

From Laptops to Petaflops

Portable Productivity and Performance
using Open-Source Development Tools

Jeff Hammond and Katherine Riley
(with plenty of help from our colleagues)



Overview

The mission of the Argonne Leadership Computing Facility is to accelerate major scientific discoveries and engineering breakthroughs for humanity by designing and providing world-leading computing facilities in partnership with the computational science community.

Breakthrough research at the ALCF aims to:

- Reduce our national dependence on foreign oil and promote green energy alternatives
- Aid in curing life-threatening blood disorders
- Improve the safety of nuclear reactors
- Assess the impacts of regional climate change
- Cut aerodynamic carbon emissions and noise
- Speed protein mapping efforts

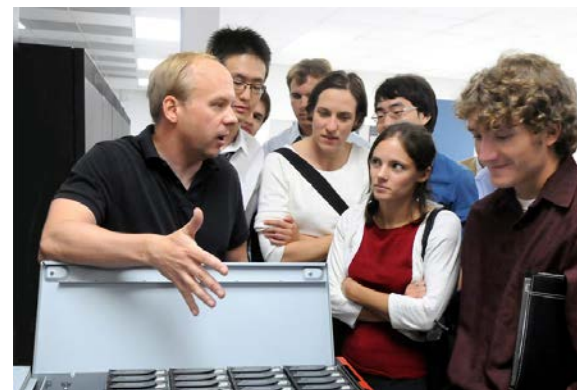


Argonne National Laboratory is a
U.S. Department of Energy laboratory
managed by UChicago Argonne, LLC



Transformational Science

The ALCF provides a computing environment that enables researchers from around the world to conduct breakthrough science.



Astounding Computation

The Argonne Leadership Computing Facility is currently home to one of the planet's fastest supercomputers. In 2012, Argonne will house Mira, a computer capable of running programs at 10 quadrillion calculations per second—that means it will compute in one second what it would take every man, woman and child in the U.S. to do if they performed a calculation every second for a year!

www.alcf.anl.gov | info@alcf.anl.gov | (877) 737-8615



Current IBM Blue Gene System, Intrepid

- 163,840 processors
- 80 terabytes of memory
- 557 teraflops
- Energy-efficient system uses one-third the electricity of machines built with conventional parts
- Ranked 13th fastest computer in the world today
- #1 on Graph500 (not currently)

The groundbreaking Blue Gene

- General-purpose architecture excels in virtually all areas of computational science
- Presents an essentially standard Linux/PowerPC programming environment
- Significant impact on HPC – Blue Gene systems are consistently found in the top ten list
- Delivers excellent performance per watt
- High reliability and availability



IBM Blue Gene/Q, Mira – Arriving 2012

- System arriving in 2012:
IBM Blue Gene/Q, Mira
 - $48 \times 1024 \times 16 = 786,432$ cores
 - 768 terabytes of memory
 - 10 petaflops



Design Parameters	Blue Gene/P	Blue Gene/Q	Change
Cores per Node	4	16	4×
Clock Speed (GHz)	0.85	1.6	1.9×
Flops per Clock per Core	4	8	2×
Nodes per Rack	1,024	1,024	--
RAM per Core (GB)	0.5	1	2×
Flops per Node (GF)	13.6	204.8	15×
Concurrency per Rack	4,096	65,536	16×
Network Interconnect	3D torus	5D torus	Smaller diameter
Cooling	Air	Water	~30% savings per watt



Programs for Obtaining System Allocations

60%	30%	10%	
Innovative and Novel Computational Impact on Theory and Experiment (INCITE)	ASCR Leadership Computing Challenge Program (ALCC)	Early Science Program (ESP)	Discretionary Projects
ALCF resources are available to researchers as part of the U.S. Department of Energy's INCITE program. Established in 2003, the program encompasses high-end computing resources at Argonne and other national laboratories. The INCITE program specifically seeks out computationally intensive, large-scale research projects with the potential to significantly advance key areas in science and engineering. The program encourages proposals from universities, other research institutions, and industry. It continues to expand, with current research applications in areas such as chemistry, combustion, astrophysics, genetics, materials science and turbulence.	Open to scientists from the research community in academia and industry, the ALCC program allocates resources to projects with an emphasis on areas directly related to the Department of Energy's energy mission, national emergencies, or for broadening the community of researchers capable of using leadership computing resources. Projects are awarded an ALCC allocation based on a peer review for scientific merit and computational readiness.	Allocations through the Early Science Program (ESP) provide researchers with preproduction hours (between system installation and full production) on the ALCF's next-generation, 10 petaflops IBM Blue Gene system. This early science period provides projects with a significant head start for adapting to the new machine and access to substantial computational time. During this shakedown period, users assist in identifying the root causes of any system instabilities, and work with ALCF staff to help develop solutions. More than four billion core hours are allocated through ESP.	Discretionary allocations are "start up" awards made to potential future INCITE projects. Projects must demonstrate a need for leadership-class resources. Awards may be made year round to industry, academia, laboratories and others, and are usually between three and six months in duration. The size of the award varies based on the application and its readiness/ability to scale; awards are generally from the low tens of thousands to the low millions of hours.

