科	学	技	術	計	算	分	科	会		選	出
						SS 積	∄ НРС	フォー	・ラノ	2010 ב	より

海外招待講演

Extreme Computing: Challenges, Constraints and Opportunities

> Anne Trefethen Oxford University

Extreme Computing: Challenges, Constraints and Opportunities

Anne E Trefethen

Preamble

This paper is a handout to support the presentation at the Fujitsu HPC Users Group Meeting on 26th August, 2010.

Abstract

Scientific applications require an ecosystem of computational infrastructure that can allow distributed collaboration and computation, large-scale simulation and analysis and appropriate consideration of data driven science. The ecosystem requires the integration of high-performance computing, in some cases to the exascale, cloud computing, databases, high bandwidth networks and the software and people infrastructure that enables the effective use of the components.

In this presentation we will consider some of the challenges and constraints that drive the development of the computational infrastructure and its components – including computational models, energy and knowhow; and the opportunities that are presented, in terms of new science applications and new algorithmic approaches. There is a particular focus on energy.

Key words

Exascale, extreme computing, energy efficiency, cloud computing

Introduction

Advanced computing is an essential tool in addressing scientific problems of national and scientific interest, including climate change, the virtual human, new materials, next-generation power sources and astrophysics, but as importantly it is equally essential to solve commercial and industrial problems in financial modelling, engineering, and real-time decision systems.

Extreme computing has moved from being a somewhat esoteric interest of a few scientists to a necessity for any computational scientist or developer who uses simulation as a tool. A shift has occurred in recent years as the processor chip designs mimic the architecture of high-performance computers, with multiple processing cores on a single chip, making efficient programming of a single processor computer as complex as it once was to develop software for a high-performance computer.

The interface to extreme computing is also changing with increasing availability of infrastructure and services through Cloud Computing enabling access to both high-performance and high-throughput computing.

The computers and computer systems, however, are still largely designed with little thought to the applications that might use them and often arrive to the users with little software infrastructure in place. The development of the required software infrastructure is becoming increasingly complex and if we are to be effective in the use of these systems there is a need to address this through appropriate abstractions, programming languages and tools to support application developers and users.

A further pressure point for extreme computing is the constraint that energy requirements are having, and will continue to have, on the developments of computer architectures and their use. This becomes increasingly important both in terms of the general computing infrastructure which is becoming a significant component of national use of electricity [1] and for the high-end computer systems, with the predictions of the exascale energy requirements up to 100MW [2] – more than would be required by a small town. Given that a standard dual-core laptop is equivalent to one of the Top500 machines of only 12 years ago we can expect these issues to be pervasive as we go forward.

This presentation provides an overview of the challenges that are faced in extreme computing, together with some of the ongoing activities to provide a roadmap for future developments. Focussing on energy as a driver there is a consideration of some of the research and activities across the ecosystem of computing and in particular we consider that required for the Square Kilometre Array [3] which provides an excellent co-design vehicle¹.

Ultimately our interest is in providing a framework to allow the development of applications and algorithms that are energy-aware and will allow the user to make choices regarding the objective optimized in the software – be it energy performance, price performance or indeed computational performance.

¹ A *co-design vehicle* is a concept described in [ref] and is essentially an application that is developed in conjunction with the design of the software and hardware infrastructure together.

Extreme Computing Challenges and Constraints

There have been many studies during the last several years that have addressed the requirements of extreme (or exascale) computing [4, 5]. Here we will provide a general overview with some specifics from a recent study in the UK focussed on the development of a roadmap for high-performance computing [6] and the roadmap of the International Exascale Software Program [7] but the reader is referred to the original texts for the full discussion and details.

The HPC/NA roadmap activity in the UK brought together computational scientists, mathematicians and computer scientists to understand what might be the key issues for the future of extreme computing both in the UK and within an international context. The activity held a number of workshops and had individual communications with scientists and groups in the UK and a review of the literature from other national activities was included to provide the international perspective. It did not address hardware challenges directly, but the whole discussion was of course motivated by the changing face of computing hardware. The IESP roadmap provides a good overview of the hardware expectations [8]. As we go forward, we can expect that power and cooling limits will continue to constrain clock speeds; we will see a shift to multi/many core with attendant hierarchical parallelism an often with the additional complexity of hardware accelerators. This can be seen by the evidence of the Top 500 machines where multi-core machines and those with GPUs and the like are an increasing proportion [9]. Memory bandwidth will continue to be a bottleneck and this is likely to get worse. Chip design and chip to chip communication capabilities are likely to change however with 3D stacking and integration designs, for example [10] and the introduction of high-speed optical communications [11].

The HPC/NA roadmap activity found there were five major themes that dominated the discussions.

Theme 1: Cultural Issues

There are gaps between disciplines that meant that often the team effort that might be required to bring a large-scale application together was not optimal and that appropriate knowledge and expertise might not be captured within such an effort. This is a common issue and many national activities are investing to ensure that community actions are supported and interdisciplinary teams created.

Differences in discipline communities were also identified in that some application domain scientists are used to sharing models and codes, and reusing other software developed by other groups; while for other domains this approach is almost completely alien with codes being entirely developed within a particular group and little use being made of libraries or other third-party software.

Similarly there is a need to ensure that the software activities are not only focussed at a national level but are within an international context. It is the case that several application codes in the UK are dependent up on software that has been developed in programmes in other countries and vice versa. These dependencies are a risk to the sustainability of the software and that risk can only be mitigated by international collaboration.

Theme 2: Applications and Algorithm development

There are many issues that need to be addressed in the area of application and algorithm development. These include ensuring that a component approach to development is taken where appropriate such that algorithms that might underpin multiple application areas are developed as such. Having noted that, many applications involve multiple models at different scales or for different elements of the application and integration of components is not possible without appropriate standards for data models and formats, interoperability of programming models, and lack of knowledge of error propagation through the integrated system. Clearly this is not a problem unique to extreme computing, but is one that as a community we have yet to address fully.

The memory hierarchy of high-performance computers is on course to get even more complex which will drive the needs for hierarchical algorithms to deal with bandwidth across the memory hierarchy together with software strategies to mitigate high memory latencies. This will in turn drive the need for algorithms to be dynamically adaptive, perhaps as components of an active library, and of course scalable much beyond the present status quo. There are few applications that can scale to petaflops levels let alone toward exascale. Such scaling can only be achieved with appropriate portioning and load balancing and effective data management. Data-intensive science [12] is increasingly important and the ability to manage large amounts of data effectively and efficiently is crucial [13]. Moving data takes energy, at every level of the system – this is a point we will return to later.

Theme 3: Software Challenges

Development of software is becoming increasingly complex and without appropriate software engineering, in the context of an exascale system, will be un-maintainable. Application scientists are not usually software engineers but there is a need to ensure that by some means appropriate practice is adopted. Frameworks and tools to support software development will be needed together with compilers and code generation tools that can provide a layer of abstraction for application scientists.

Theme 4: Sustainability

The HPC/NA Roadmapping activity identified a general concern regarding the sustainability of application codes, software libraries and skills (we consider skills in the next section). This issue is integral to that of programming models, interoperability and also the cultural issues above. Scientists are naturally loathed to invest a great deal of effort in the development of software that will last only the lifetime of a particular computer architecture, compiler or other dependent component of the environment. We can only address this issue as a community as we agree in the adoption of standards and open source mechanisms.

Theme 5: Knowledge Base

It is important to ensure that computational scientists have the right set of skills, and this set of skills needs constantly updating. Within the UK activity it was also discovered that there was a lack of awareness by some members of the community of existing libraries/packages. Maintaining a strong knowledge base will require education for graduate students as well as the opportunity for life-long learning. The report on exascale computing for energy and environment [14] notes "The current belief is that the broad market is not likely to be able to adopt multi-core systems at the 1000-processor level without a substantial revolution in software and programming techniques for the

hundreds of thousands of programmers who work in industry and do not yet have adequate parallel programming skills."

Other issues:

Since the last of the roadmapping workshops Cloud Computing [15] has matured and is now very much part of the computational environment. It offers a different interface to distributed, and in some cases, high-performance computing and as such is part of the evolving ecosystem for application development. As data-intensive applications dominate it is possible that the business models underpinning Cloud computing might offer benefits to scientists. These issues will likely be resolved as computational scientists become engaged with the various vendors but there is no question that the scale of Cloud computing infrastructure may offer benefits in terms of energy efficiency and from that point of view we will return to them as part of the ecosystem later.

Co-Design Vehicles and the Square Kilometre Array

During the workshops and activities of the Department of Energy in the USA addressing the requirements for exascale computing, the notion of a co-design vehicle has been adopted. At present, other than in a few specific cases [15, 16, 17], a computer system is designed with little thought for the applications that will use it, and the software is generally designed to try to fit the hardware. In the co-design model the two are seen as part of the same system and the design of each is done through collaboration. In the extreme this might result in computer systems that are only useful for a given application such as [15, 16, 17] but the model can be adopted to enable more generic progression with the evolution of the computer architect and algorithms being better aligned and providing better capability and specifically energy performance.

We are interested in the computing challenge provided by the development of the Square Kilometre Array [?], the next generation radio telescope. The telescope is in the design phase with the anticipation that construction of the first phase will begin in 2016 with the full telescope completed and in operation by 2022. The SKA will likely be located in either Australia or South Africa, in a desert so as to have little or no interference, but is a collaborative effort of over 50 groups in 19 countries.

The present design [18,19,20] has a combination of aperture arrays in the core and up to 3000 phased array feeds on dishes and a collecting area of approximately one square kilometre with receptors extending out to a distance of 3000km from the centre of the telescope (figure 1). It will allow a sensitivity of more than 50 times that of existing telescopes, and 10,000 times the survey speed and will provide data to answer fundamental science questions on gravitation and magnetism, galaxy formation and even the question of life on other planets. The design of the SKA is developing through design studies based on the science requirements, Pathfinder telescopes that provide experience of design options, and technology capability considerations.



