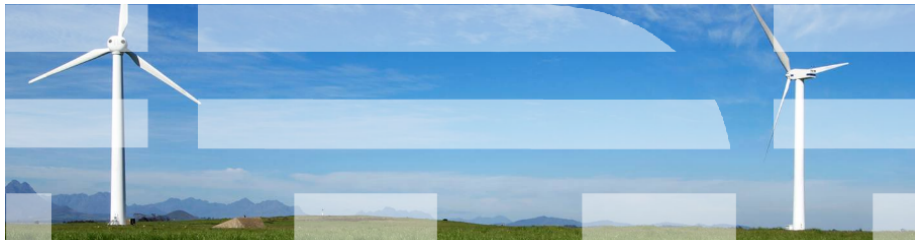


Introduction to the Practical Aspects of Programming for IBM Blue Gene/P

Alexander Pozdneev

Research Software Engineer, IBM

June 22, 2015 — International Summer Supercomputing Academy



IBM Science and Technology Center

Established in 2006

- Research projects
 - ▶ HPC simulations
 - ▶ Security
 - ▶ Oil and gas applications
- ESSL math library for IBM POWER
- IBM z Systems mainframes
 - ▶ Firmware
 - ▶ Linux on z Systems
 - ▶ Software
- IBM Rational software
- Smarter Commerce for Retail
- Smarter Cities



Outline

- 1 Blue Gene architecture in the HPC world
- 2 Blue Gene/P architecture
- 3 Conclusion
- 4 References

Outline

1 Blue Gene architecture in the HPC world

- Brief history of Blue Gene architecture
- Blue Gene in TOP500
- Blue Gene in Graph500
- Scientific applications
- Gordon Bell awards
- Blue Gene installation sites

2 Blue Gene/P architecture

3 Conclusion

4 References

Brief history of Blue Gene architecture

- **1999** — US\$100M research initiative, novel massively parallel architecture, protein folding
- **2003, Nov** — Blue Gene/L first prototype — TOP500, #73
- **2004, Nov** — 16 Blue Gene/L racks at LLNL — TOP500, #1
- **2007** — Second generation — Blue Gene/P
- **2009** — National Medal of Technology and Innovation
- **2012** — Third generation — Blue Gene/Q
- Features:
 - ▶ Low frequency, low power consumption
 - ▶ Fast interconnect balanced with CPU
 - ▶ Light-weight OS



Blue Gene in TOP500

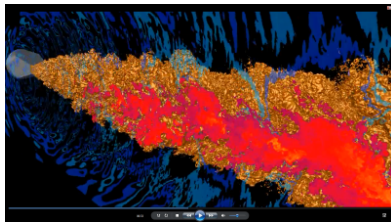
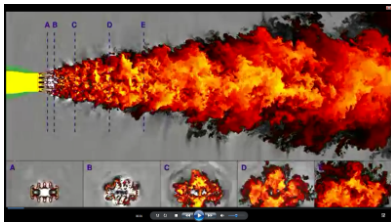
Date	#	System	Model	Nodes	Cores
Jun'14 / Nov'14	3	Sequoia	Q	96k	1.5M
Jun'13 / Nov'13	3	Sequoia	Q	96k	1.5M
Nov'12	2	Sequoia	Q	96k	1.5M
Jun'12	1	Sequoia	Q	96k	1.5M
Nov'11	13	Jugene	P	72k	288k
Jun'11	12	Jugene	P	72k	288k
Nov'10	9	Jugene	P	72k	288k
Jun'10	5	Jugene	P	72k	288k
Nov'09	4	Jugene	P	72k	288k
Jun'09	3	Jugene	P	72k	288k
Nov'08	4	DOE/NNSA/LLNL	L	104k	208k
Jun'08	2	DOE/NNSA/LLNL	L	104k	208k
Nov'07	1	DOE/NNSA/LLNL	L	104k	208k
Jun'07	1	DOE/NNSA/LLNL	L	64k	128k
Jun'06 / Nov'06	1	DOE/NNSA/LLNL	L	64k	128k
Nov'05	1	DOE/NNSA/LLNL	L	64k	128k
Jun'05	1	DOE/NNSA/LLNL	L	32k	64k
Nov'04	1	DD2 beta-system	L	16k	32k
Jun'04	4	DD1 prototype	L	4k	8k