For more information:
Email Jeff Hammond
jhammond@alcf.anl.gov

The U.S. Department of Energy's **INCITE** Program

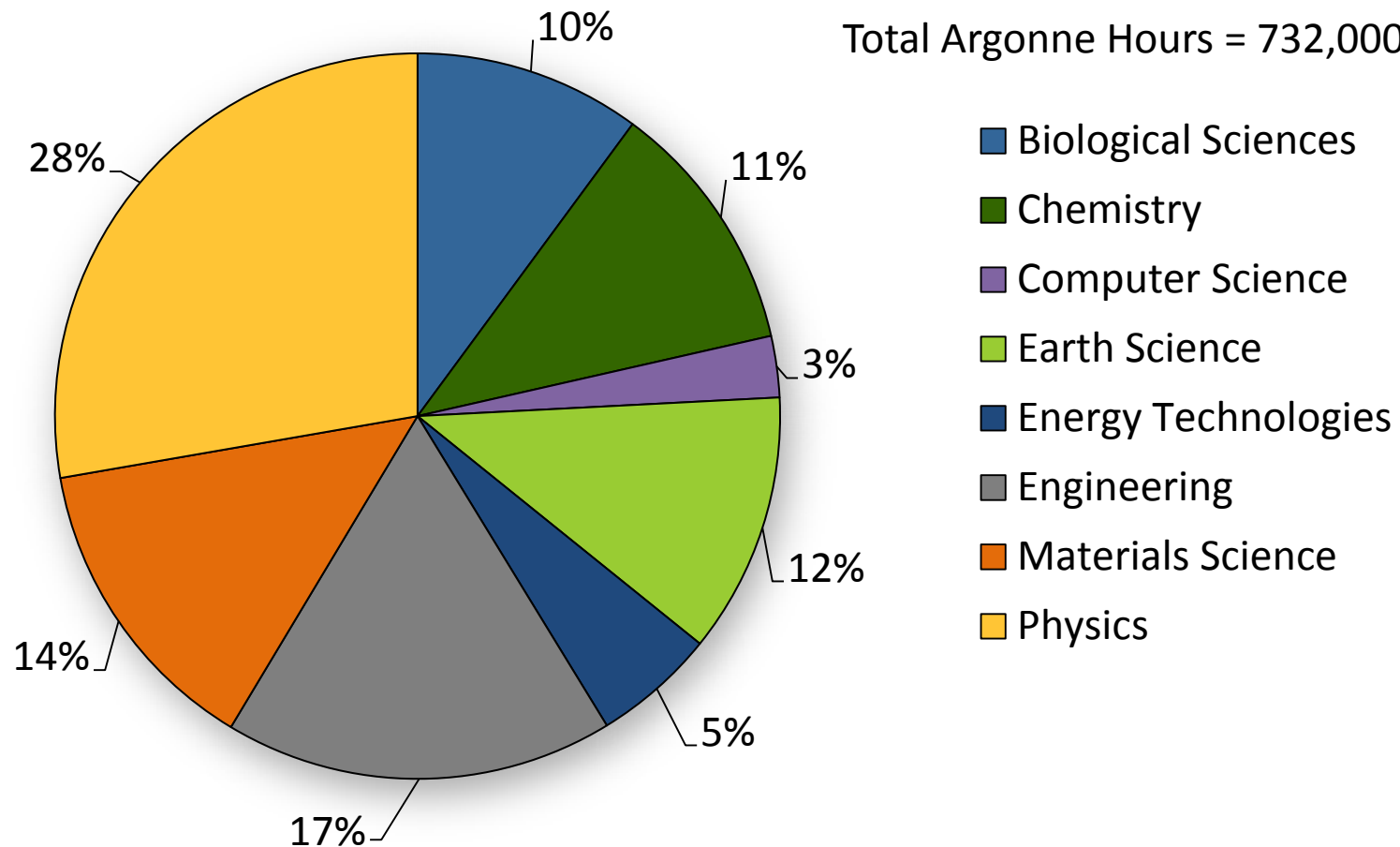
INCITE seeks out large, computationally intensive research projects and awards more than a billion processing hours to enable high-impact scientific advances.

- Open to researchers in academia, industry, and other organizations
- Proposed projects undergo scientific and computational readiness reviews
- More than a billion total hours are awarded to a small number of projects
- Sixty percent of the ALCF's processing hours go to INCITE projects
- Call for proposals issued once per year

Innovative and
N
ovel
C
omputational
Impact on
T
heory and
E
xperiment

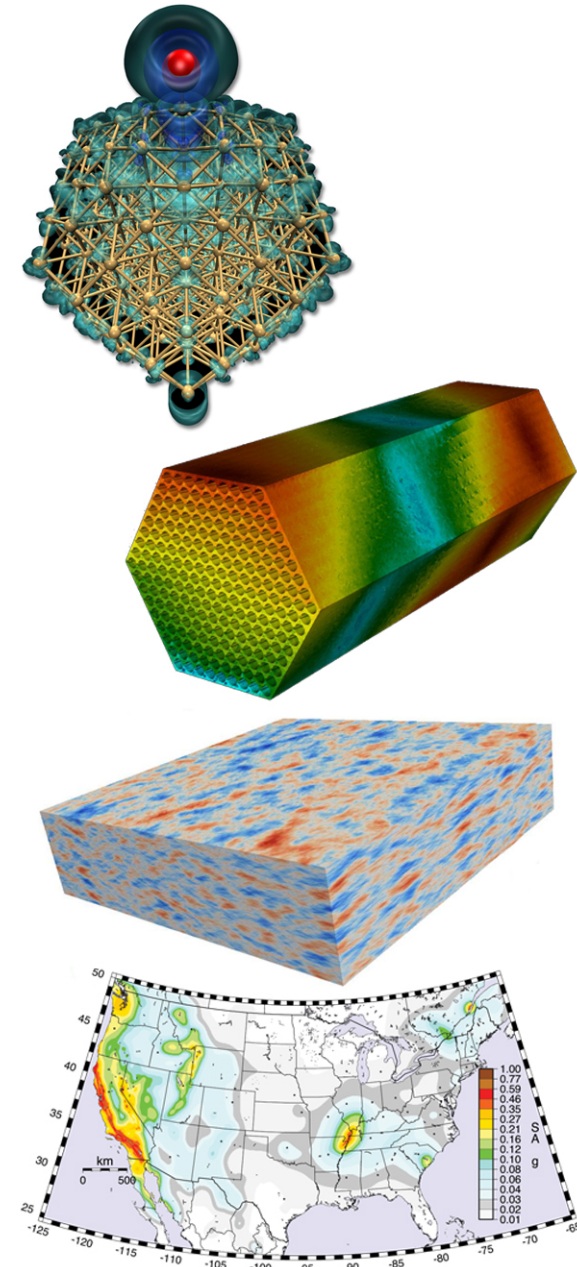


2011 INCITE Allocations by Discipline



Science Underway at the ALCF

- Research that will lead to improved, emissions-reducing catalytic systems for industry (Greeley)
- Enhancing public safety through more accurate earthquake forecasting (Jordan)
- Designing more efficient nuclear reactors that are less susceptible to dangerous, costly failures (Fischer)
- Accelerating research that may improve diagnosis and treatment for patients with blood-flow complications (Karniadakis)
- Protein studies that will apply to a broad range of problems, such as finding a cure for Alzheimer's disease, creating inhibitors of pandemic influenza, or engineering a step in the production of biofuels (Baker)
- Furthering research to bring green energy sources, like hydrogen fuel, safely into our everyday lives, reducing our dependence on foreign fuels (Khoklov)