Figure 1 Possible configuration of SKA receptors and artist's impression of SKA core(from [21])

The SKA provides a fabulous information technology challenge with a typical data rate from each dish antenna on the order of 100Gbs⁻¹ aggregating to over 100Tbs⁻¹ [21] and need for Exaflop computation [23] for post-processing. The infrastructure required to support the various science cases will need to range from real-time capability to transport and analyse the data at these high-data rates and the capacity to store and "publish" the data for later analysis and interpretation by the global astrophysics community. The computational systems will likely range from specifically designed FPGA-like units to exascale computing and Cloud-like data centres. The communications infrastructure will range from intra-chip and inter-chip with optical fibre to the correlator and on to a high-performance computer, to trans-oceanographic with the latter having data rates of at least 100Gbs⁻¹ over the general network providers. The SKA will succeed or not depending on both the physical implementation of the telescope design and the software infrastructure that will enable it. The software infrastructure required to realise this information technology challenge is itself has been identified as > 2000 person year task [22] but even this may not take full account of the complexity of the task.

The fundamental algorithms in the SKA data pipeline include beamforming, gridding, convolution and filtering algorithms. Cornwell et al [23] estimate the computational requirements in the context of the Top500 and show them to be beyond the scale of projected performance in the timescale required (figure 2).

Alongside the challenges of the computational infrastructure are the related challenges of powering the infrastructure. This includes the power



Figure 2: SKA computational requirements in context of Top500 [23]

required for the core antennas (~30MW) and the remote stations ~0.5MW) and the various computational components including high-performance computing (~40MW) data transmission (~?MW) and Cloud (or equivalent) provision (~?MW). Energy provision is a major design factor in the delivery of the telescope – plans to mitigate the energy constraints include renewable energy sources (sun) providing the station power and it is easy to see that the SKA could possibly be the largest Green IT project ever to be considered.

Counting Joules

As indicated above the ecosystem of computational resources required to enable the SKA provides a plethora of challenges and opportunities where energy efficiency is concerned.

Moving data takes energy whether the data is moving from L2 cache on chip or within a transatlantic Data Cloud. Effective management of that data communication is a major component of any There has been a great deal of research on wireless network optimal energy model. communications and sensor networks, where devices are most often low-energy devices with battery constraints. Indeed a lot of research has been done in general on low-energy devices including computational algorithms from which we might learn. The existing communications network across the globe relies on hundreds of thousands of switches and routers and these, unlike wireless or mobile infrastructure, do not power down when idle. The network is a complex system of different technologies and it is difficult to predict the energy consumption under different circumstances so sending data 10 times as fast might use interfaces that use 100 times the amount energy or in some cases (e.g. newer optical ones) 1000 times less [24]. A good overview of network energy costs for realistic configurations can be found in [25]. There are a number of research efforts considering the issues around Green network communications including INTERNET [26] and others can be seen at recent conferences [27,28]. Of course as energy-aware systems are developed there are still issues of what they optimise - usage or cost? Qureshi et al [29] illustrate effective ways in which energy-aware data centres can optimise cost by moving computation to nearby states where electricity costs are less.

Optimizing energy usage in large-scale data centres and Clouds is almost a science in itself. McKinsey and the UpTime Insitute [1] indicate that the energy used by data centres in the US is becoming a significant percentage and is likely to overtake airlines in terms of carbon emissions. The report states that the average data centre uses as much power as 25,000 households, but that estimate is probably somewhat out of date as in the last few years Cloud provisioners have built very –large scale data centres across the US [30]. On a somewhat smaller scale at the Oxford Supercomputing Centre (OSC) we have developed software that intelligently powers down components of the systems at times of under utilisation. The indications are that this will provide significant savings.

The McKinsey report identifies a number of issues around the effectiveness of data centres including siloed organisations and limited transparency that match very well with the findings of the HPC/NA for that community. McKinsey make the recommendation that metrics be defined that are not only measuring the facility but are linked to the applications using it and the processes integral to it. A consortium called the Green Grid is now in place with the aim to develop standards and best

practice for data centres. The recommendation regarding metrics is one that we, as a community, should also take on board as we develop exascale technology.

At the computing system level there are many approaches to energy efficiency. Again there is a large knowledge base in the low-energy systems community that we should consult. Indeed an approach that might be considered for some application areas builds directly on the low-energy embedded microprocessor technology to build petascale systems. This approach is under development for the Green Flash system [31, 32] at Berkeley National Laboratory where researchers are collaborating directly with Tensilica, Inc. to explore the use of Tensilica's Xtensa processor cores in building a computer to model clouds (that is clouds in the weather system not as a computer system!). They believe that a specially designed core could get 10–100 times better performance per watt. The Green Flash design will have on the order of 20 million processors consuming less than 4 MW whereas the equivalent in conventional microprocessors would require around 200MW.

With appropriate systems support it is possible to provide hooks into the operating system to allow spindown policies to be enacted, adaptive placement of memory blocks, agile use of component devices and of course energy-aware routing [33]. To make use of this level of support it would be helpful to have tools that can provide a mapping from the higher-level action to the underlying activity.

The power cost for a device is proportional to the frequency cubed. This means that a multi-core device that is running at a lower clock rate is bound to provide the potential for higher performance at a lower energy cost. Interesting issues arise when there are multiple such devices to hand and the algorithm designer can chose to optimise the energy usage for a given computational performance. At this point in time the SKA community are considering the use of GPUs, FPGAs and multi-core chips [34, 35]. Nieuwpoort et al [34] show some interesting results in this area where they have used a variety of chipsets to underpin the computation to correlate signals. The results are particularly interesting in respect to the achieved efficiency vs Gflops/Watt where although a given device might underperform in terms of percentage efficiency it may still perform better in Gflops/Watt. The Green500 provide good reference indicators for feasible platforms at this time [36].

Within the Oskar project [37] we tackled the Digital beamforming for the aperture array components of the SKA that pose considerable computational challenges [38]. The proposed algorithm provides a hierarchical algorithm for beamforming using a simplified and flexible computational approach of direct matrix-vector multiplies rather than FFTs that provides a reduced data rate and a computational cost of forming beams of 1TFlop as opposed to 20TFlops for the FFT approach.

Within PrepSKA [39] Savlini et al [40] have developed a Pipeline for Extensible Lightweight Imaging and Calibration (PELICAN). This framework for parallel quasi-real time data processing offers two deployment options either with the server supplying multiple pipelines or the pipeline connecting to the data stream directly. PELICAN allows reuse of modular components and will be deployed on LOFAR [41] interferometer stations to allow all sky calibration and imaging, and pre-processing for pulsar searching. The modular framework allows appropriate computational components to be implemented on appropriate devices, in this case GPGPUs.

Tools, Benchmarks, metrics

While benchmarks for datacentres and systems are now relatively mature [42, 43] there is still work to be done at the numerical algorithms level. On the exascale roadmap [8] it is suggested that standards to support energy-aware algorithms should be agreed in 2014/2015 with energy-aware libraries available in 2016. We have our work cut out!

Although for the SKA specific algorithms of interest we can model and implement across different platforms there is a lack of agreed metrics, tools and benchmarks to provide the community with clear evidence to support design features. This is true across the board for computational algorithms although there are some efforts in this direction.

There are a number of tools for a variety of platforms that will allow measurements of device and sub-device power usage. Kirk Cameron and his group have created a profiling tool that is an open source software system together with hardware power measurement devices, PowerPack [44]. Using the profiler they have analysed the HPC Challenge Benchmarks [45] and been able to provide very clear analysis of the behaviour of the algorithms in terms of energy used.

It may be time to consider again what benchmarks should stress. As we consider adaptive algorithms that are able to autotune to heterogeneous sets of processors we may need to consider a different set of characteristics and metrics that are captured by existing suites. There may be instances where replacing data movement by computation provides a saving in energy while maintaining computational speed.

Conclusions

Energy constraints provide one of the major challenges for extreme computing in the future. There are opportunities at every level of the computational infrastructure to address the challenges through both efficient physical hardware but also through more sophisticated use of the physical infrastructure. An intelligent software infrastructure could reduce the energy use on chip, within a processor, between processors, between computers, across Cloud platforms and indeed over international boundaries.

Ultimately our own interest lies in the development of energy-aware algorithms for extreme computing. We believe that to achieve such algorithms and to see the benefits there will need to be significant investment in the support of tools and low-level mechanism to allow the profiling of energy use across platforms. Equally important are standard metrics and computational benchmarks that will allow agreement on measures of "success". Just as the linpack benchmark has come to be seen as the measure of the performance of a computer system we need the equivalent application/algorithm benchmarks that capture the energy characteristics of any given system.

The future holds continuing complexity in computing systems from the combinations of computing within the whole ecosystem to the heterogeneity of chips on our laptops. The most effective algorithms at any given point in time for any specific application are going to depend upon the choices that can be made given the hardware configurations. In general, at present we cannot even articulate those choices let alone provision the most appropriate algorithm. This challenge of

articulation and optimal algorithm design can only be addressed through a co-design approach with hardware, software and application scientists working together.

While we have focussed here on the energy aspects of the problem, and specifically as driven by SKA, there are any one of several other key issues that could and should be addressed including the complexity of the software development for such an infrastructure, the usability of the software systems, the data management and related semantic issues for such colossal data systems or the provision of imaging and data analysis in the Cloud. Each of these is a talk of its own – for next time.

Acknowledgements

My thanks to colleagues on the prepSKA project Steve Rawlings, Aris Karastergiou, Stef Salvini, Ben Mort, Chris Williams, Fred Dulwich and Andy Faulkner. The HPC/NA project was in collaboration with Nick Higham, Iain Duff, Peter Coveney, Mark Hylton and Stef Salvini. I am grateful to Jeyan Thiyagalingam, Simon McIntosh-Smith, Jon Crowcroft and Jaafar Elmirghani for their input and suggestions.