Blue Gene in Graph500

Date	#	System	Model	Nodes	Cores	Scale	GTEPS
Nov'14	1	Sequoia	Q	96k	1.5M	41	23751
Jun'14	2	Sequoia	Q	64k	1M	40	16599
Nov'13	1	Sequoia	Q	64k	1M	40	15363
Jun'13	1	Sequoia	Q	64k	1M	40	15363
Nov'12	1	Sequoia	Q	64k	1M	40	15363
Jun'12	1	Sequoia/Mira	Q	32k	512k	38	3541
Nov'11	1	BG/Q prototype	Q	4k	64k	32	253
Jun'11	1	Interpid/Jugene	P	32k	128k	38	18
Nov'10	1	Interpid	P	8k	32k	36	7

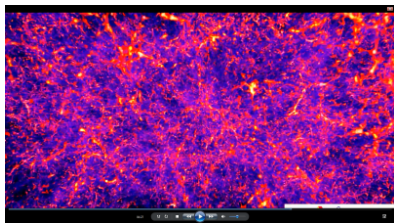
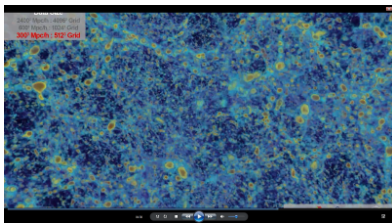
Jet Engine Noise CFD Simulation

- Supersonic jet noise simulation, effects of chevrons on jet noise
- 128k Blue Gene/P cores — \approx 100 hours
- 1M Blue Gene/Q cores — \approx 12 hours
- IBM Blue Gene Q Sequoia, Lawrence Livermore National Laboratory
- <https://youtu.be/cjoz5tncRUs>
- <https://youtu.be/uxT-VmY30Wc>
- Video: Joseph W. Nichols, Stanford Center for Turbulence Research



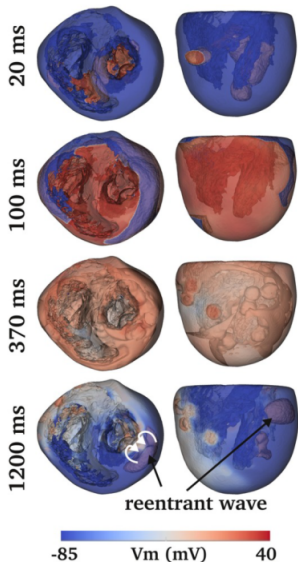
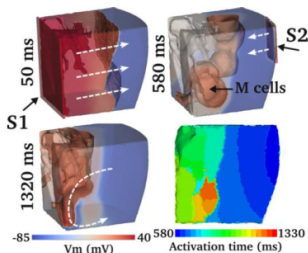
Secrets of the Dark Universe

- The evolution of the Universe simulation, understanding the physics of the dark matter and energy
- 1 BG/Q rack — 68B particles
- 32 BG/Q racks — 1.1T particles
- <http://www.youtube.com/watch?v=tdv8yrJk4VE>
- http://www.youtube.com/watch?v=-S-T_iTiAxQ
- Video: ANL, LANL, LBNL



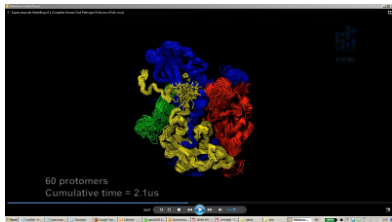
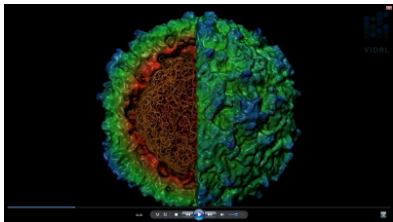
Real-Time Modeling of Human Heart Ventricles

- Simulation of drug-induced arrhythmias
- Resolution — 0.1 mm
- 768k Blue Gene/Q cores
- 43% peak
- <http://dl.acm.org/citation.cfm?id=2388999>
- LLNL, IBM Research, IBM Research Collaboratory for Life Sciences



