	256 ³	201.5	165.7	100.4	87	0.006149

(1): with HPGMG-CPU-8a2f0e1, a HPGMG executable binary built with gcc/4.9.3 and cray-mpich/7.3.3 (commit: 8a2f0e1)

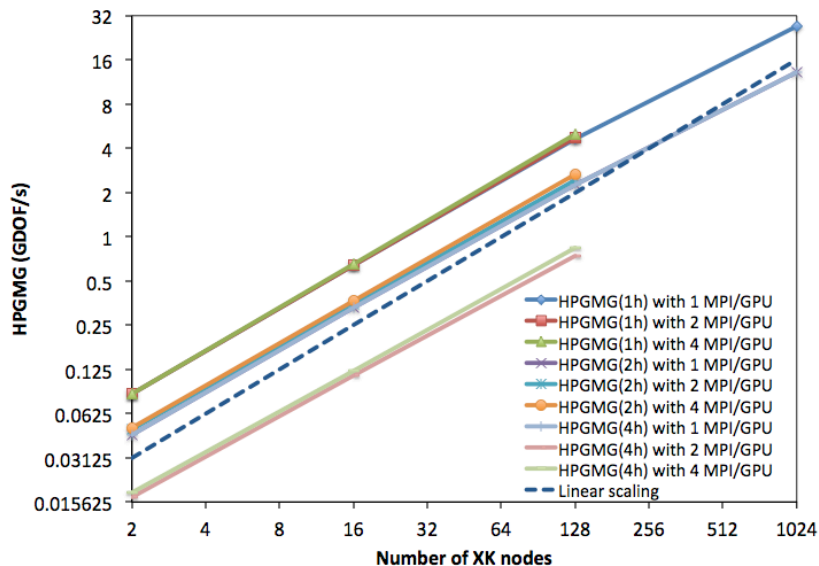
HPGMG-GPU-02c7ea2 ⁽¹⁾ on 16 XKs



nMPI /GPU ⁽²⁾	K ⁽³⁾	nB ⁽⁴⁾	XE nodes	nMPI ⁽⁵⁾	Boxes / MPI ⁽⁶⁾	Total boxes	HPGMG (1h)
1	8	1	16	16	1	512 ³	0.3304
1	8	2	16	16	1.688	768 ³	0.5639
1	8	3	16	16	1.688	768 ³	0.5636
2	8	1	16	32	1.688	768 ³	0.589
3	8	1	16	48	0.562	768 ³	0.5909
1	7	8	16	16	7.812	640 ³	0.545
1	7	9	16	16	7.812	640 ³	0.5447
1	7	14	16	16	13.5	768 ³	0.6195
1	7	16	16	16	13.5	768 ³	0.6194
7	7	2	16	112	1.929	768 ³	0.5837
7	7	3	16	112	1.929	768 ³	0.5829
1	6	112	16	16	108	768 ³	0.6433
1	6	108	16	16	108	768 ³	0.6433
2	6	54	16	32	54	768 ³	0.6442
4	6	27	16	64	27	768 ³	0.646
8	6	14	16	128	13.5	768 ³	0.5937

- (1): a HPGMG executable binary built with CUDA 7.5 and cray-mpich/7.3.3 (commit: 02c7ea214ce8),
 (2) Number of MPI rank assigned to each GPU, (3) the log base 2 of the dimension of each box on the finest grid,
 (4) Target number of boxes per process, (5) Number of MPI rank, (6) Number of boxes assigned to process