References

- 1. Revolutionizing Data Center Efficiency, McKinsey Company and the Uptime Institute, www.**mckinsey**.com/.../Revolutionizing_Data_Center_Efficiency.pdf
- 2. <u>www.exascale.org</u>
- 3. The SKA project: <u>www.skatelescope.org/</u>
- 4. The DARPA Exascale Study Report: <u>http://users.ece.gatech.edu/mrichard/ExascaleComputingStudyReports/exascale_final_repo</u> <u>rt_100208.pdf</u>
- 5. The Department of Energy workshop on extreme computing: <u>http://extremecomputing.labworks.org/index.stm</u>
- 6. The HPC/NA Roadmap: www.oerc.ox.ac.uk/research/hpc-na/roadmap
- 7. The International Exascale Software Project: <u>www.exascale.org</u>
- 8. The IESP roadmap: <u>http://www.exascale.org/mediawiki/images/4/42/IESP-roadmap-1.0.pdf</u>
- 9. The Top500 machines: <u>www.top500.org</u>
- 10. Intel <u>http://www.eetimes.com/electronics-news/4204959/Intel-Silicon-Optics</u>
- 3D DRAM Design and Application to 3D Multicore Systems, Hongbin Sun, Jibang Liu, Rakesh S. Anigundi, Nanning Zheng, Jian-Qiang Lu, Kenneth Rose, Tong Zhang, IEEE Design and Test of Computers, vol. 26, no. 5, pp. 36-47, September/October, 2009.
- 12. The Fourth Paradigm: Data Intensive Research, http://research.microsoft.com/enus/collaboration/fourthparadigm/
- The Data Deluge: An e-Science Perspective, A.J.G. Hey and A.E. Trefethen, in *Grid Computing* - *Making the Global Infrastructure a Reality*, Berman, Fox, Hey, eds., pp 809-824, Wiley and Sons, 2003.
- 14. DOE Report on <u>Modeling and Simulation at the Exascale for Energy and the Environment</u>, June 2007
- 15. Above the Clouds: A Berkeley View of Cloud Computing, Michael Armbrust, Armando Fox, Rean Griffith, Anthony D. Joseph, Randy Katz, Andy Konwinski, Gunho Lee, David Patterson,

Ariel Rabkin, Ion Stoica, and Matei Zaharia, UC Berkeley Reliable Adaptive Distributed Systems Laboratory . http://radlab.cs.berkeley.edu/ February 10, 2009

- 16. The QCDOC project: *ukqcd.epcc.ed.ac.uk/community/qcdoc/*
- 17. The MD-Grape3 machine: <u>http://www.peta.co.jp/index-en.html</u>
- 18. The D.E.Shaw Machine: www.deshawresearch.com/
- 19. The Square Kilometre Array, <u>http://www.skatelescope.org/</u>
- 20. The Square Kilometre Array, Peter E. Dewdney, Peter J. Hall, Richard T. Schilizzi, and T. Joseph L. W. Lazio, Proceedings of the IEEE | Vol. 97,No. 8, August 2009
- The Square Kilometer Array (SKA) Radio Telescope: Progress and Technical Directions, P.J. Hall, Richard T. Schilizzi, and T. Joseph L. W. Lazio, U.R.S.I, The Radio Telescope Bulletin, No 326, September 2008
- SKA Memo 100: Preliminary Specifications for the Square Kilometre Array, R. T. Schilizzi, P. Alexander, J. M. Cordes, P. E.Dewdney, R. D. Ekers, A. J. Faulkner, B. M. Gaensler, P. J. Hall, J. L. Jonas, K. I. Kellermann, http://www.skatelescope.org/PDF/memos/100 Memo Schilizzi.pdf

nttp://www.skatelescope.org/PDF/memos/100_Wemo_Schilizzi.pdf

- 23. Scaling Mount Exaflop: from the pathfinders to the Square Kilometre Array, T.J. Cornwell, *Ger van Diepen, <u>www.atnf.csiro.au/people/tim.cornwell/publications/MountExaflop.pdf</u>*
- 24. Green Optical Communications—Part II: Energy Limitations in Networks, Rodney S. Tucker, IEEE JOURNAL OF SELECTED TOPICS IN QUANTUM ELECTRONICS
- 25. Private communication with Jon Crowcroft
- 26. First ACM SIGCOMM Workshop on Green Networking, August 2010, <u>http://conferences.sigcomm.org/sigcomm/2010/gncfp.php</u>
- 27. INTERNET: Intelligent Energy Aware Networks, <u>http://www.internet-project.org.uk/</u>
- 28. E-energy: First international conference on energy-efficient computing and networking, April 2010, <u>http://www.e-energy.uni-passau.de/</u>
- 29. Cutting the Electric Bill for Internet-Scale Systems, A. Qureshi, R. Weber, H. Balakrishnan, J. Guttag and B. Maggs. SIGCOMM 2009.
- 30. Data Centres in the News: <u>http://www.datacenterknowledge.com/archives/2010/02/05/virginia-nc-battling-for-</u> <u>microsoft-data-center/</u>
- *31.* The Green Flash computer project, <u>http://www.lbl.gov/cs/html/greenflash.html</u>
- *32.* Towards Ultra-High Resolution Models of Climate and Weather, Michael Wehner, Leonid Oliker, and John Shalf, *International Journal of High Performance Computing Applications May 2008 22: 149-165, doi:10.1177/1094342007085023*
- 33. Every Joule is Precious: The Case for Revisiting Operating System Design for Energy Efficiency, Amin Vahdat, Alvin Lebeck, and Carla Schlatter Ellis, Proceedings of the 9th workshop on ACM SIGOPS European workshop: beyond the PC: new challenges for the operating system Pages: 31 – 36, 2000 ISBN:1-23456-789-0
- Using many-core hardware to correlate radio astronomy signals, Rob V. van Nieuwpoort, John W. Romein, Proceedings of the 23rd international conference on Supercomputing, pages: 440-449, 2009, ISBN:978-1-60558-498-01
- 35. Evaluating Multi-Core Platforms for HPC Data-Intensive Kernels, Alexander S. van Amesfoort, Ana L. Varbanescu, Henk J. Sips, Delft University of Technology, The Netherlands
- 36. The Top Green 500 Computers, <u>www.green500.org</u>
- 37. The OSKAR project: http://www.oerc.ox.ac.uk/research/oskar

- OSKAR: Simulating Digital Beamforming for the SKA Aperture Array, Fred Dulwich, Benjamin J. Mort, Stefano Salvini, Kristian Zarb Adami, and Mike E. Jones, Widefield Science and Technology for the SKA, SKADS Conference 2009
- 39. Preparation for SKA: www-astro.physics.ox.ac.uk/~sr/prepska.html
- 40. The Pelican framework: <u>https://wiki.oerc.ox.ac.uk/svn/pelican/slides/pelican-slides2.ppt</u>
- 41. The Lofar Pathfinder: <u>www.lofar.org/</u>
- 42. The SPEC power benchmark, <u>http://www.spec.org/benchmarks.html#power</u>
- 43. JouleSort: A Balanced Energy-Efficiency Benchmark, Suzanne Rivoire, Mehul A. Shah, Parthasarathy, Ranganathan, Christos Kozyrakis, SIGMOD'07, June 11–14, 2007, Beijing, China
- 44. PowerPack: Energy Profiling and Analysis of High-Performance Systems and Applications, Rong Ge, Xizhou Feng, Shuaiwen Song, Hung-Ching Chang, Dong Li, Kirk W. Cameron, IEEE Transactions on Parallel and Distributed Systems, 23 Apr. 2009. IEEE computer Society Digital Library. IEEE Computer Society,

<http://doi.ieeecomputersociety.org/10.1109/TPDS.2009.76>

45. Energy Profiling and Analysis of the HPC Challenge Benchmarks Source, Shuaiwen Song, Rong Ge, Xizhou Feng, Kirk W. Cameron International Journal of High Performance Computing Applications Volume 23, Issue 3 (August 2009), Pages: 265-276, 2009



10PetaFLOPS は 生命科学で何を実現できるか

理化学研究所

姫野 龍太郎

10PetaFLOPSは生命科学で何を実現できるか

理化学研究所 次世代生命体統合シミュレーション研究グループ グループディレクタ 姫野龍太郎

Keywords

スーパーコンピュータ、生命科学、シミュレーション、分子、細胞、臓器、全身、脳神経系、マルチ スケール

1. はじめに

次世代スーパーコンピュータは来年度から部分的に稼働が始まり、再来年度には全システムの稼働が予定されている。この稼働に合わせて、生命科学分野でのグランドチャレンジとも呼ばれる次世代生命体統合シミュレーション研究開発プロジェクト(ISLiM)が、2006年秋から理研を中核拠点として行われている。このプロジェクトはペタスケールのシミュレーション技術によって、ライフサイエンスに仮説・検証型の新たな研究手段を提供し、生命現象を定量的かつ統合的に理解・予測・解明することを目指すと共に、創薬・ヘルスサイエンスサイエンスへの貢献、新規医療技術の実用化を図るものである。



2. 推進体制とソフトウェア開発の状況

このような目的の実現のため、分子・細胞・臓器全身という3つの階層に分かれた研究チームと、 大量の実験データから法則に迫るデータ解析融合チーム、更に脳神経系チーム、全てのチームを支 える HPC チームの6チームを設置して取り組んでいる。

これらのチームで現在開発中のソフトウェアは量子化学や分子動力学から構造流体連成計算、脳の局所回路モデル、可視化ソフトウェア、並列化のためのミドルウェアなど、合計34本にのぼる。 これまでの開発で、並列化に関して1,024を超えるところまでテストが終わったものが18本、8,000 並列までのテストが終わったものは9本という状況である。今後は次世代スーパーコンピュータで の実効性能と実行時間の予測を行い、どの程度のモデルでどのくらいの計算をするかを検討してゆ く予定である。

3. 現在の見通し

昨年から開発中のソフトウェアの並列性能などの開発状況を睨みながら、10PetaFLOPSの次世 代スーパーコンピュータを使ってどんな問題を解くかを検討してきた。また、どの問題をどの時期 に解いてゆくか、問題間の優先順位をどのようにつけるかを検討した。その結果、それぞれのソフ トウェアを3つに分類、計算を行う時期に合わせて、第一走者、第二走者、第三走者として、波状 的に成果を出してゆくこととした。開発中のソフトウェアの数が多いので、第一走者と第二走者の 一部について、そのソフトウェアの概要と解こうとしている問題を紹介する。





		次世代	この研究 えパコン	開発の位 開発プロ・	置づけ ジェクトの	一部 現在	S		RIKE
		平成18年度 (2006)	平成19年度 (2007)	平成20年度 (2008)	平成21年度 (2009)	平成: (20	2 年度 10)	平成23年度 (2011)	平成24年度 (2012)
5 2	7 4	概念	1928+ / 8¥1	細設計	/	試作·	平価・製	造	性能 チューニ ング
ソフトワン	次世代ナノ統合 シミュレーション		1	開発・製作・評価				実証	供用開始
F 10	次世代生命体統合 シミュレーション			開発・製作・	評価			実	äÆ
	計算機棟		設計	建	設				
愛	研究棟			設計	建設				