Modelling of a Complete Human Viral Pathogen Poliovirus

- Reconstruction and simulation of poliovirus
- Antiviral drugs, virus infection, modelling related viruses
- 3.3M–3.7M atoms
- Blue Gene/Q, Victorian Life Sciences Computing Initiative
- <http://www.youtube.com/watch?v=Nih0Qa673FY>



Parallel Sparse Matrix–Matrix Multiplication

- Semiempirical Molecular Dynamics (SEMD) I: Midpoint-Based Parallel Sparse Matrix–Matrix Multiplication Algorithm for Matrices with Decay
- Inspired by the midpoint method
- Approaching a perfect linear scaling
- Up to 185 193 processes on BG/Q
- doi: 10.1021/acs.jctc.5b00382

ACS Publications
Most Trusted Most Cited Most Read

ACS Journals | ACS Chemistry | ACS eBooks

JCTC Journal of Chemical Theory and Computation

Search
Search text
J. Chem.

Home Browse the Journal Articles ASAP Current Issue Submission & Review Subscri

Article

Semiempirical Molecular Dynamics (SEMD) I: Midpoint-Based Parallel Sparse Matrix–Matrix Multiplication Algorithm for Matrices with Decay

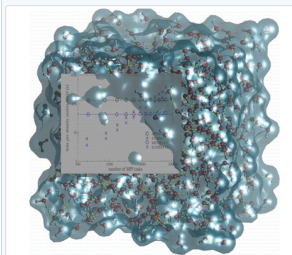
Václav Vřeber¹, Teodoro Laino¹, Alexander Probstheev¹, Irina Fedulova¹, and Alessandro Curioni¹

¹ IBM Research–Zurich, Säumerstrasse 4, 8803 Rüschlikon, Switzerland
² IBM Systems Lab Services Russia/CIS, Priesnenskaya Nab., 10, 123317, Moscow, Russia
[†] IBM Science and Technology Center, Priesnenskaya Nab., 10, 123317, Moscow, Russia

J. Chem. Theory Comput., Article ASAP
 DOI: 10.1021/acs.jctc.5b00382
 Publication Date (Web): June 3, 2015
 Copyright © 2015 American Chemical Society

*Tel.: +41 (0) 44 724 8200; Fax: +41 (0) 44 724 8508; E-mail: vve@zurich.ibm.com.

Abstract

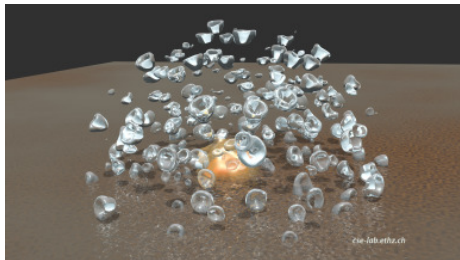


Gordon Bell prize — 7 awards

- **2005** — 100+ TFlop Solidification Simulations on Blue Gene/L
- **2006** — The Blue Gene/L supercomputer and quantum ChromoDynamics
- **2006** — Large-scale Electronic Structure Calculations of High-Z Metals on the Blue Gene/L Platform
- **2007** — Extending stability beyond CPU millennium: a micron-scale atomistic simulation of Kelvin-Helmholtz instability
- **2008** — Linearly scaling 3D fragment method for large-scale electronic structure calculations
- **2009** — The cat is out of the bag: cortical simulations with 10^9 neurons, 10^{13} synapses
- **2013** — 11 PFLOP/s simulations of cloud cavitation collapse

11 PFLOP/s simulations of cloud cavitation collapse

- ACM Gordon Bell Prize 2013
- IBM Blue Gene/Q Sequoia
 - ▶ 96 racks, 20.1 PFLOPS peak
 - ▶ 1.5M cores (6.0M threads)
 - ▶ 73% peak (14.4 PFLOPS)
 - ▶ 13T grid points
 - ▶ 15k cavitation bubbles
- Destructive capabilities of collapsing bubbles
 - ▶ high pressure fuel injectors and propellers
 - ▶ shattering kidney stones
 - ▶ cancer treatment
- doi: 10.1145/2503210.2504565
- https://youtu.be/zfG9solC6_Y



- IBM Research Zurich
- ETH Zurich
- Technical University of Munich
- Lawrence Livermore National Laboratory (LLNL)