Blood Flow: Multi-scale Modeling and Visualization

Leopold Grinberg, George Karniadakis (Brown University), Dmitry A. Fedosov (Forschungszentrum Juelich)
Bruce Caswell (Brown University), Vitali Morozov, Joseph A. Insley, Michael E. Papka, (Argonne National Laboratory)

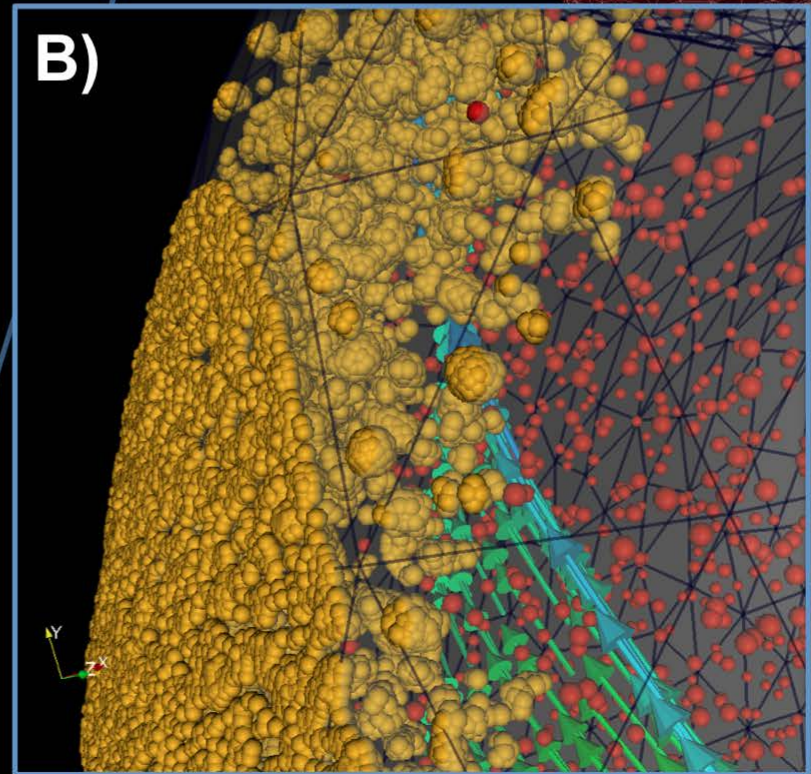
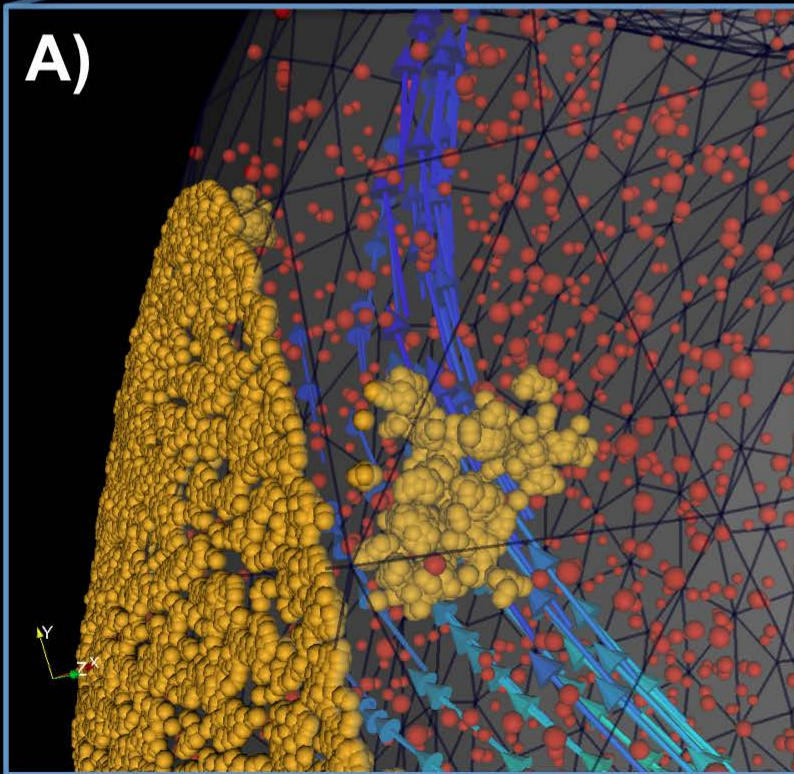
**Right Interior
Carotid Artery**