HPGMG scalability test⁽¹⁾ with up to 24% Blue Waters XKs



XK nodes	MPI ranks	Cells/ MPI rank	Number of cells	HPGMG (GDOF/s)			Efficiency	HPGMG /MPI rank
				1h	2h	4h		
2	2	108 × 64 ³	384 ³	0.08475	0.04496	0.0144	100%	0.042375
16	16	108 × 64 ³	768 ³	0.6433	0.3325	0.1043	95%	0.040206
128	128	108 × 64 ³	1536 ³	4.611	2.271	0.697	85%	0.036023
1024	1024	108 × 64 ³	3072 ³	27.24	13.11	4.104	63%	0.026602
2	4	54 × 64 ³	384 ³	0.08539	0.0486	0.01677	100%	0.021348
16	32	54 × 64 ³	768 ³	0.6446	0.3463	0.1133	94%	0.020144
128	256	54 × 64 ³	1536 ³	4.734	2.419	0.7401	87%	0.018492
1024	2048	54 × 64 ³	3072 ³	malloc failed - create_level/level->my_boxes				
2	8	27 × 64 ³	384 ³	0.08531	0.05022	0.01801	100%	0.010664
16	64	27 × 64 ³	768 ³	0.6492	0.366	0.1217	95%	0.010144
128	512	27 × 64 ³	1536 ³	4.972	2.675	0.8368	91%	0.009711
1024	4096	27 × 64 ³	3072 ³	malloc failed - create_level/level->my_boxes				

(1): with HPGMG-GPU-02c8ea2, a HPGMG executable binary built with CUDA 7.5 and cray-mpich/7.3.3 (commit: 02c7ea214ce8)

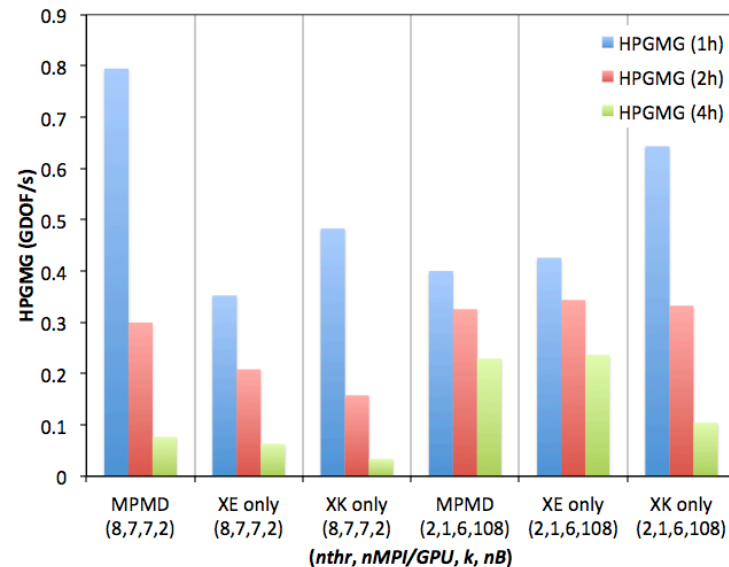
HPGMG in MPMD mode on 16 XE and 16 XK nodes

With synchronized per-MPI process performance

Run type (<i>nthr, nMPI/GPU, k, nB</i>)	Number of cells	Boxes /MPI	HPGMG (GDOF/s)			HPGMG /MPI rank
			1h	2h	4h	
16 XE only (8,7,7,2)	640 ³	1.953	0.3525	0.2084	0.063	0.005508
16 XK only (8,7,7,2)	768 ³	1.929	0.4828	0.1578	0.03365	0.004311
MPMD (8,7,7,2)	896 ³	1.949	0.7949	0.2997	0.0768	0.004516

With the configuration for the best GPU performance

Run type (<i>nthr, nMPI/GPU, k, nB</i>)	Number of cells	Boxes /MPI	HPGMG (GDOF/s)			HPGMG /MPI rank
			1h	2h	4h	
16 XE only (2,1,6,108)	768 ³	85.75	0.426	0.3437	0.2363	0.001664
16 XK only (2,1,6,108)	1792 ³	108	0.6433	0.3325	0.1043	0.040206
MPMD (2,1,6,108)	1792 ³	80.706	0.4002	0.3257	0.2294	0.001471