Some of Blue Gene installation sites

- USA:
 - ▶ DOE/NNSA/LLNL (2)
 - ▶ DOE/SC/ANL (3)
 - ▶ University of Rochester
 - ▶ Rensselaer Polytechnic Institute
 - ▶ IBM Rochester / T.J. Watson (4)
- Germany:
 - ▶ Forschungszentrum Juelich (FZJ)
- Italy:
 - ▶ CINECA
- Switzerland:
 - ▶ Ecole Polytechnique Federale de Lausanne
 - ▶ Swiss National Supercomputing Centre (CSCS)
- Poland:
 - ▶ Interdisciplinary Centre for Mathematical and Computational Modelling, University of Warsaw
- France:
 - ▶ CNRS/IDRIS-GENCI
 - ▶ EDF R&D
- UK:
 - ▶ University of Edinburgh
 - ▶ Science and Technology Facilities Council — Daresbury Laboratory
- Japan:
 - ▶ High Energy Accelerator Research Organization/KEK (2)
- Australia:
 - ▶ Victorian Life Sciences Computation Initiative
- Canada:
 - ▶ Southern Ontario Smart Computing Innovation Consortium/University of Toronto

Outline

1 Blue Gene architecture in the HPC world

2 Blue Gene/P architecture

- Air-cooling

- Hardware overview

- BG/P at Lomonosov MSU

- Computing node

- Execution process modes

- Memory subsystem

- Networks

- Compilers

- Software

- Software stack

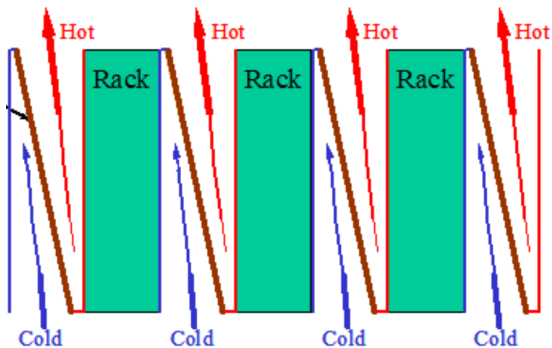
- I/O nodes psets

- Double Hammer FPU

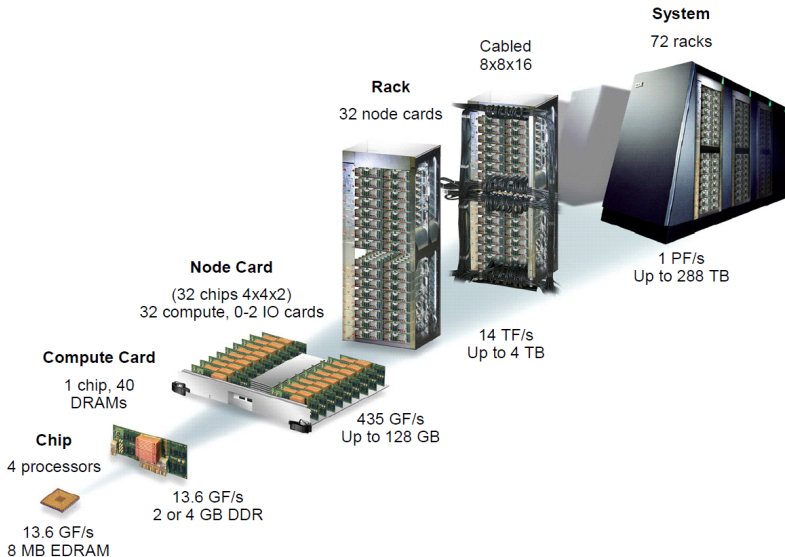
3 Conclusion

4 References

Blue Gene/P air-cooling



Blue Gene/P hardware overview

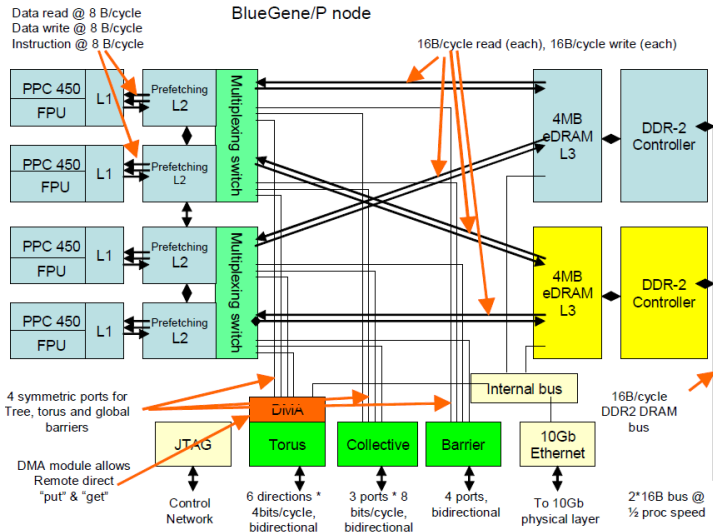


Blue Gene/P at Lomonosov MSU

- Two racks
- 4 TB RAM
- 27.2 TFLOPS peak
- 23.2 TFLOPS on LINPACK (85% peak)

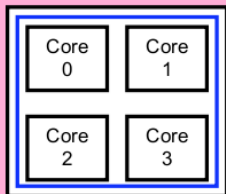


Blue Gene/P computing node

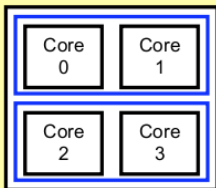


Execution process modes

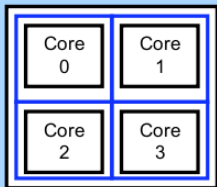
SMP Mode
1 Process
1-4 Threads/Process



Dual Mode
2 Processes
1-2 Threads/Process