≫ 「生命	向体基盤	ソフトウ	ェア開発	·高度化	; 7 -4
<mark>ISL</mark> 分子	細胞	臓器全身	<i>デー</i> タ 解析融合	その他	各機関(备=N チーム)の
#	通基盤ライ	高性能化支			基盤ソフト開 発・高度化支
GUI/ワ-	-クフロー等	直感的ないユーザ	操作等を実 一の利用を	そう そ現し幅広 を可能に	援 実証利用 支援
_{実現象} 次世代	スパコン用	シミュレー	ションソフト	ウェア群	グランドチャレ ンジ達成目標
して フィード 自ら改良 バック		商用	ノフトウェア	への改良 ~	ソフトウェア ハウス
	学行	析利用・ 産∮	業利用		















	2	3		現	在	の	開発	状況	ļ			2
T	48	アプリケーション名			略种		開発賞	任者				RIKEN
	A1-1	マルチコピー・マルチスケール分子ショ 基盤となるクラスライブラリ	ュレーション	法開発の	Platys	us-MM/CG	木寺					
	A1-2	レプリカ交換分子動力学計算インター	1=17		Distin	RCM	2010	X48*				2010868148
	A2-1	全原子分子動力学計算	現Phase	到建設列集	2	進み(月)	ハイブリッド 並列化	FJコンパイラ 対応(RECC)	FJコンパイラ 対応(FX1)	実行時間の 測定	清算量の測 定	时建蓝列数
쓝		Real Ballin Frider Margar	I-1	1	4000	-7	×	0	0	*	*	8000
Ŧ		CANERAL B X7777EL	2-1		1024			0		÷.		1024
	A3-1	🎒 管理:開発フェー	ズを4	っに	计规	i						2048
	A8-1	Phase I:プロ	トタイス	プ開発	と覚	列化0	ロテスト	(想定2)	006年~	-2008年	E)	8000
	A6-2	Phase II 大規	模並	列化(想定	2009	年~20	10年)			· /	0000
	A5-3	Phase III:招·	た相構	前词-	⊥重			.。」) バ(相定	2011年)		8000
-	81-1	Dhase IV. 7	一代は			での木	_ 故的た	ア(心た	2011- DD	/ 相定20	12年)	128
			ጠ ተ	- +2	1-1	2007	コロロジみ		に自及い	10.0000	12 4)	8000
	01-1		041	05	51-1	5 507.	1100	~~~ <i>~</i> ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	ここでに、	I OFX PE		2000
-	C1-2	(1)		_								128
	C1-3	重粒子線治療シミュレーション	8-1	1	1024	-5	×	×		×	ж	1024
R	02-1	気快装酒産シミュレーション (ボクセル超音波伝播プログラム)		0	000	-	~	0	0	0	0	8000
	C3-1	肺呼吸・肺循環シミュレーション	R	<u>発</u> い/-	7.5	דדל	全34才	-				8000
-	06-1	マルチスケール・マルチフィジックス心	1-0	175 / -	< 1	/ .	±044	`				512
2		リケーマップニー	8-1	18	本:	1024	並列以	上、107	5:800	0並列。	ᆝ上	1024
1	ア月 、問	調査があれば機動	8-2	可	視化	シソフト	ŧ1000	GPGPL	I並列			8000 24
的	りに	対処	8-2	8	000	4	×	0	0	×	ж	8000























	アブリケーション名	絡称	開発責任者	含語	並列化の方法	
	密度汎開数法に基づくタンパク賞全電子波動関数計算	ProteinDF	佐藤(東大)	C++	MPL OpenMP	
	粗視化モデル計算	CafeMol	高田(京大)	Fortran90	MPI, OpenMP	
	全原子分子動力学計算	MARBLE	池口(横浜市大)	с	MPI	
	マルチコピー・マルチスケール分子シミュレーション法開発の	Platypus-MM/CG	木寺(横浜市大・理研)	C++	MPI	
分子	ハイブリッドQM/MM反応自由エネルギー計算	Platypus-QM/MM-FE	林(京大)	FORTRAN77	GAMESSのsocksラ イブラリ	
	レプリカ交換分子動力学計算インターフェイス	Platypus-REIN	杉田(理研)	Fortma90	MPI	
	組視化モデル計算/分子動力学計算	Platypus-CGM/MM	中村(版大)	C++,	MPI	
	量子化学計算/分子動力学計算	Platypus-QM/MM	中村(版大)	FORTRAN77,	MPI	
	量子化学計算	Platypus-GM	中村(版大)	FORTRAN77,	MPI	
細葉	細胞シミュレーションプラットフォーム	RICS	横田(理研)	Fortran, C,	Sphere(MPI)	
	全身ボクセルシミュレーション(ボクセル構造液体達成解析プログラム)	SPH3D	高木(東大・理研)	Fortran. C.	Sphere(MPI)	
	マルチスケール・マルチフィジックス心臓シミュレーション	UTHeart	久田(東大)	Fortran	MPS&OpenMP	
読器	低侵装治療シミュレーション (ボクヤル総合演伝導プログラム)	HIFU	高木・松本(東大)	Fortran90	MPI&OpenMP	
全日	数小循環シミュレータ (使め込み後界はにとみ取小原理プログラム)	ZZ-RBC	高木(東大・理研)	Fortran90	Sphere(MPI)	
R	重粒子線治療シミュレーション	ZZ-DOSE	高木(東大·理研)	Fortran90	MPI	
	林祥盛・林嶺擅シミュレーション	77-11000	\$0 (D) (N) (+)	Fortran, C.	a	

アプリのリスト(2)

ハブロタイブ関連解析に於ける統計検定を行うためのソフト ウェア	ParaHaplo	鎌谷(理研)	c	MPL OpenMP
大規模遺伝子制御ネットワーク推定プログラム	SIGN	宮野(東大)	C99	MPL EP (SGE)
再帰的正則化法による生体内分子の大規模ネットワーク推定 ゴログラム	LIGN	宮野(東大)	R	EP (SGE)
状態空間モデルによる時系列データからの遺伝子ネットワー 2推定プログラム	SSM	宮野(東大)	C, C++. JAVA	MPI, OpenMP
データ解析融合プラットフォームの開発	SBIP	宮野(東大)		ブラットフォーム
県羅的タンパク質ドッキング解析プログラム	MEGADOCK	秋山(東工大)	C++	MPI
生命体データ同化プログラム	LISDAS	4月四(統数研)	Fortran90, C, C++	MPL OpenMP
文世代シークエンス解析プログラム	NG5 analyzer	鎌谷(理研)	pert, shell	EP
拡張RAT法による25NP組合せの全ゲノム関連解析ソフトウェ ア	ExRAT	鎌谷(理研)	C++	MPI
Neural Simulation Tool	NEST	Diesman(理研)		
Cortical Microcircuit Developed on NEST - Interneuron	CMDN	深井(理研)	SLL C++	MPI
神経細胞形態シミュレーションキット	NeuroMorphoKit	石井(京大)	MATLAB, C/C++	逐次
模覚系シミュレーションのための共有ブラットフォーム	VSM	臼井(理研)		ブラットフォーム
昆虫嗅覚系全脳シミュレータ	IOSSIM	神崎(東大)		実験系アプリ
大規模並列用MDコアプログラム	cppmd	泰地(理研)	C++	MPI
大規模仮想化合物ライブラリー	VLSVL	船津(東大)		EP
分散並列大規模データ可視化システム	LSV	小野(理研)	C++, C	MPI
アプリケーションミドルウェア	SPHERE	(195(1948)	C++ C	MPI

SLIM











世界の状況



- ・ 並列化の状況に関しては、アプリケーションによってまちまち。
- MDの場合:

Desmond (D. E. Shaw Research), NAMD(UIUC), Blue Matter (IBM) などが数千並列以上の並列を実 現。

何れのコードも、高性能計算の専門家が中心に参加し て開発を進めている。

ライフアプリケーションでの10⁴以上の並列度については、世界的に見てもこれまでの蓄積は非常に少ない。















51 Eppler, Helias, Muller, Diesmann, and Gewaltig (2008) 51 Frontiers in Neuroinformatics 2:12. doi:10.3389/neuro.11.012.2008













8	まとめ	
 開発ン - その - 10本 - 可視 今後、 次世代 計算時 し、利 	フトウェア:全34本 うち18本:1024並列以上 :8000並列以上 化ソフトも100GPGPU並列 並列性能の改善 スパコンにフィットした最近 間に合わせた適正な計算 目計画を立案、フル稼働に	^{適化} 「規模を設定 -備える
		58



科	学	技	術	計	算	分	科	会		選	出
					2010 쇼	∓度 科	学技術	i計算分	·科ɗ	<u> </u>	より

原子力研究開発を推進する スーパーコンピュータシステムと 原子力アプリケーション

日本原子力研究開発機構

平山 俊雄

原子力研究開発を推進するスーパーコンピュータシステムと 原子力アプリケーション

平山 俊雄 ・ 清水 大志 ・ 久米 悦雄 日本原子力研究開発機構・システム計算科学センター

[アブストラクト]

原子力機構では、計算科学を活用した原子力の研究開発を加速するために、本年、3月には、Altix3700 Bx2 (13TFlops) に替えて、国内最高性能となるピーク性能 200 TFlops の Linux クラスターシステム (BX900) とピーク性能 12 TFlops の次世代計算機プロトタイプ機 (FX1) からなる2種類のスーパーコンピュータを導入し、それらの運用を開始した。 BX900 は、旺盛な計算需要に応えるために、また FX1 は次世代計算機 (京コンピュータ) の利用に向けた原子力アプリケーションのチューニングを目的としている。

BX900 については、その利用開始後3日目で90%を超える利用率を達成し、その後も高い利用状況が続いている。これは、原子力計算科学研究者が数千コアを用いた大規模並列計算コードの開発能力と利用技術力を既に兼ね備えていることを表している。次世代計算機向け原子力コードのチューニングについては、FX1を用いたプロファイリングデータに基づく性能予測を行い、1PFlopsまでの良好な並列性能を推定した。また、新旧3機種のスーパーコンピュータの基本性能を比較した結果についても言及する。

[キーワード]