Aneurysm

**Platelet
Aggregation**

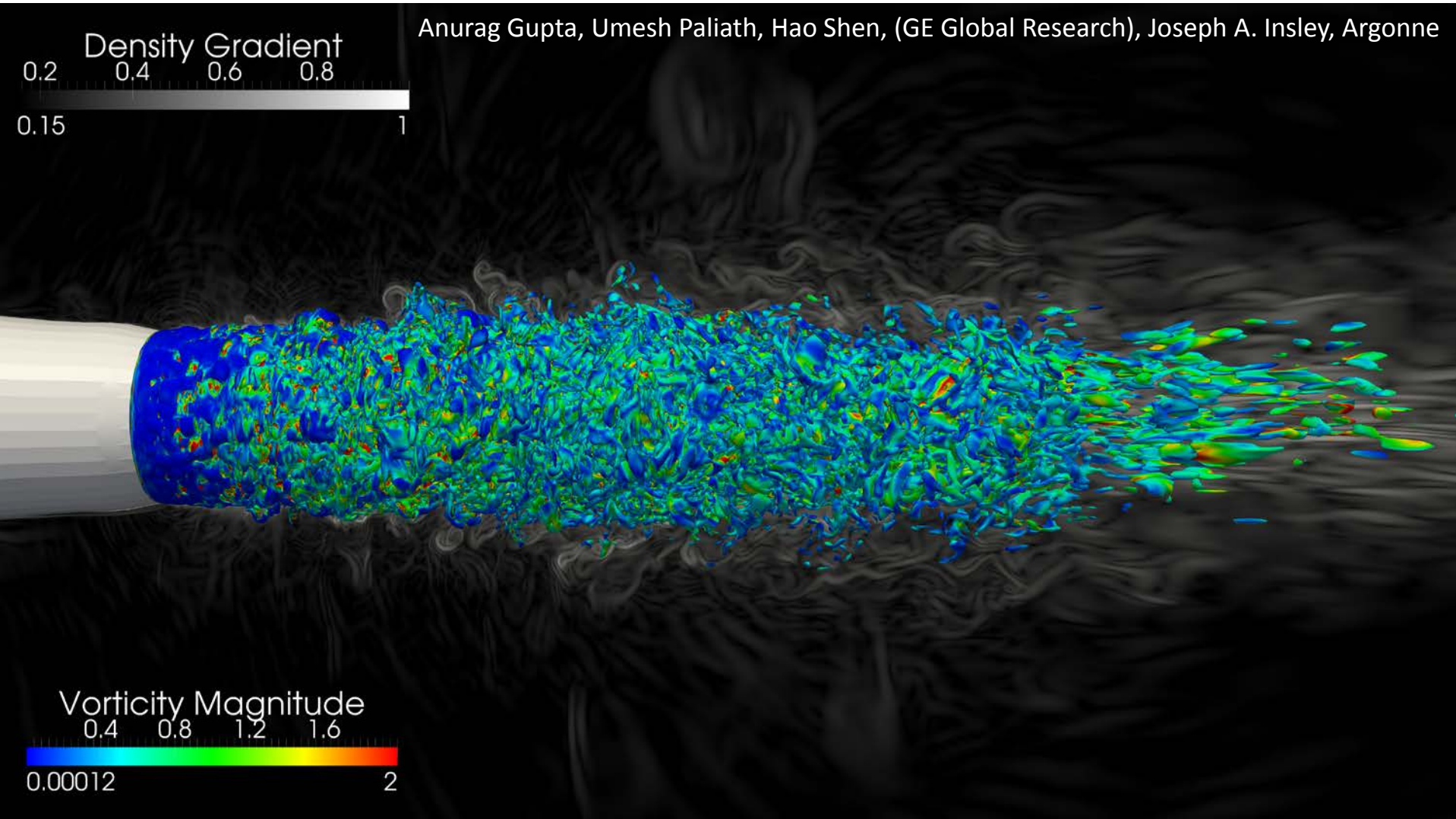
A)

B)



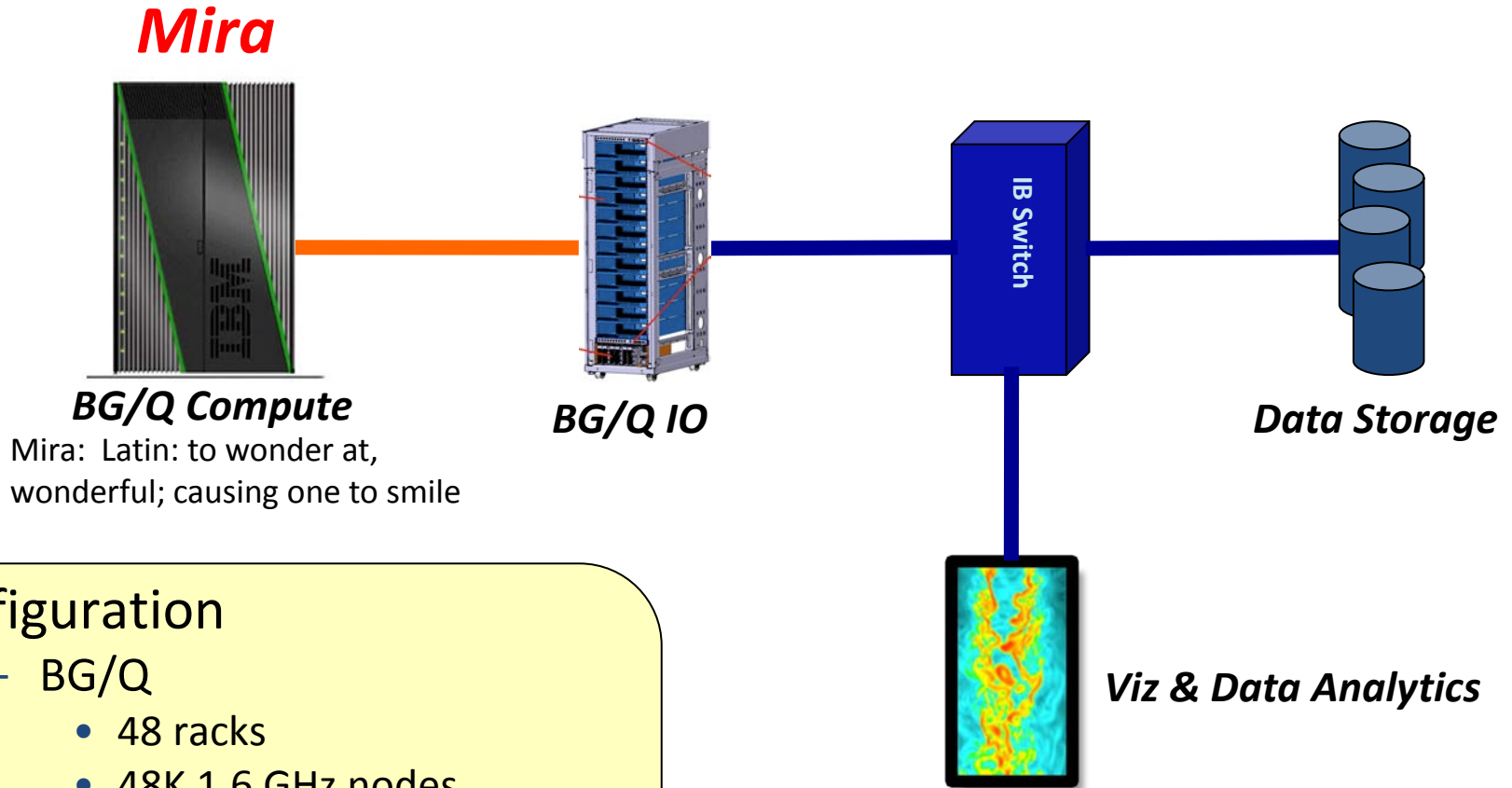
Turbulent Mixing Noise from Jet Exhaust Nozzles