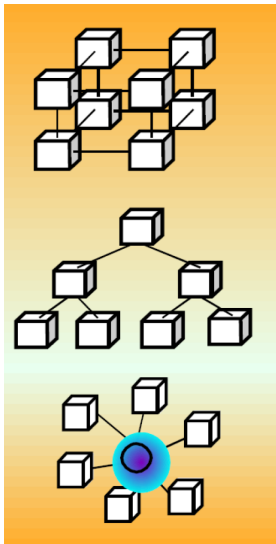
Quad Mode (VNM)
4 Processes
1 Thread/Process



Blue Gene/P memory subsystem

Cache	Total per node	Size	Replacement policy	Associativity
L1 instruction	4	32 KB	Round-Robin	<ul style="list-style-type: none"> ▶ 64-way set-associative ▶ 16 sets ▶ 32-byte line size
L1 data	4	32 KB	Round-Robin	<ul style="list-style-type: none"> ▶ 64-way set-associative ▶ 16 sets ▶ 32-byte line size
L2 prefetch	4	14 x 256 bytes	Round-Robin	<ul style="list-style-type: none"> ▶ Fully associative (15-way) ▶ 128-byte line size
L3	2	2 x 4 MB	Least recently used	<ul style="list-style-type: none"> ▶ 8-way associative ▶ 2 bank interleaved ▶ 128-byte line size
Double data RAM (DDR)	2	<ul style="list-style-type: none"> ▶ Minimum 2 x 512 MB ▶ Maximum 4 GB 	N/A	<ul style="list-style-type: none"> ▶ 128-byte line size

Blue Gene/P networks



- 3D torus
 - ▶ 3.4 Gbit/s per each of 12 ports (5.1 GB/s per node)
 - ▶ Hardware latency: 0.5 ms (nearest neighbours), 5 ms (round-trip)
 - ▶ MPI latency: 3 ms, 10 ms
 - ▶ Main communication network
 - ▶ Many collectives take benefit of it
- Collectives (tree)
 - ▶ Global one-to-all communications (broadcast, reduction)
 - ▶ 6.8 Gbit/s per port
 - ▶ Round-trip latency: 1.3 ms (hardware), 5 ms (MPI)
 - ▶ Connects computing nodes and I/O nodes
 - ▶ Collectives on `MPI_COMM_WORLD`
- Global interrupts
 - ▶ Worst-case latency: 0.65 ms (hardware), 5 ms (MPI)
 - ▶ `MPI_Barrier()`

Blue Gene/P compilers

Compiler		C/C++		Fortran			
		C	C++	Fortran 77	Fortran 90	Fortran 95	Fortran 2003
GCC		mpicc	mpicxx	mpif77	—	—	—
IBM XL	non-threaded	mpixlc	mpixlcxx	mpixlf77	mpixlf90	mpixlf95	mpixlf2003
	thread-safe	mpixlc_r	mpixlcxx_r	mpixlf77_r	mpixlf90_r	mpixlf95_r	mpixlf2003_r