原子力計算科学、次世代計算機、ペタフロップス計算、BX900、FX1、ハイブリッド並列

1. はじめに

計算科学技術は「理論」及び「実験」と並ぶ第3の研究手法として、21世紀の先端的研究のフロン ティアを切り開くための重要な基盤技術となっている。特に、原子力のような巨大技術においては、安 全面や時間・空間の制約により実験が困難な場合が多く、計算科学技術は従来から重要な研究手段とな っている。すなわち、計算科学技術は、原子力研究開発の効率化、原子力施設の安全評価、国際競争力 強化のカギを握る技術となっている。原子力委員会においても、その原子力政策大綱の中で「シミュレ ーション技術の高度化による大規模な技術システム開発の効率化」を求めている。

原子力関連二法人の統合と独立行政法人化を経て、旧二法人が所有していたスーパーコンピュータの 資源について、整理統合と合理化を果たし、膨大な計算需要に応えるための新たな設備を導入すると共 に、京コンピュータの利用に向けた準備を整えた。

2. 新スーパーコンピュータの性能比較

BX900 の Linpack 性能は 191.40 TFlops、実行効率 95.66 %であり、2010 年 6 月時点の TOP500 によれば世界 22 位、実行効率世界 2 位であった。FX1 については、11.66 TFlops、実行効率 90.37 % である。

第1図に、新旧3機種のHPC Challenge Benchmark の性能を示す。ここで、 Altix3700Bx2(Itanium2)は更新前の機種で、ピ ーク性能13 TFlops である。Altix3700Bx2 に比 較して、BX900は演算性能、通信性能ともにコ アあたりのピーク性能差(~2倍)だけの改善が みられる。ただし、プロセス間通信のバンド幅 および遅延については、ソケットあたりのコア 数の増加やインターコネクトの構造の違いのた め、単純に比較することができない。

FX1 については、実効メモリバンド幅及びプ





ロセス間通信性能が BX900 と比較して低位にある。特に、実効メモリバンド幅については、BX900 が カタログ性能の 70%程度出ているのに対し、FX1 では 30%以下となっている。この問題は、メモリコ ントローラチップが CPU の外にあることに起因しているようだが、京コンピュータでは CPU 内蔵 と なっており、本問題は解決しているとのことである。

3. 京コンピュータの利用に向けたチューニングと性能予測

将来において膨大な計算需要を擁している核融合コードについて、 PFlops 級の計算を想定した並列性 能および実行性能を BX900 および FX1 を用いて評価した。

BX900 の 16,384 コアを用いて並列性能を評価した結果、全てのコアに対して MPI プロセスを割り 当てる FlatMPI では MPI プロセス間通信量の増大により、全体処理性能は大幅に低下する。したがっ て、大規模並列計算ではスレッド並列を併用する Hybrid 並列化を用いることで、MPI プロセス数を削 減することが必須となる。なお、富士通 Fortran の自動並列化は核融合コードについてはかなりの部分 で有効であった。また、メモリ使用量を軽減する上でも、Hybrid 並列化は有効である。

第2図は、京コンピュータ4096並列(~500TFlops) 上での性能予測結果である。FX1の256CPU(1024 コア)を用いて性能予測をするため、4096並列の1プ ロセスあたりの問題規模を保ったまま、weak scaling を仮定して4096プロセス並列から256プロセスにス ケールダウンをする。次に、FX1の256CPUを用いて 取得したプロファイリング情報に対して、FX1から京 コンピュータへの性能変換表を用いて演算時間を推定 した。通信時間については、当該コードの通信の種類、 データ長、通信回数を調査し、インターコネクトのト ポロジー(仮定)をもとに机上計算により推定してい る。



演算性能については、京コンピュータにおいてプロセスあたり 2.4 倍程度の高速化が期待できる。京 コンピュータのコア数が、4 コア (FX1)から8 コア (京)と2倍になっていることから、演算性能は 十分にスケールしているといえる。しかし、1 コアあたりの性能はすでに飽和しつつあることに留意す る必要がある。したがって、京コンピュータの利用にあたっては、コア数が増えた時のコードの並列化 性能が特に重要になる。

第3図は、京コンピュータの1PFlops について、 上記推定手法に基づいた2段階の Strong Scaling で ある。1PFlops においても良好な並列化性能が期待で きることがわかる。

4. まとめ

日本原子力研究開発機構では、「シミュレーション 技術の高度化による大規模な技術システム開発の効 率化」を求めた原子力政策大綱に基づき、計算科学的 手法を用いた原子力研究開発を推進してきた。その結



図3:1PFlopsまでのスケーラビィリティ

果、年間の論文数の 20%以上が計算科学的手法を用いた研究に移行している。しかし、数 PFlops に達 する膨大な潜在需要に応えるためには、計算機資源の不足は深刻な状況であり、京コンピュータの積極 的利用が不可欠となっている。

京コンピュータのプロトタイプ・アーキテクチャの FX1 を用いたコード分析とチューニングは、京 コンピュータの利用に向けた準備作業として極めて有効であることを確認した。また、原子力コードに ついて FX1 を用いたチューニングを行い、1PFlops までの良好な並列性能を確認した。









3 原子力研究における計算科学の役割 🕬

- ◆原子力の研究開発において、計算科学技術は理論、実験と並ぶ 第3の柱
- ◆研究開発の効率化の原動力
 - > 原子力分野においては、大規模な施設・設備を用いる実験または長期間を要する実 験や、観測が困難な現象について、計算科学技術はその解明や予測を可能にする
- ◆原子力施設の安全評価に必要不可欠
 - > 計算科学技術は、原子炉材料・核燃料の経年変化、原子力施設の耐震性等の詳細 を評価・予測可能にする
 - ▶ 安全性・健全性評価に最新知見を取り入れるための必須技術
- ◆ 国際競争力強化のカギを握る技術
 - > 原子力分野への計算科学技術の応用により、「実験・設計の繰り返しによる開発」から 「高精度シミュレーション予測に基づく開発の効率化」へと技術革新を可能にする
 - ※ 米DOE, 欧EURATOMにおいても、原子炉の設計、安全、材料に係るシミュレーション 技術開発等を実施

apan Atomic Energy Agenc

原子力委員会からの要請 → 計算科学技術を活用した原子力研究開発の必要性 ● シミュレーション技術の高度化等による大規模な技術システム開発の効率化 > 原子力政策大綱 ◆ 最高水準の計算機シミュレーション能力を整備・維持・発展 ● 開発システムの適合性を早い段階から計算機によるモデリング・シミュレーション技術等を駆使して多面的に評価 > 第5回原子力委員会資料第1号(独立行政法人日本原子力研究開発機構の次期中期目標の策定について(見解)) ◆ 近年急速に進歩してきた情報技術、モデリングをシミュレーション 技術を駆使して、性能目標を達成するシステム実現を、設計、 製造・建設、運転、廃棄に至るすべてのプロセスについて、研究

開発の可能な限り初期の段階から多面的かつ徹底的に検討す るフロントローディング型の研究開発活動の導入を検討 > 原子力委員会研究開発専門部会資料原子力政策大綱に示している原子力研 究開発に関する取組の基本的考え方の評価

an Atomic Energy













<section-header><section-header><section-header><figure><figure><figure>









1	スーノ	独立行政派 ペーコン	ま人化ない ピュー	。 びに二法 ・ 夕資派	人統合によ の整理	≈ 里•統合	(AEA
	H16年度	H17年度	H18年度	H19年度	H20年度	H21年度	H22年度~
原科研	VPP5000			Altix3700Bx2			BX900
那珂研	Origin3800	統合·合理化					
関西研	ITBL# 光量子シミュレー	十算機PRIMEPOWEI ーション専用計算機(R AlphaServer)	廃止			
大洗研			HPC25	00		\rightarrow	統合·合理化
250 200 50 50 0 ***	BX900 Altix3700Bx2 HPC2500 AlphaServer PRIMEPOWER Origin3800 VPPS000 6.4 18.2 18.2 18.2 18.2 18.2	: 15.5 15.5 1 \$ H34\$ R20\$\$ R20\$	200 5.5 14ġ H22年ġ~	◆原子力 及び大: ・平成1: ・平成1: した関 ・平成2: を統合 ・統合・	機構では、原 先研の4地 6年度末に属 8年度末にに 西研のシス・ 1年度末に属 1年度末に属 子力機 合理化計画	原料研、那毎 「マスパコン 原料研と那毎 は法定耐用4 テムを廃止。 原料研と大労 構における を完了。	I研、関西研 を運用 J研を統合。 F数を経過。 た研 (旧JNC 唯一のシス
			anon Atomi				







次 ⁻ Application D	世代コ Vevelopment Unit	ード開 for the Next Gen	月発部 eration Supercomputer
次世代コード開発部「FX1」は、12TFLOPS され、各ノードには富士道クアッドコアプロ 40GB/sという高メモリバンド幅を実現してい インターコネクトスイッチ (ハードウェ烈)を有1 いる特徴を生かし、次世代スーパーコンピュー	の理論演算性能 (ビーク) を セッサSPARC64™ VIIを1 ます。また、並列プログラム: こています。本システムでは、 <u>オ</u> - 夕の利用に向けたアプリケー	有するクラスタンステムです。: プロセッサ、主記値16GBを こおけるプロセス間の同期や着 にシステムの技術が以世代スー -ション開発を行います。	本システムは、300/ードから構成 活義化、専用チップセットにより 急和処理などを高速化する高機能 パーコンゼュータのペースとなって
	EX1		
	101		
	AND ALL ALL ALL ALL ALL ALL ALL ALL ALL AL	12mm 4.6m	ノード開ネットワーク構成器
		12 4.6 300 300 (1,200)	ノード製ネットワーク構成器
	12 ARRAN	12 4.6 300 300 (1,200) 8748C84	ノード読ネットワーク機道器 INSERTION
	1 2	12mm 4.6m 300mm 300mm 11,200mm 8PMC04 ^m W 4.0m	ノード開ネットワーク機成面 Market Water
	XXXXXXX XXXXXXX XXXXXXX XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX	12mm 4.6m 300mm 300mm 87MRCe4* W 4mm 10mm 40m	
	Normalization 1 1 2 <	12nus 4.6n 300num 11,200num 5940044" W 4.0n 10anus 40an 128nus	ノード菜ネットワーク電道室
	NUMBER NU	12 mins 4.6- 300	ノード菜ネットワーク菜店名 With With With With With With With With
		12-curs 4.6- 300 300 300 10 10 10 40 10 10 10 10 10 10 10 10 10	ノード菜ネットワーク菜店芝
	1 2	12 mins 4.6 m 300	J-HERTY-D-DERES
		12 mm 4.6 m 300mm 100000000	J-HRATHO-PRES
	11 2000 12 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000 2000	12 4.5 300 300 300 90 40 40 40 15	
	11 AUST 12 AUST 2011 AUST	12 mm 4.5 m 300m 300m 100mm (1,200m) 100mm (1,200m) 40m 128mm 64mm 128mm 64mm 128mm 128mm 15m 16m 16m 16m 16m	J-HRAYHO-ORES