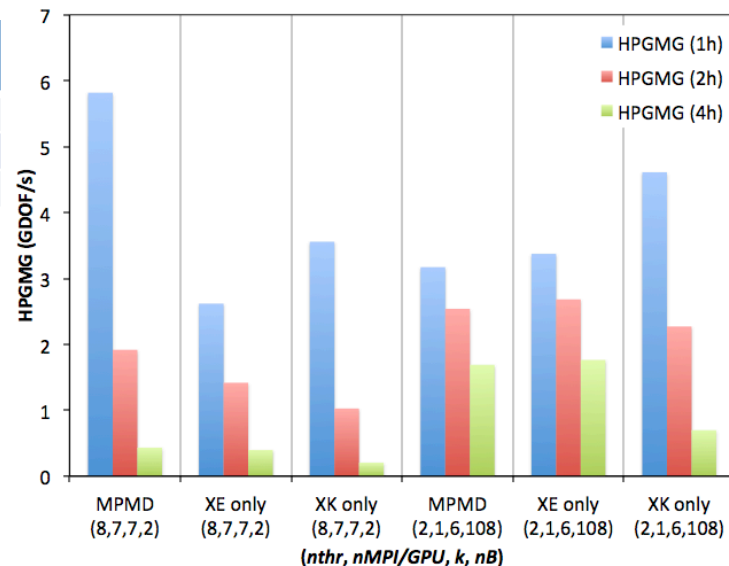
HPGMG in MPMD mode on 128 XE and 128 XK nodes

With synchronized per-MPI process performance

Run type (<i>nthr, nMPI/ GPU, k, nB</i>)	Number of cells	Boxes /MPI	HPGMG (GDOF/s)			HPGMG /MPI rank
			1h	2h	4h	
128 XE only (8,7,7,2)	1280 ³	1.953	2.618	1.417	0.3976	0.005113
128 XK only (8,7,7,2)	1536 ³	1.929	3.557	1.026	0.2039	0.003970
MPMD (8,7,7,2)	1792 ³	1.949	5.818	1.917	0.433	0.004132

With the configuration for the best GPU performance

Run type (<i>nthr, nMPI/ GPU, k, nB</i>)	Number of cells	Boxes /MPI	HPGMG (GDOF/s)			HPGMG /MPI rank
			1h	2h	4h	
128 XE only (2,1,6,108)	1536 ³	85.75	3.373	2.682	1.764	0.001647
128 XK only (2,1,6,108)	3584 ³	108	4.611	2.271	0.697	0.036023
MPMD (2,1,6,108)	3584 ³	80.706	3.171	2.539	1.69	0.001457



Concluding remarks - HPGMG

- HPGMG benchmarking on CPU-based XEs
 - HPGMG performance rapidly increases with number of cells/MPI, and then converges.
 - Very effective OpenMP implementation in a NUMA core
 - Very impressive parallel efficiency: $\geq 87\%$ with up to 36% Blue Waters XE nodes (=8192 XEs)
- HPGMG benchmarking on GPU-enabled XKs
 - HPGMG on 16 XKs is around 34% faster than HPGMG on 16 XEs
 - Parallel efficiency rapidly drops to 64% on 24% Blue Waters XK nodes (=1024 XKs)
- HPGMG on MPMD mode
 - Synchronization of per-MPI performance
 - Increasing number of OpenMP threads on XEs
 - Increasing number of MPI ranks per GPU on XKs
 - The synchronization results in 100% improvement in performance compared to others

QUESTIONS ?

Acknowledgment

This study is part of the Blue Waters sustained-petascale computing project, which is supported by the National Science Foundation (awards OCI-0725070 and ACI-1238993) and the state of Illinois. Blue Waters is a joint effort of the University of Illinois at Urbana-Champaign and its National Center for Supercomputing Applications. We thank Massimiliano Fatica and Mauron Bisson at NVIDIA for providing HPCG binaries used in this paper. We also thank Samuel Williams at LBNL, Nikolai Sakharnykh at NVIDIA and Vladimir Marjanovic at HLRS for sharing your precious experience on HPGMG benchmarking with us.