Anurag Gupta, Umesh Paliath, Hao Shen, (GE Global Research), Joseph A. Insley, Argonne



Turbulent structures in free shear layer flow from dual flow conic nozzle

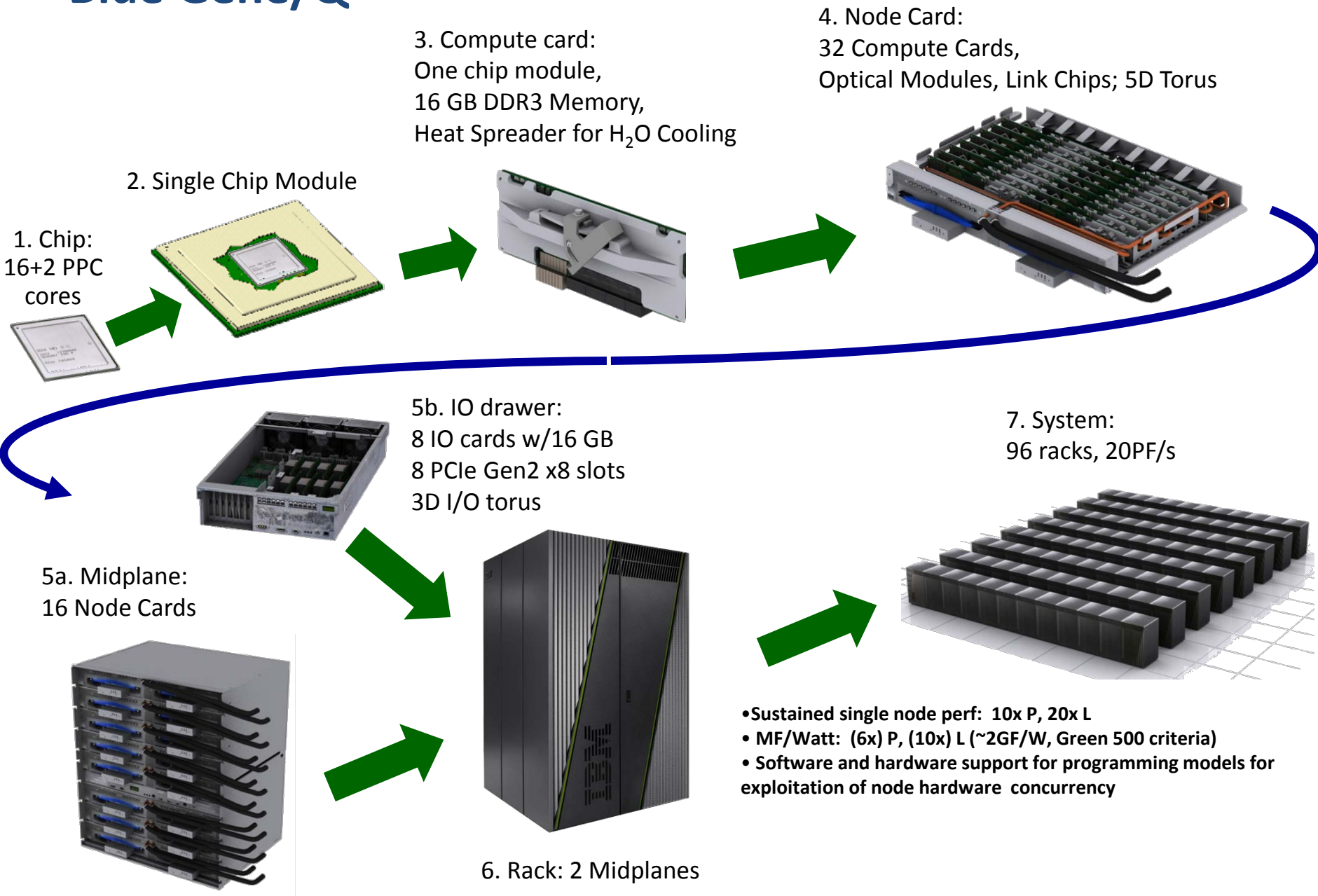
ALCF-2 architecture: not just Blue Gene/Q...



Configuration

- BG/Q
 - 48 racks
 - 48K 1.6 GHz nodes
 - 768K cores & 786TB RAM
 - 384 I/O nodes
- Storage
 - 240 GB/s, 35 PB

Blue Gene/Q



Systems Overview

	Vesta	Cetus	Mira
Racks	1	1	48
Purpose	General Test & Development	INCITE Testing	INCITE Production
IO Arrangement	128:1 top-hat wiring	128:1 top-hat wiring	128:1 IO racks
Login Nodes	1 Power 740 8 A2 IO nodes	1 Power 740 8 A2 IO nodes	4 Power 740s 16 A2 IO nodes

This is analogous to Surveyor, Challenger and Intrepid in the ALCF-1 (Blue Gene/P) ecosystem.

Having independent systems allow us to test system software changes before deploying them on the production research, to do “dangerous” computer science research, and to give users a nice sandbox to explore the architecture.



ALCF's BG/Q System - Mira



Mira Out Of The Box



Hardware:

- 48 racks – 10 Petaflops peak
- 48k compute nodes / 768k cores
 - PowerPC A2 with 16 cores/node at 1.6 GHz
 - 4 way multi-threading
 - Quad FPU – 4 wide double precision SIMD
- 384 IO nodes
- 768 TB Memory – 16 GB/node
- 5D Torus network
- Login, Service, & Management Nodes

Software:

- Compute Node Kernel OS + IO Node Linux
- XL Fortran, C, C++ Compilers
- GNU Compilers + GNU Toolchain
- Messaging Stack – SPI, PAMI, MPI
- ESSL (BLAS/LAPACK)
- Control System – MMCS, CIOD, CDTI

What's not in the box

The BG/Q software stack is a good start but users want more, including:

- Performance Tools:
 - PAPI, TAU, HPCToolkit, mpiP,
- Debuggers:
 - TotalView, DDT
- Libraries:
 - FFTW, ScaLAPACK, ParMetis, P3DFFT
- Programming Model:
 - UPC, Charm++, Global Arrays, ...