	Linpack性能 [TFLOPS]	ピーク性能 [TFLOPS]	効率 [%]	実行時間 [h]	問題規模 [元]
BX900	191.40	200.1	95.66	13.40	2,401,344
FX1	11.66	12.9	90.37	5.68	709,632
BX90 原子力 地球シミュし 宇宙航空研究開発	OのLinpack性能は 機構 一々	♥、 国内1位(実行) (Top500	時事1位)、世行 2010.6) 110.6	22位(案行 前年 122.4	2112) 191.4
BX90 原子力: 地球シミュレ 全安航空研究開発	OのLinpack性能は 機構	ŧ、 国内1位(実行) (Top 500	地車1位)、世代 2010.6)	22位(実行地平 122.4	212) 191.4
BX90 原子力 地球シミュし 宇宙航空研究開発 東京大学(1 理化学研究	0のLinpack性能は 機構 	#、 国内1位(実行) (Top500		22位(案行 前平 122.4 〕	2位)
BX90 原子力 地球シミュレ 宇宙航空研究開発 東京大学(1 理化学研 東京工業	0のLinpack性能は 機構 	・ 国内1位(実行) (Top500	2010.6) 2010.6) 110.6 101.7 97.9 87.0	22位(実行発平 122.4 〕	2位)
BX900 原子力 地球シミュレ 宇宙航空研究開発 東京大学(1 理化学研 東京工業 施 植動会科学研	0のLinpack性能に 機構 一々 一々 一線構 [2次) 二次版 二次版 大学 二、 大学	• 国内1位(実行) (Top500	2010.6) 2010.6) 110.6 101.7 97.9 87.0 77.3	22位(実行発平 122.4 〕	2位)





- 国内	内機関のスーノ	パーコンピュータ	初用支援体制	ij —
原子力		「ラムの最適化(高速化 支援体制の下、研究者	と)や利用相談対応 行のサポートと計算機	など、充実した 残資源の有効利用
機関	利用相談	プログラム最適化	プログラム作成	備考
JAEA	(5)	(7)	(6)	()内は要員数
A機関	•			
B機関	•			
C機関	•	•	•	プログラム最適化は相談・ チューニング支援レベル
D機関	•	•		
E機関	•			
F機関	•	•		
G機関	٠			
H機関	•			
		Japan Atomic Energy	Agency	28

5 f	曾大	、する計算需要への対応 📟
(_	大型計算機利用課題評価制度 =
大型計算4 「大型計算 分	幾資源 「機利」	の効率的・効果的利用を促進するために、平成17年度より、 用委員会」を設置し、大規模利用課題を評価、計算時間を配
対象計1 課題由請	<mark>泉殿</mark> 条件	BX900(資源の約90%を承認課題に配分) 年間 250,000コア時間超、または256並列超
1010/02-1-1113	1	計算コードの実行効率、並列化効率
isti kur	2	「研究成果の期待度」、「必要性及び緊急性の程度」、「計算時間 設定の妥当性」の3項目について委員会による5段階評価
评恤	3	年度末に利用結果「得られた成果の重要度、達成度及び費用 (CPU時間)対効果」について委員会による5段階評価。
		次年度への継続課題については、この結果を次年度評価に反映
		Jopan Atomic Energy Agonov







1	原子力	コプロク	ブラノ	の特徴
コード	対象	手法	並列化	特徴
GT5D	プラズマ乱流	差分法 有限要素法 FFT	Hybrid	高演算効率 (メモリアクセスがボトルネック)
NOPIC	レーザー伝播	差分法	FlatMPI	1対1通信(exchange)
exact-diag	固体中の電子	固有値問題 (CG)	Hybrid	集団通信(scatter)
TPFIT	圧縮性流体	行列解法 (ILUCGS)	FlatMPI	低キャッシュヒット率 (メモリアクセスがボトルネック)
AQUA	非圧縮性流体	差分法 (ICCG)	OpenMP	(共有メモリ型)
演算性能	、メモリバンド幅	及び通信バン	ド幅のバラ	ランスのとれた計算機が必要
		Japan Atomic	Energy Agen	cy







T5Dによるハイブリ	Jッド処理性能 べ	シチマーク		
同題リイス・240x	ED/11 79CE/00m	19-9MD D/E-0	5) 4006 ラマ 唐田	1
インテルFortran:-03,	LF (11,12GF/COFE フラットMPI(コード全体	ットラーのMD,D/F-0 にの完全OpenMP化は	・5/1,4050ユノ (安片 未対応)	1
富士通Fortran:-Kfast	t -Kocl -x200 -Kpara	allel,OMP, ハイブリッド	並列(MPI,自動並列(≿, OpenMP)
	BX900 Intel-SMP1	BX900 Fujitsu-SMP1	BX900 Fujitsu-SMP4	BX900 Fujitsu-SMP8
	10231.35	10397.27	9286.57	10928.65
全体/step (msec)				
全体/step (msec) 通信/step (msec)	2414.72	2608.57	1516.63	1533.82
全体/step (msec) 通信/step (msec) メモリ使用量(GB)	2414.72 3641.14	2608.57 3754.09	1516.63 1545.49	1533.82 1249.91
全体/step (msec) 通信/step (msec)	2414.72	2608.57	1516.63	1533.82













	國内 (5)	00TflopsQ(L)				外
	九州大学	スーパーコンピュータ			イリノイ大	BlueWaters
1)		システム		1)	総演算性能:1~10Pflops	稼働開始:2011年
	総演算性能:500Tflops	稼働開始:2011年/5月以降		-	ローレンス・リバモア	Sequoia
				2)	総演算性能:20Pflops	稼働開始:2012年
	東京大学	大規模超並列 スーパーコンピューター システム	大規模超並列		NASA	Pleiade
2)				3)	総演算性能:10Pflops	稼働開始:2012年
	総演算性能:1Pflops	稼働開始:2011年/10月以降			ロスアラモス	Cielo
				4)	総演算性能:1Pflops	稼働開始:2010年
a)	東京工業大学	クラウド型グリーン		5)	ロスアラモス	Roadrunner ?
3)		X-X-1761-9			総演算性能:20Pflops	稼働開始:2012年
	総演算性能:2.4Pflops	稼働開始:2010年		6)	EU	PRACE
					総演算性能:?Pflops	稼働開始:2010年
	理化学研究所	京		7)	EU	ITER-BA
I)					総演算性能:1Pflops	稼働開始:2012年
	総演算性能:10Pflops	稼働開始:2012年/6		8)	中国	Nebulae
					総宗質析体·3Dflone 9	彩励明44·2010年





科	学	技	術	計	算	分	科	会		選	出
					2010 쇼	∓度 科	学技術	i計算分	·科ɗ	<u></u>	より

3 次元 MHD コードによる FX1 の性能評価

名古屋大学太陽地球環境研究所

荻野 竜樹

3次元MHDコードによる FX1の性能評価

荻野竜樹

名古屋大学太陽地球環境研究所

[アブストラクト]

並列型スパコン FX1 や HX600 の利用により、地球磁気圏のグローバル構造とメソスケールの境界層渦乱流の自己無撞着 な大規模高精度シミュレーションができる時代が訪れてきたが、その実現にはさらに高効率の並列計算コードの開発が不可 欠である。ここに、キャッシュヒット率向上と多次元領域分割/一括転送による通信量最小化を可能にした高効率の 3 次元 MHD コード開発に成功し、FX1 の 3072 コア数まで 15・20%以上の実行効率を得ることができた。Flat MPI と Hybrid (Impact)の比較では、Flat MPI が最速だったが、コア数によらず Hybrid もよい実行効率を得た。通信速度と計算速度の比 は 3072 コアでも 4%以下であり、通信速度は 10 万コア程度までボトルネックにならずにスケーラビリティが伸びることが 期待できる。もちろん、実際にテストして実証する必要がある。

[キーワード]

地球磁気圏のMHD シミュレーション、高効率の並列コード開発、多次元領域分割と通信量最小化、Flat MPI と Hybrid の 比較、スーパーコンピューティング

1. はじめに

スペースプラズマや核融合プラズマの最も特異な性質はその異常輸送現象にあるといえる。特に顕著な例と して温度の輸送や磁場の拡散は粒子間衝突から得られる古典的な輸送係数の千倍から 1 万倍に達することもあ ることが観測や実験から知られている。しかし、その基本的な物理過程の説明は今日でも一致した考えに達し ていないといえる。プラズマの異常輸送の有力な説明には二つある。一つは、マクローメソーミクロスケール が必ず結合していてメソスケール現象が決定的な役割を果たしているとする説である。もう一つは、マクロと ミクロスケールの直接結合が本質的でメソスケールは重要な役割を果たしていないとする説である。ここに、 最近の並列型スパコンの急速な進歩により MHD (電磁流体力学的) モデルからマクローメソスケール現象を自 己無撞着に解くことができるような時代が実際に訪れてきた。即ち、メソスケール現象がプラズマの異常輸送 現象に本質的な役割を果たしているかどうかに決着をつけるような大規模な 3 次元 MHD シミュレーションを 実行できる段階に達してきたといえる。もちろん、ミクロスケール現象は MHD モデルで扱うことができない ので、超粒子モデルやブラソフモデルなどの運動論的モデルによるシミュレーションが不可欠である。