National Science Foundation
WHERE DISCOVERIES BEGIN



ILLINOIS

References

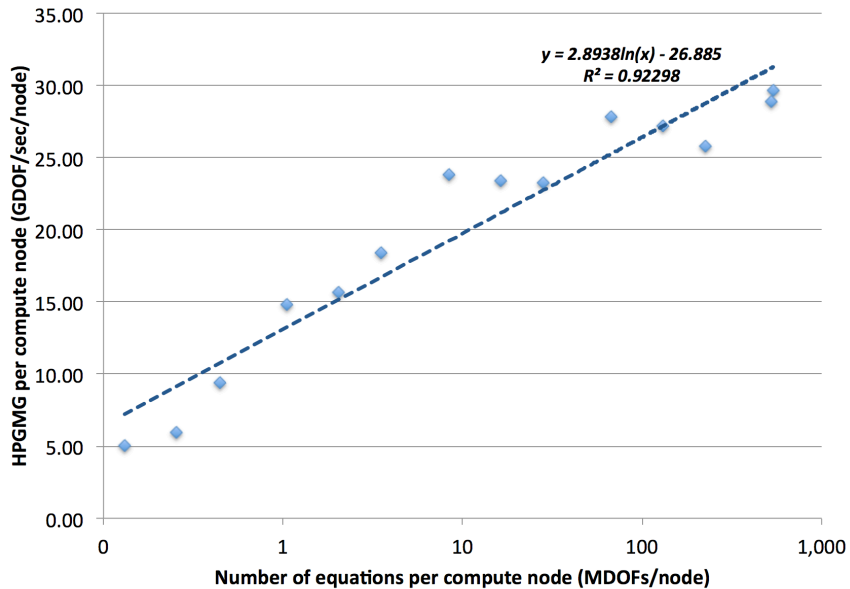
- J. Dongarra, J. Bunch, C. Moler and G.W. Stewart “LINPACK Users Guide,” SIAM, Philadelphia, PA, 1979.
- M. Heroux, J. Dongarra and P. Luszczek “HPCG Technical Specification,” Sandia National Laboratories Technical Report, SAND2013-8752, October, 2013.
- J. Dongarra, M. Heroux and P. Luszczek “HPCG Benchmark: a New Metric for Ranking High Performance Computing Systems,” Technical Report, Electrical Engineering and Computer Science Department, Knoxville, Tennessee, UT-EECS-15-736, November, 2015.
- M. Adams, J. Brown, J. Shalf, B. Straalen, E. Strohmaier and S. Williams, “HPGMG 1.0: A Benchmark for Ranking High Performance Computing Systems,” LBNL Technical Report, 2014, LBNL 6630E.
- S. Williams, “4th Order HPGMG-FV Implementation,” HPGMG BoF, Supercomputing, November 2015.
- J. Kwack, G. Bauer and S. Koric, “Performance Test of Parallel Linear Equation Solvers on Blue Waters – Cray XE6/XK7 system,” Proceedings of the Cray Users Group Meeting (CUG2016), London, England, May 2016.
- B. Bode, M. Butler, T. Dunning, W. Gropp, T. Hoefler, W. Hwu, and W. Kramer (alphabetical), “The Blue Waters Super-System for Super-Science,” Contemporary HPC Architectures, Jeffery Vetter editor. Sitka Publications, November 2012. Edited by Jeffrey S. Vetter, Chapman and Hall/CRC 2013, Print ISBN: 978-1-4665-6834-1, eBook ISBN: 978-1-4665-6835-8.
- W. Kramer, M. Butler, G. Bauer, K. Chadalavada and C. Mendes, “Blue Waters Parallel I/O Storage Sub-system,” High Performance Parallel I/O, Prabhat and Quincey Koziol editors, CRC Publications, Taylor and Francis Group, Boca Raton FL, 2015, Hardback Print ISBN 13:978-1-4665-8234-7.