IBM XL compilers flags

- `-qarch=450 -qtune=450`
 - ▶ exploit only one of two FPUs
 - ▶ use when data is not 16-byte aligned
- `-qarch=450d -qtune=450`
 - ▶ alignment!
- `-O3 (-qstrict)`
 - ▶ minimal optimization level for Double Hammer FPU
- `-O3 -qhot`
- `-O4 (-qnoipa)`
- `-O5`
- `-qreport -qlist -qsource`
 - ▶ a lot of useful information in `.lst` file
- `-qsmp=omp, -qsmp=auto`

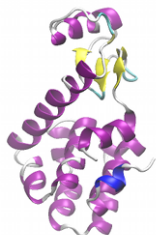
Blue Gene/P software

System software

- IBM XL and GCC compilers
- ESSL mathematical library (BLAS, LAPACK, etc.)
- MASS mathematical library
- LoadLeveler job scheduler

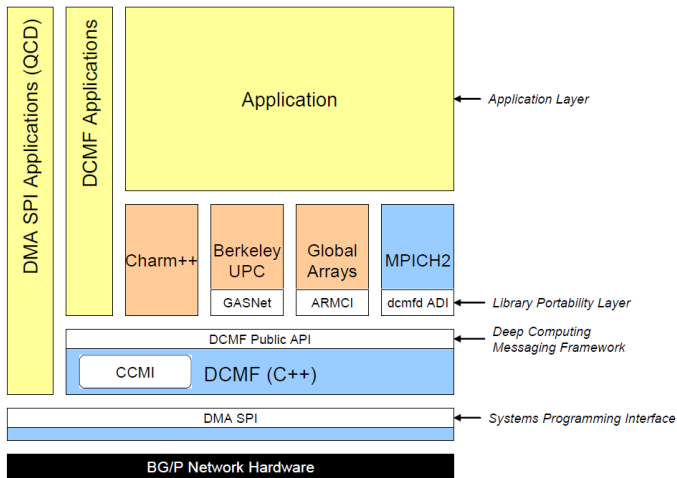
Application software and libraries

- FFTW
- GROMACS
- METIS
- ParMETIS
- szip
- zlib



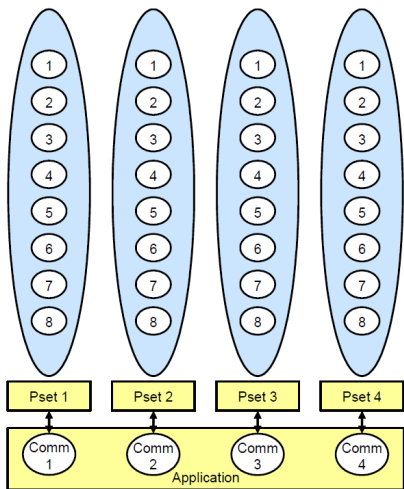
<http://hpc.cmc.msu.ru/bgp/soft>

Blue Gene/P software stack

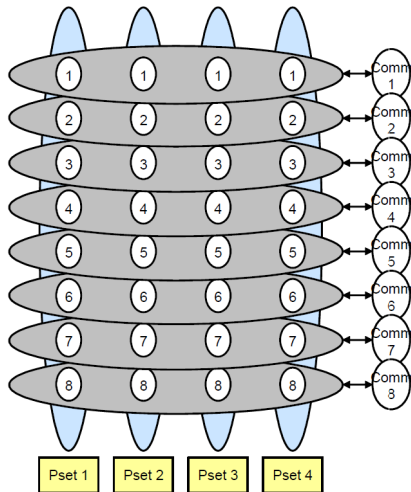


- IBM supported software
- Externally supported software

Grouping processes in respect to I/O nodes

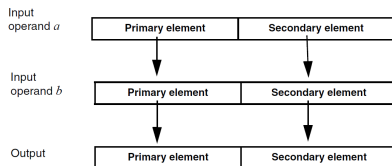
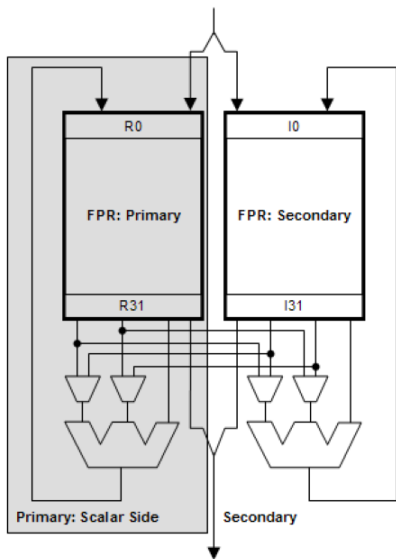


```
int MPIX_Pset_same_comm_create (MPI_Comm *pset_comm);
```



```
int MPIX_Pset_diff_comm_create (MPI_Comm *pset_comm);
```

Double Hammer FPU



```
#define MY_ALIGN __attribute__((aligned(16)))
```

```
int main() {
    double MY_ALIGN a[2] = { 2, 3 };
    double MY_ALIGN b[2] = { 4, 5 };
    double MY_ALIGN c[2] = { 6, 7 };
    double MY_ALIGN r[2]; // 2 * 4 + 6 = 14
                        // 3 * 5 + 7 = 22

    double _Complex ra = __lfpd(a);
    double _Complex rb = __lfpd(b);
    double _Complex rc = __lfpd(c);

    double _Complex rr = __fpmadd(rc, rb, ra);

    __stfpd(r, rr);
}
```

Outline

- 1 Blue Gene architecture in the HPC world
- 2 Blue Gene/P architecture
- 3 Conclusion
- 4 References

Conclusion

- Brief history of Blue Gene
 - ▶ BG/L, BG/P, BG/Q
 - ▶ National Medal of Technology and Innovation
 - ▶ TOP500, Graph500
 - ▶ Seven Gordon Bell awards
- Blue Gene/P architecture
 - ▶ Four cores per node
 - ▶ Two GB RAM per node
 - ▶ 1024 nodes per rack
 - ▶ Two BG/P rack at Lomonosov MSU
 - ▶ Compilers and software