そのマクローメソスケール結合を調べるには、太陽風と地球磁気圏相互作用のグローバル3次元 MHD シミ ュレーションが一つの典型的な対象となる。なぜなら、磁気圏のマクロ構造の中に物理量が急峻に空間変化す る境界層があり、そこにメソスケールのプラズマ乱流が発生してプラズマの異常輸送を引き起こしている可能 性が高く、グローバルモデルから境界層のプラズマ不安定とその非線形発展と結果としての渦乱流発生を自己 矛盾なく解ける段階に達したからである。しかし、その大規模 MHD シミュレーションを実行するためには大 規模並列型スパコンの効率を最大化するプログラムの開発が不可欠である。

2. MHD コードの並列化と超並列スパコンの有効利用

MHD モデルでは流体の方程式と Maxwell 方程式を高精度の数値計算法の一つである Modified Leap-Frog 法 で初期値境界値問題として解く。その3次元 MHD コードを MPI Fortran を用いて並列化する。その際、CPU 数が1千個を越えて1万個、10万個程度まで計算速度が劣化しないようなよいスケーラビリティが得られる並 列計算コードの開発を目指す。 私達は、これまで富士通の VPP5000/64、PRIMEPOWER HPC2500 1536core などで MPI を利用してそれ ぞれ 60%、20%を超える実行効率を出す 3 次元 MHD コードを開発してきた。その知識と経験を活かし、高効 率の並列計算 3 次元 MHD コードを開発して、平成 21 年度に名古屋大学情報基盤センターに導入された富士通 の FX1 3072core と HX600 2048core を用いてテストした。作成した MHD コードは基本的に 4 種類に分類で きる。それらは、1 次元領域分割法(1Da)、2 次元領域分割法(2Da)、3 次元領域分割法(3Da)およびキャ ッシュヒットを上げるために MHD 方程式の配列成分を第 1 変数に移動した 3 次元領域分割法(3Db)である。 2、3 次元領域分割法を用いると分散した CPU 間の通信量が小さくなり、これまでのテスト計算からは、スカ ラー並列機 HPC2500 では(3Db)が最速、ベクトル並列機 Earth Simulator では(2Da)が最速であった。

3. FX1 と HX600 における 3 次元 MHD コードの計算速度

FX1 は計算速度とメモリバンド幅が同程度との特徴があるスパコンで、Flat MPIと Hybrid (Impact) 両方 の利用ができて、Hybrid は 1 ノード 4 コアを共有メモリとしてスレッド並列化し、ノード間をプロセス並列 (MPI) で利用する。そのどちらが高速計算できるかが最初の注目点であったが、ほとんどの場合 Flat MPI が 最速になっている。ただし、Hybrid の場合もそれほど劣化しない。また、かなりバラつきはあるが、3072 コア 利用までスケーラビリティはよく保たれていた。前章の 4 種類の MHD コードを比較した場合、計算速度の速 い順に並べると (3Db) > (3Da) > (2Da) > (1Da) となり、(3Db) で 15-20%以上の実行効率を、(1Da) で 5-15%程度の実行効率を得ることができた。特に、(3Db) ではコア数が 3072 まで多くなってもほとんど計 算速度の劣化は起らず、バランスが良い場合は若干ではあるがむしろ計算効率が向上する場合もあった。JAXA の FX1 でテストした場合は、多数コア (4096、5760 コア) で Hybrid (Impact) が Flat MPI よりも高速にな る場合があった。これは大変興味ある結果であるが、MHD 変数の大きさや 3 次元方向の並列分割数がより適切 に取られていたかの疑問もあり、継続的なテストが必要である。また、計算時間と通信時間を分離して計測し たところ、通信時間と計算時間の比は 256-3072 コアの範囲で 3-4%であり、しかもその比は 3072 コアまであ まり変らないので、スケーラビリティは 1 万コアを超えてもよく成り立つことが期待される。しかし、実際は より多くのコア数で実証確認する必要がある。

HX600 も Flat MPI が最速であり、スレッド並列数を増やして Hybrid にすると計算速度がかなり劣化する傾向が明確であった。私のテストでは、最適値を求めた場合 HX600 は FX1 より常に少し計算速度が低めであった。一方、共同研究者のほとんどは逆の結果を得ていて、HX600 が FX1 よりも並列効率がより良いと結論付けている。その結果の違いがどこからきているかはまだ明確ではないが、一つの可能性として 1 コア当りのメモリ使用量、即ち粒度が関係していると予想される。

4. おわりに

学際大規模情報基盤共同研究・共同拠点のプロジェクトで他の大学の情報基盤センターのスパコンを利用す る機会を得た。それらはSX-9、HA8000、SR16000、TSUBAMEである。4 種類の MHD コードでテストした が、どのスパコンも期待したとおりのよい結果が得られ、スケーラビリティもテストしたコア数(3072 コア) の範囲でよく成り立っていた。気がついた点として、T2K のスパコン(HX600, HA8000)は(2Da)が最速 で(3Da)がそれに続き、(3Db)は逆に遅かった。一方、SR16000 は FX1 と同様に(3Db)が最速で両者は似 た傾向を示した。これらの結果は FX1 と SR16000 ではキャッシュのヒット率向上が重要であることを意味し ている。また、SR16000 は SMT が非常によく効き約 2 倍の速度がでた。一方、FX1 では SMT による計算速 度向上の効果は僅かであった。T2K のスパコンで(2Da)と(3Da)が速いのは「擬似ベクトル機能が効くため と推測しているが、理由究明にはさらなる系統的なテストが必要である。



名古屋大学太陽地球環境研究所 荻野竜樹

共同研究者

梅田隆行(名古屋大学太陽地球環境研究所) 深沢圭一郎(九州大学大学院理学研究院地球惑星科学部門)

アウトライン

- 1. 背景
- 2. 地球磁気圏の3次元グローバルMHDモデル
- 3. PRIMEPOWER HPC2500からFX1とHX600への移行
- 4. FX1とHX600おける3次元MHDコードの計算速度
- 5. 種々のスパコンの計算速度の比較と考察
- 6. 地球磁気圏の3次元MHDシミュレーションの将来
- 7. まとめ





太陽表面爆発現象と惑星間擾乱の伝播

太陽フレアなどの太陽活動はジオスペース環境変化の源





「ようこう」による太陽の軟X線画像

惑星間空間を伝播するCME(コロナ質量放出)

磁気圈·電離圈·熱圈結合:領域間結合

太陽風のエネルギーは磁気圏・電離圏・熱圏の領域間結合を介 して地球内部に運ばれ、極域でオーロラを輝かせる。宇宙天気研 究はその領域間結合を調べ、予報に結びつける。



















基礎方程式:MHD方程式とMaxwell方	程式
$\frac{\partial \rho}{\partial t} = -\nabla \cdot (\mathbf{v}\rho) + D\nabla^2 \rho$	(1)
$\frac{\partial \mathbf{v}}{\partial t} = -(\mathbf{v} \cdot \nabla)\mathbf{v} - \frac{(\nabla p)}{\rho} + \frac{\mathbf{J} \times \mathbf{B}}{\rho} + g + \frac{\mathbf{\Phi}}{\rho}$	(2)
$\frac{\partial p}{\partial t} = -(\mathbf{v} \cdot \nabla) p - \gamma p \nabla \cdot \mathbf{v} + D_p \nabla^2 p$	(3)
$\frac{\partial \mathbf{B}}{\partial t} = \nabla \times (\mathbf{v} \times \mathbf{B}) + \eta \nabla^2 \mathbf{B}$	(4)
$\mathbf{J} = \nabla \times (\mathbf{B} - \mathbf{B}_d)$	(5)
	15

スカラー並列機 並列計算のcr 通信量を大幅	で高速計算を ou数が大幅に に減らす	を実現するために 曽加
3次元MHDコー	トの領域分割	制法による亚列化
	計算時間	通信時間
1次元領域分割	Ts=k ₁ N ³ /P	Tc=k ₂ N ² (P-1)
2次元領域分割	$Ts=k_1 N^3 / P$	Tc=2k ₂ N ² (P ^{1/2} -1)
<mark>3次元領域分割</mark>	$Ts=k_1 N^3 / P$	Tc=3k ₂ N ² (P ^{1/3} -1)
k1とk2:一定の N:3次元配列(係数 の1方向の変数量	
P: 並列CPUの	数	16





	3次元領域分割の方法	
СС	CMPI START 一括転送の2次元バッファを導入	•
	parameter (npex=2,npey=2,npez=2) parameter (npe=npex*npey*npez,npexy=npex*npey) integer itable(-1:npex,-1:npey,-1:npez)	
С	parameter(nzz=(nz2-1)/npez+1)	
	parameter(nyy=(ny2-1)/npey+1)	
	parameter(nxx=(nx2=1)/npex=1) parameter(nxx3=nxx+2,nyy3=nyy+2,nzz3=nzz+2)	
С	dimension $f(n) = 0$ (ny) (+1 0 (ny) (+1 0)	
	dimension ftemp1x(nb,nyy3,nzz3),ftemp2x(nb,nyy3,nzz3) dimension ftemp1y(nb,nxx3,nzz3),ftemp2y(nb,nxx3,nzz3) dimension ftemp1z(nb,nxx3,nyy3),ftemp2z(nb,nxx3,nyy3)	
С		19









名古屋大学情報基盤センターの新	所スーパーコンピュータ
 システムS1 M9000 大規模な共有メモリ 3/ードx128コア (1/ードはFX1のフロントエンド) ノード当り:演算性能:1.28TFlops,メモリ:1TB シミュレーション結果の解析と画像処理に利用 	大規模シミュレーション はシステムS3 (FX1) Hybrid (Impact)の利用 を推奨
 システムS2 HX600 クラスター型計算機(T2K型)ノード内は共有メモリ 160ノード x 16コア ノード当り:演算性能:160GFlops,メモリ:64GB 中規模シミュレーションに利用 	総ノード/コア数:768/3072 総演算性能:30.72TFlops 総メモリ:24TB
 システムS3 FX1 大規模分散並列型、次世代スパコンとの連携 768ノード×4コア ノード当り:演算性能:40GFlops,メモリ:32GB メモリバンド幅:40GB/s 大規模シミュレーションに利用 	FX1利用の計算規模 コア数 256 512 1024 メモリ(TB) 2 4 8 演算性能 2.56 5.12 10.24 (Tflops)
	24