Bonus slides

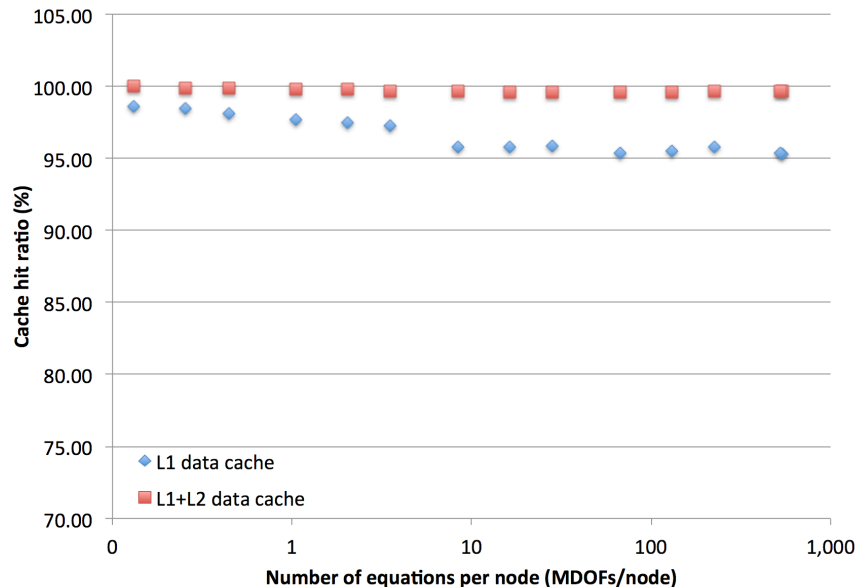
PROFILING HPGMG (WHY HPGMG PERFORMANCE IS PROPORTIONAL TO LOCAL BLOCK SIZE?)

Profiling HPGMG (why HPGMG is proportional to local problem size?)

DOFS/node vs. HPGMG

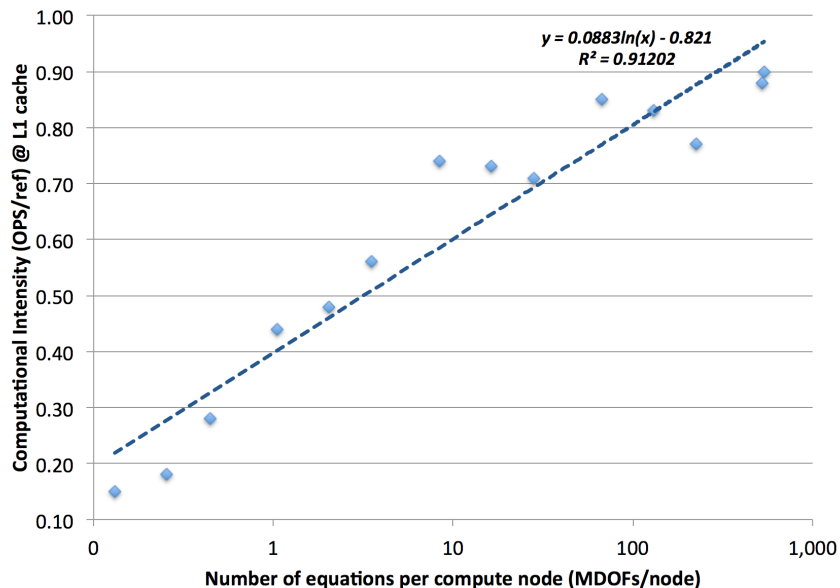


DOFS/node vs. Cache hit ratio



Profiling HPGMG (why HPGMG is proportional to local problem size?)

DOFS/node vs. Computational intensity @ L1 cache



Computational intensity @ L1 cache vs. HPGMG