A collaborative “Early Software” project was initiated at Argonne to make this software available when Mira comes online

Building the Blue Gene/Q scientific software ecosystem

IBM has provided OSS all the way down:

- MPI sits on top of PAMI sits on top of SPI, all are OSS
- CNK, which is Linux-like, is OSS
- Python, GCC, glibc, etc. all OSS (obviously, due to GPL for GNU codes)

What IBM does not provide is readily ported by learning from these.

For example, GASNet ported by learning PAMI and MPI design.

We are working to enable a seamless transition between your laptop and Blue Gene/Q.

Why not just build an x86 supercomputer?

Listed below are the The Green500's Top 10 most energy-efficient supercomputers in the world as of June 2012.

[Column Separated Values \(CSV\)](#)
[Microsoft Excel Spreadsheet \(XLS\)](#)

Green500 Rank	MFLOPS/W	Site*	Computer*	Total Power (kW)
1	2100.88	DOE/NNSA/LLNL	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	41.10
2	2100.88	IBM Thomas J. Watson Research Center	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	41.10
3	2100.86	DOE/SC/Argonne National Laboratory	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	82.20
4	2100.86	DOE/SC/Argonne National Laboratory	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	82.20
5	2100.86	Rensselaer Polytechnic Institute	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	82.20
6	2100.86	University of Rochester	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	82.20
7	2100.86	IBM Thomas J. Watson Research Center	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	82.20
8	2099.56	University of Edinburgh	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	493.10
9	2099.50	Science and Technology Facilities Council - Daresbury Laboratory	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	575.30
10	2099.46	Forschungszentrum Juelich (FZJ)	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	657.50

* Performance data obtained from publicly available sources including [TOP500](#)

<u>11</u>	2099.39	CINECA	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	821.90
<u>12</u>	2099.14	High Energy Accelerator Research Organization /KEK	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	246.60
<u>13</u>	2099.14	EDF R&D	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	328.80
<u>14</u>	2099.14	IDRIS/GENCI	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	328.80
<u>15</u>	2099.14	Victorian Life Sciences Computation Initiative	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	328.80
<u>16</u>	2099.14	IBM - Rochester	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom	164.40
<u>17</u>	2099.14	IBM - Rochester	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom	164.40
<u>18</u>	2099.14	DOE/NNSA/LLNL	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	164.40
<u>19</u>	2069.04	DOE/SC/Argonne National Laboratory	BlueGene/Q, Power BQC 16C 1.60GHz, Custom	3945.00
<u>20</u>	2069.04	DOE/NNSA/LLNL	BlueGene/Q, Power BQC 16C 1.60 GHz, Custom	7890.00
<u>21</u>	1380.67	Intel	Intel Cluster, Xeon E5-2680 8C 2.600GHz, Infiniband FDR, Intel MIC	72.90
<u>22</u>	1379.79	Nagasaki University	DEGIMA Cluster, Intel i5, ATI Radeon GPU, Infiniband QDR	47.00
<u>23</u>	1266.26	Barcelona Supercomputing Center	Bullx B505, Xeon E5649 6C 2.53GHz, Infiniband QDR, NVIDIA 2090	81.50
<u>24</u>	1151.91	Center for Computational Sciences, University of Tsukuba	Xtream-X GreenBlade 604, Xeon E5-2670 8C 2.600GHz, Infiniband QDR, NVIDIA 2090	366.00
<u>25</u>	1050.26	Los Alamos National Laboratory	Xtreme-X , Xeon E5-2670 8C 2.600GHz, Infiniband QDR, NVIDIA 2090	226.80

<u>26</u>	1010.11	CEA/TGCC-GENCI	Bullx B505, Xeon E5640 2.67 GHz, Infiniband QDR	108.80
<u>27</u>	994.14	CSIRO	Nitro G16 3GPU, Xeon E5-2650 8C 2.000GHz, Infiniband FDR, NVIDIA 2050	110.30
<u>28</u>	958.35	GSIC Center, Tokyo Institute of Technology	HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU, Linux/Windows	1243.80
<u>29</u>	953.29	Virginia Tech	SuperServer 2026GT-TRF, Xeon E5645 6C 2.40GHz, Infiniband QDR, NVIDIA 2050	126.30
<u>30</u>	932.83	University of Chicago	iDataPlex DX360M4, Xeon E5-2670 8C 2.600GHz, Infiniband FDR	73.10
<u>31</u>	932.78	Centro Euro-Mediterraneo per i Cambiamenti Climatici	iDataPlex DX360M4, Xeon E5-2670 8C 2.600GHz, Infiniband FDR	159.50
<u>32</u>	932.71	Indian Institute of Technology Madras	iDataPlex DX360M4, Xeon E5-2670 8C 2.600GHz, Infiniband FDR	97.70
<u>33</u>	932.62	Automotive Company	iDataPlex DX360M4, Xeon E5-2670 8C 2.600GHz, Infiniband FDR	67.30
<u>34</u>	932.53	Durham University	iDataPlex DX360M4, Xeon E5-2670 8C 2.600GHz, Infiniband FDR	139.60
<u>35</u>	932.41	Science and Technology Facilities Council - Daresbury Laboratory	iDataPlex DX360M4, Xeon E5-2670 8C 2.600GHz, Infiniband FDR	170.20
<u>36</u>	932.19	Universidad de Cantabria - SSC	iDataPlex DX360M4, Xeon E5-2670 8C 2.600GHz, Infiniband FDR	79.80
<u>37</u>	919.44	Institute of Process Engineering, Chinese Academy of Sciences	Mole-8.5 Cluster, Xeon X5520 4C 2.27 GHz, Infiniband QDR, NVIDIA 2050	540.00
<u>38</u>	910.90	Max-Planck-Gesellschaft MPI/IPP	iDataPlex DX360M4, Xeon E5-2670 8C 2.600GHz, Infiniband FDR	205.70
<u>39</u>	908.83	Leibniz Rechenzentrum	iDataPlex DX360M4, Xeon E5-2680 8C 2.70GHz, Infiniband FDR	2841.00
<u>40</u>	901.61	Georgia Institute of Technology	HP ProLiant SL390s G7 Xeon 6C X5660 2.8Ghz, nVidia Fermi, Infiniband QDR	117.90