Outline

- 1 Blue Gene architecture in the HPC world
- 2 Blue Gene/P architecture
- 3 Conclusion
- 4 References

References



IBM System Blue Gene Solution: Blue Gene/P Application Development

Understand the Blue Gene/P programming environment

Learn how to run and debug MPI programs

Learn about Bridge and Real-time APIs



Carlos Sosa
Branit Knudsen

Redbooks

ibm.com/redbooks

<http://www.redbooks.ibm.com/abstracts/sg247287.html>

The screenshot shows the website hpc@cmc with the following content:

- Системные системы**
 - Blue Gene/P
 - Ресурсы
- Справочная информация**
 - Blue Gene/P
 - Новые возможности
 - Внешний вид
 - Подключение
 - Настройка программ
 - Запуск машин
 - Загрузка информации
 - Образование документации
 - FAQ
 - Служба Blue Gene в мире
 - Интервью
 - Публикации
 - Ресурсы
 - Новости
- Тема CNC**
 - Справка по установке
 - Новости
 - Загрузка информации
 - Публикации
- Общие вопросы**
 - Доступ на IBM
 - Выбор машин
 - API-интерфейсы
- Ресурсы**
 - Ресурсы загрузки
 - Тренинг-материал
 - Вопросы/ответы
 - Справочная служба
- Поддержка**
 - Служба
- Полезные ссылки**
 - Векторный компьютерный программист и инженер IBM
 - Настройка суперкомпьютера
 - Суперкомпьютерный комплекс
 - Настройка суперкомпьютера
 - развития
 - Дизайн суперкомпьютерных архитектур

Суперкомпьютер IBM Blue Gene/P на факультете ВМК МГУ
 С 2008 года на Факультете ВМК МГУ имени М. В. Ломоносова работает суперкомпьютер IBM Blue Gene/P, который является одной из первых систем данной серии, созданной специально в мире. Архитектура Blue Gene была разработана компанией IBM в рамках проекта по созданию возможностей достижения новых рубежей в суперкомпьютере. Более тысячи машин данной серии в настоящее время работают на территории факультета в составе первой суперкомпьютерной среды TurboSP, а машина Blue Gene/P, установленная на ВМК МГУ, в раздате работы от 18 ноября 2008 года оказалась на 120-м месте (в раздате от 16 ноября 2009 года – 348-м месте). В список самых высокопроизводительных суперкомпьютеров стран СНГ, опубликованный 22 октября 2009 года, она вошла на 4-й позиции.