119	410.51	Interdisciplinary Centre for Mathematical and Computational Modelling, University of Warsaw	Power 775, POWER7 8C 3.836GHz, Custom	156.70
120	410.44	IBM Poughkeepsie Benchmarking Center	Power 775, Power7 3.836 GHz	172.40
121	408.97	DOE/Bettis Atomic Power Laboratory	Atipa Cluster, Xeon X56xx (Westmere-EP) 2.66 GHz, Infiniband QDR	163.02
122	408.97	Knolls Atomic Power Laboratory	Atipa Cluster, Xeon X56xx (Westmere-EP) 2.66 GHz, Infiniband QDR	163.02
123	406.87	Slovak Academy of Sciences (SAV)	Power 775, POWER7 8C 3.836GHz, Custom	188.10
124	405.38	Swiss Scientific Computing Center (CSCS)	Cray XE6, Opteron 6272 16C 2.10 GHz, Cray Gemini interconnect	780.00
125	404.45	CLUMEQ - McGill University	iDataPlex DX360 M3, Xeon 2.66, Infiniband	337.00
126	398.07	United Kingdom Meteorological Office	Power 775, POWER7 8C 3.836GHz, Custom	313.40
127	378.77	King Abdullah University of Science and Technology	Blue Gene/P Solution	504.00
128	378.77	EDF R&D	Blue Gene/P Solution	252.00
129	378.76	IDRIS	Blue Gene/P Solution	315.00
130	377.99	Rice University	BlueGene/P, PowerPC 450 4C 850 MHz, Proprietary	189.00
131	377.48	DOE/SC/Oak Ridge National Laboratory	Cray XK6, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA 2090	5142.00
132	377.22	IBM Poughkeepsie Benchmarking Center	Power 775, Power7 3.836 GHz	423.10
133	369.09	Environment Canada	Power 775, POWER7 8C 3.84 GHz, Custom	501.50
134	369.09	Environment Canada	Power 775, POWER7 8C 3.84 GHz, Custom	501.50
135	368.64	Automotive	BladeCenter HS22 Cluster, WM Xeon 6-core 2.93Ghz, Infiniband	240.40

166	318.69	NOAA/Oak Ridge National Laboratory	Cray XT6-HE, Opteron 6100 12C 2.1GHz	610.00
167	311.76	NASA/Ames Research Center/NAS	SGI Altix ICE X/8200EX/8400EX, Xeon 54xx 3.0/5570/5670/E5-2670 2.93/2.6/3.06/3.0 Ghz, Infiniband QDR/FDR	3987.00
168	305.33	Electronics	xSeries x3550M3 Cluster, Xeon X5650 6C 2.66 GHz, Infiniband QDR	251.10
169	305.33	Electronics	xSeries x3550M3 Cluster, Xeon X5650 6C 2.66 GHz, Infiniband QDR	251.10
170	305.33	Electronics	xSeries x3550M3 Cluster, Xeon X5650 6C 2.66 GHz, Infiniband QDR	241.80
171	297.44	National Institute for Computational Sciences/University of Tennessee	Cray XT5-HE Opteron Six Core 2.6 GHz	3090.00
172	297.26	National Computational Infrastructure National Facility (NCI-NF)	Sun Blade x6048, Xeon X5570 2.93 Ghz, Infiniband QDR	425.22
173	288.73	Lawrence Livermore National Laboratory	Xtreme-X , Xeon E5-2670 8C 2.600GHz, Infiniband QDR	1203.20
174	286.74	Defence	BladeCenter HX5, Xeon E7-4870 10C 2.40 GHz, Infiniband QDR	226.20
175	281.56	IT Service Provider	Cluster Platform 3000 280c G7, Xeon X5672 4C 3.20GHz, Infiniband QDR	216.03
176	278.89	DOE/NNSA/LANL/SNL	Cray XE6, Opteron 6136 8C 2.40GHz, Custom	3980.00
177	277.76	Lawrence Livermore National Laboratory	Dell DCS Xanadu 2.5, Xeon E55xx 2.4Ghz, Infiniband DDR	260.69
178	277.29	University of Alaska - Arctic Region Supercomputing Center	Cray XE6 8-core 2.3 GHz	320.68
179	271.71	Centre for High Performance Computing	Blade X6275/ PowerEdge C6100 Cluster, Xeon X5570/X5670/ 4C 2.93 GHz, Infiniband QDR	225.72
180	268.22	Cray Inc.	Cray XE6, Opteron 6272 16C 2.1/2.2/2.3 GHz, Cray Gemini interconnect	467.90

Rank	Site	Computer/Year Vendor	Cores	R _{max}	R _{peak}	Power
1	DOE/NNSA/LLNL United States	Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom / 2011 IBM	1572864	16324.75	20132.66	7890.0
2	RIKEN Advanced Institute for Computational Science (AICS) Japan	K computer , SPARC64 VIIIfx 2.0GHz, Tofu interconnect / 2011 Fujitsu	705024	10510.00	11280.38	12659.9
3	DOE/SC/Argonne National Laboratory United States	Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM	786432	8162.38	10066.33	3945.0
4	Leibniz Rechenzentrum Germany	SuperMUC - iDataPlex DX360M4, Xeon E5- 2680 8C 2.70GHz, Infiniband FDR / 2012 IBM	147456	2897.00	3185.05	3422.7
5	National Supercomputing Center in Tianjin China	Tianhe-1A - NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010 NUDT	186368	2566.00	4701.00	4040.0
6	DOE/SC/Oak Ridge National Laboratory United States	Jaguar - Cray XK6, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA 2090 / 2009 Cray Inc.	298592	1941.00	2627.61	5142.0
7	CINECA Italy	Fermi - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM	163840	1725.49	2097.15	821.9
8	Forschungszentrum Juelich (FZJ) Germany	JuQUEEN - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM	131072	1380.39	1677.72	657.5
9	CEA/TGCC-GENCI France	Curie thin nodes - Bullx B510, Xeon E5- 2680 8C 2.700GHz, Infiniband QDR / 2012 Bull	77184	1359.00	1667.17	2251.0
10	National Supercomputing Centre in Shenzhen (NSCS) China	Nebulae - Dawning TC3600 Blade System, Xeon X5650 6C 2.66GHz, Infiniband QDR, NVIDIA 2050 / 2010 Dawning	120640	1271.00	2984.30	2580.0
11	NASA/Ames Research Center/NAS United States	Pleiades - SGI Altix ICE X/8200EX/8400EX, Xeon 54xx 3.0/5570/5670/E5-2670 2.93/2.6/3.06/3.0 Ghz, Infiniband QDR/FDR / 2011 SGI	125980	1243.00	1731.84	3987.0
12	International Fusion Energy Research Centre (IFERC), EU(F4E) - Japan Broader Approach collaboration Japan	Helios - Bullx B510, Xeon E5-2680 8C 2.700GHz, Infiniband QDR / 2011 Bull	70560	1237.00	1524.10	2200.0