Система IBM Blue Gene/P предоставляет и новую модель суперкомпьютера, обладающая высокой производительностью, масштабируемостью, возможностью обрабатывать данные большого объема, позволяя при этом значительно меньше времени и денег на единицу площади по сравнению с традиционными системами.

[Читать далее](#)

Система BlueGene/P доступна из интернета без OpenVPN
 Система BlueGene/P была переведена на доступ по SSH-соединению. Это решает проблему с простыми паролями, которые BlueGene/P не может принять и повлечет безопасность доступа.

В связи с переводом на ssh и отныне парольный доступ на систему теперь открыт по сети Интернет, а не только по сети факультета ВМК. Команда технической поддержки надеется, что вы найдете работу и простоту входа на систему.

Обновление прикладного программного обеспечения на Blue Gene/P
 На Blue Gene/P в частном режиме работает система модабл. В настоящее время через нее доступны следующие библиотеки и программы:

- FFTW – высокопроизводительная библиотека, реализующая функции быстрого преобразования Фурье
- OpenMPI – пакет для взаимодействия высокопроизводительных систем методами межсетевой доставки.

Свои замечания и предложения по установке и настройке дополнительных программного обеспечения вы можете сообщить службе технической поддержки факультета "Суперкомпьютеры".

<http://hpc.cmc.msu.ru>

Disclaimer

All the information, representations, statements, opinions and proposals in this document are correct and accurate to the best of our present knowledge but are not intended (and should not be taken) to be contractually binding unless and until they become the subject of separate, specific agreement between us.

Any IBM Machines provided are subject to the Statements of Limited Warranty accompanying the applicable Machine.

Any IBM Program Products provided are subject to their applicable license terms. Nothing herein, in whole or in part, shall be deemed to constitute a warranty.

IBM products are subject to withdrawal from marketing and or service upon notice, and changes to product configurations, or follow-on products, may result in price changes.

Any references in this document to “partner” or “partnership” do not constitute or imply a partnership in the sense of the Partnership Act 1890.

IBM is not responsible for printing errors in this proposal that result in pricing or information inaccuracies.

Правовая информация

IBM, логотип IBM, BladeCenter, System Storage и System x являются товарными знаками International Business Machines Corporation в США и/или других странах. Полный список товарных знаков компании IBM смотрите на узле Web: www.ibm.com/legal/copytrade.shtml.

Названия других компаний, продуктов и услуг могут являться товарными знаками или знаками обслуживания других компаний.

(c) 2015 International Business Machines Corporation. Все права защищены.

Упоминание в этой публикации продуктов или услуг корпорации IBM не означает, что IBM предполагает предоставлять их во всех странах, в которых осуществляет свою деятельность, информация о предоставлении продуктов или услуг может быть изменена без уведомления. За самой свежей информацией о продуктах и услугах компании IBM, предоставляемых в Вашем регионе, следует обращаться в ближайшее торговое представительство IBM или к авторизованным бизнес-партнерам.

Все заявления относительно намерений и перспективных планов IBM могут быть изменены без уведомления. Информация о продуктах третьих фирм получена от производителей этих продуктов или из опубликованных анонсов указанных продуктов. IBM не тестировала эти продукты и не может подтвердить производительность, совместимость, или любые другие заявления относительно продуктов третьих фирм.

Вопросы о возможностях продуктов третьих фирм следует адресовать поставщику этих продуктов.

Информация может содержать технические неточности или типографические ошибки. В представленную в публикации информацию могут вноситься изменения, эти изменения будут включаться в новые редакции данной публикации. IBM может вносить изменения в рассматриваемые в данной публикации продукты или услуги в любое время без уведомления.

Любые ссылки на узлы Web третьих фирм приведены только для удобства и никоим образом не служат поддержкой этим узлам Web. Материалы на указанных узлах Web не являются частью материалов для данного продукта IBM.