Rank	Installation Site	Machine	Number of nodes	Number of cores	Problem scale	GTEPS
1	DOE/SC/Argonne National Laboratory	Mira/BlueGene/Q	32768	524288	38	3541.00
1	LLNL	Sequoia/Blue Gene/Q	32768	524288	38	3541.00
2	DARPA Trial Subset, IBM Development Engineering	Power 775, POWER7 8C 3.836 GHz	1024	32768	35	508.05
3	Information Technology Center, The University of Tokyo	Oakleaf-FX (Fujitsu PRIMEHPC FX 10)	4800	76800	38	358.10
4	GSIC Center, Tokyo Institute of Technology	HP Cluster Platform SL390s G7 (three Tesla cards per node)	1366	16392	35	317.09
5	Brookhaven National Laboratory	BLUE GENE/Q	1024	16384	34	294.29
6	DOE/SC/Argonne National Laboratory	Vesta/BlueGene/Q	1024	16384	34	292.36
7	NASA-Ames / Parallel Computing Lab, Intel Labs	Pleiades - SGI ICE-X, dual plane hypercube FDR infiniband, E5-2670 "sandybridge"	1024	16384	34	270.33
8	NERSC/LBNL	XE6	4817	115600	35	254.07
9	NNSA and IBM Research, T.J. Watson	NNSA/SC Blue Gene/Q Prototype II	4096	65536	32	236.00
10	GSIC Center, Tokyo Institute of Technology	TSUBAME 2.0 (CPU only)	1366	16392	36	202.68
11	Intel Dupont / Parallel Computing Lab, Intel Labs	Endeavor; Dual-socket 2.6GHz SNB-EP	320	5120	32	115.94
12	LBL	Hopper	1800	43200	37	105.00

There is no such thing as an x86 supercomputer

Power is the limiting factor and traditional machines like Jaguar are 10 times less power efficient than Blue Gene/Q.

Intel is moving towards x86+MIC. MIC has the x86 toolchain but it is still an attached coprocessor and is not binary compatible with x86 CPUs.

AMD is moving towards CPU+APU. No one uses this hardware in HPC...

NVIDIA provides a disruptive software development environment and has all the issues associated with a heterogeneous system.

Even when platforms resemble commodity hardware, they still require extensive custom engineering to scale. Cray Gemini is substantially more scalable than Infiniband, for example.

IBM PERCS (formerly Blue Waters), is much less efficient than BG/Q...

What do you want for development?

The same:

- Compilers
- MPI
- Performance tools
- IDE
- Languages
- What else?

We have:

- GCC, LLVM
- MPICH2
- TAU, HPCToolkit, PAPI
- Eclipse
- Fortran, C/C++, Python

Do you take it for granted that you can use Python?

Do the portable compilers – i.e. not IBM XL or Cray – support the features?

OSS Compilers on BG/Q

ALCF has been frustrated by the compiler options on Blue Gene/P.

LLVM is an OSS compiler project driven by Apple and used by NVIDIA, Cray, Google, etc.

Hal Finkel (CSGF@ALCF) ported LLVM to Blue Gene/P and Blue Gene/Q and it supports features previously only available from the IBM compiler (auto and manual vectorization). (switch slide decks)

By the time you have access to BG/Q, we will have an implementation of QPX intrinsics in LLVM so that you can write explicitly vectorized code on your laptop and run it locally.

Collaborators are working on transactional memory in LLVM.

Should be able to have TM and SE support in GCC (since these are OpenMP oriented)

BG/Q Tools status

Tool Name	Source	Provides	Status
bgpm	IBM	HPC	Available in GA driver
gprof	GNU/IBM	Timing (sample)	Available in GA driver
TAU	Univ. Oregon	Timing (inst), MPI	Ported and available on BG/Q
HPCToolkit	Rice University	Timing (sample), HPC (sample)	Ported and available on BG/Q
IBM HPCT	IBM	MPI, HPC	IBM product, preliminary version available on BG/Q
mpiP	LLNL	MPI	Beta version available on BG/Q
PAPI	UTK	HPC API	Ported and available on BG/Q
Darshan	ANL	IO	Ported and available on BG/Q
Open Speedshop	Krell Institute	Timing (sample), HCP, MPI, IO	Preliminary version available on BG/Q
MPE/Jumpshot	Argonne	MPI	Ported and available on BG/Q
Scalasca	Juelich	Timing (inst), MPI	Ported and available on BG/Q
DynInst	UMD/Wisc./IBM	Binary rewriter	Preliminary version available on BG/Q
ValGrind	ValGrind/IBM	Memory & Thread Error Check	Development planned for summer

ALCF Supported BG/Q Libraries on VEAS

- Located in `/soft/libraries/alcf`
- Maintained in-house, frequently updated
- GNU and XL built versions of each library

Library Name	Source	Provides	VEAS status
ESSL	IBM	Dense Linear Algebra & FFT Kernels	5.1.1-0 beta version
BLAS	NETLIB (UTK) & ESSL	Dense Linear Algebra Kernels	Veas version is based on ESSL GEMM. All level-3 routines are GEMM-based.
BLIS	Univ. Texas & ANL	Framework for a successor to GotoBLAS	Under design; plans for hand-tuned BG/Q version
CBLAS	UTK	C wrappers to BLAS	On Veas
LAPACK	UTK	Dense Linear Algebra Solver	3.4.1
ScaLAPACK	UTK	Parallel Dense Linear Algebra Solver	2.0.2
ARPACK & PARPACK	Rice Univ.	Eigenvalues & Eigenvectors	2.1 (last version released in 2001)
FFTW	MIT	Fast Fourier Transform	2.1.5, 3.3.1. No hand-tuning for BG/Q yet.
METIS	UMN	Graph Partitioning	5.0.2
ParMETIS	UMN	Graph Partitioning	4.0.2