Flat MPIとHybrid (user並列+自動並列) はどちらが最速か?	
Hybridでは	
「Node内はベンダーが最高効率の自動並列化を	
提供する	
Node間はユーザーが最高効率のプロセス並列	
プログラムを作成する	
この結果、最高効率のプログラムが得られるは	
ずである???」	
この命題は正しいか?	
それともFlat MPIが最速か?	
	2





スーパースカラー並列スパコンで最高効率の 並列プログラムを作成するには	
1. 先ず、プログラムの構造の確立が最も重要	
 ・計算での変数参照を元にプログラム構造を決定 配列変数定義の最小化 計算の明確なブロック化 計算と通信の完全な分離 定数と変数は意識して区別(dependencyを解く) a(i)=a(i+j)の形で j は定数で定義 	
 ・配列の大きさ、do loopの回数(初数と終数)、 配列内の変数は定数で定義 定数定義はparameter文を利用 	

29

スーパースカラー並列スパコンで最高効率の 並列プログラムを作成するには(その2)

- ・領域分割法の導入(2、3次元領域分割) 通信では一括転送を実施
- ・Dependencyを可能な限り取り除く 計算のメイン部分ではIf文や条件文は絶対に使用不可 定数定義の利用
- Mask、List vector、Gather and scatterの最適な使用 ・変数の一括変換と一括転送
 - これを条件文やdo loopを導入せずに実行

If文や条件文が並列化(擬似ベクトル化)できたとの表示に だまされてはいけない。→→ その効率は実際非常に悪い。 (必ず実測すること)

30

コンピュータ	グリッド数	総CPU CORE 数	ブロセス 並列数	自動 並列 数	CPU 時間 (sec)	1グリッド 当りの 計算時 間 (nsec)	120% 当りの 計算時間 × 総CPU CORE数
HPC2500	1024*1024*1024	512	512	1	2.487	2.316	1186
HPC2500	2048*2048*2048	1024	512	2	10.763	1.253	1283
VPP5000	1002*1002*1120	56	56	1	5.794	5.152	289
SX6	512*256*256	16	16	1	0.693	20.662	331
Earth Sim	2048*1024*1024	1024	1024	1	0.62	0.289	296
HX600	512*512*512	64	64	1	1.451	10.813	692
HX600	1024*1024*1024	512	512	1	1.906	1.775	909
HX600	2048*2048*2048	512	512	1	19.654	2.288	1171
FX1	512*512*512	64	64	1	1.337	9.964	638
FX1	1024*1024*1024	512	512	1	1.377	1.283	657
FX1	1024*1024*1024	1024	1024	1	0.752	0.701	684
FX1	2048*2048*2048	512	512	1	10.525	1.225	627
FX1	2048*2048*2048	1024	1024	1	5.215	0.607	622
SR16000	508*508*508	64	64	1	1.236	9.424	603
SR16000SMT	508*508*508	32	64	1	1.239	9.448	605
SR16000	1018*1018*510	256	256	1	1.256	2.376	608
SR16000SMT	1018*1018*1018	256	512	1	1.29	1.223	626



FX1の計算速度 1											
コン ピュータ	プログラ ム	グリッド数	総コア 数	プロセス 並列数	自動 並列数	CPU時間 (sec)	計算速度 (Gflops)	効率 (%)			
FX1	1Da	1024**3	256	64	4	152.0	9	0.4			
FX1	1Da	1024**3	128	128	1	232.2	6	0.5			
FX1	1Da	1124*1024**2	128	128	1	28.4	52	4.1			
FX1	2Da	1024**3	128	128	1	17.1	79	6.2			
FX1	2Da	1124*1024**2	128	128	1	8.23	180	14.1			
FX1	3Da	1024**3	128	128	1	7.32	185	14.4			
FX1	3Db	1024**3	128	128	1	5.22	259	20.2			
FX1	2Da	3072*2048**2	3072	3072	1	17.2	1004	3.3			
FX1	2Da	3072*2048**2	3072	768	4	19.7	822	2.7			
FX1	3Da	3072*2048**2	3072	3072	1	3.91	4145	13.5			
FX1	3Da	3072*2048**2	3072	768	4	4.36	2782	9.1			
FX1	3Db	3072*2048**2	3072	3072	1	3.29	4932	16.1			
FX1	3Db	3072*2048**2	3072	768	4	3.76	4321	14.1			
1Da:1次元 グリッド数0	記領域分割、 の選び方に、	2Da:2次元領域: よっては極端に計	分割、3Da 算効率が	a:3次元領: 劣化する場	<mark>域分割、3</mark> 合がある	Db:3次元領 が、少し調節	域分割+f(m すると改良て	n,i,j,k) きる。 ₃₃			

コンピュータ	プログラ ム	グリッド数	総コア 数	プロセス 並列数	自動 並列数	CPU時間 (sec)	計算速度 (Gflops)	効率 (%)
FX1	3Db	1024**3	1024	1024	1	0.645	2096	20.5
FX1 (SMT)	3Db	1024**3	1024	1024	1	0.661	2043	20.0
FX1	3Db	1024**3	1024	258	4	0.732	1846	18.0
FX1 (SMT)	3Db	1024**3	1024	258	4	0.783	1810	17.7
FX1	3Db	1024**3	2048	2048	1	0.355	3813	18.6
FX1	3Db	1024**3	2048	512	4	0.380	3555	17.4
FX1 (SMT)	3Db	2048**3	1024	1024	1	5.09	2127	20.8
FX1	3Db	2048**3	2048	2048	1	2.72	3979	19.4
FX1	3Db	2048**3	2048	512	4	2.88	3752	18.3
FX1	3Db	3072* 2048**2	3072	3072	1	2.70	5997	19.5
FX1	3Db	3072* 2048**2	3072	768	4	2.92	5556	18.1

コンピュー タ JAXAJAX A JAXA	プログ ラム	グリッド数	総コア 数	プロセス 並列数	自動 並列数	CPU時間 (sec)	計算速度 (Gflops)	効率 (%)				
FX1 (SMT)	3Db	2048**3	256	256	1	19.8	541	21.3				
FX1 (SMT)	3Db	2048**3	1024	1024	1	5.09	2127	20.8				
FX1	3Db	2048**3	2048	2048	1	2.72	3979	19.4				
FX1	3Db	3072*2048**2	3072	3072	1	2.70	5997	19.5				
FX1	3Db	2048**3	512	512	1	10.53	1028	20.1				
	3Db	1024**3										
JAXA FX1	3Db	1024**3	1024	1024	1	0.701	1929	18.8				
JAXA FX1	3Db	1024**3	4096	4096	1	0.220	6160	15.0				
JAXA FX1	3Db	1024**3	4096	1024	4	0.211	6406	15.6				
JAXA FX1	3Db	1000*1200**2	5760	5760	1	0.247	7338	12.7				
JAXA FX1	3Db	1000*1200**2	5760	1440	4	0.228	7966	13.8				
1Da:1次元領 多コア数の言	域分割、 †算も高效	2Da:2次元領域分 動率、JAXA FX1で	う割、3Da では多コア	:3次元領域 でHybridカ	成分割、3D Flat MPI	Db:3次元領は よりも高効率	載分割+f(m, の場合もあっ	i,j,k) ot: 35				

HX600の計算速度									
コンピュータ	プログラ ム	グリッド数	総コア 数	プロセス 並列数	自動 並列数	CPU時間 (sec)	計算速度 (Gflops)	効率 (%)	
HX600	2Da	1024**3	256	256	1	4.95	273	10.7	
HX600	3Da	1024**3	256	256	1	3.56	379	14.8	
HX600	3Db	1024**3	256	256	1	5.22	259	10.1	
HX600	3Db	1024**3	2048	2048	1	0.464	2912	15.9	
HX600	3Db	1024**3	2048	512	4	0.685	1972	9.6	
HX600	3Db	1024**3	2048	128	16	0.894	1511	7.4	
HX600	3Da	1024**3	2048	2048	1	0.406	3327	16.3	
HX600	2Da	1024**3	2048	2048	1	0.674	2005	9.8	
HX600	2Da	2048**3	2048	2048	1	6.27	1724	8.4	
HX600	3Da	2048**3	2048	2048	1	3.54	3047	14.9	
HX600	3Db	2048**3	2048	2048	1	4.99	2217	10.6	
HX600	3Da	2048**3	2048	128	16	3.54	151	1.5	
2Da:2次元 大粒度計算 効率がかが	元領域分割 算ではHX6 なり低くなる	、3Da:3次 00はFX1よ 。	:元領域分 :り計算効)割、3Db:)率は低い	3次元領 、また、自	域分割+f(r 目動並列数を	n,i,j,k) を上げると言	ł算	

















	◆ 種々のスーパーコンピュータの計算速度の比較										
並列ベクトル計算機から超並列スカラ計算機まで - ベクトル機: Fujitsu: VPP5000、NEC: SX-6、SX-8R、SX-9、ES - スカラ機(x86除く): Fujitsu: HPC2500、FX1、PRIMQUEST、Hitachi: SR16000 - X86機: Hitachi HA8000、Fujitsu HX600、Sun: TSUBAME、Fujitsu PRIMERGY											
		Core數 /CPU數		理論性能 (GFlops)		領域分割					
S	K-9	64/64	2058	6553	31	2次元	ベクトル				
H.	48000	8192/1024	10038	75366	13	3次元A	Opteron				
H	X600	1024/256	2166	10240	21	3次元A	Opteron				
F)	K 1	1024/256	2081	10240	21	3次元B	SPARC64VII				
SI	R16000	1344/672	5375	25267	21	3次元B	POWER6				
		SX-9	HA8000	HX600	FX1	SR16000					
	@core	1.0000	0.0400	0.0658	0.0632	0.1244					
	@CPU	1.0000	0.1595	0.2631	0.2528	0.2487	45				

新しい技術について

SMT(Simultaneous Multithreading)は有効か? Core内で複数の実行スレッドを同時に実行するプロ セッサの機能 Node – Multi-Core – Multi-Thread

日立SR16000ではSMT効果が顕著に現れる 富士通FX1ではSMT効果は微少

SMTやコア内複数演算機がいつも機能するかは疑問

99.99%を超える並列化率と20%程度以上 の絶対効率を得るためのヒント

- 1. プログラムのメイン部分で全てのdo loopを並列化 (1個でも残したらダメ)
- 2. 全てのdo loopで並列化OKの表示が出ても安心しては いけない(if文や条件文は低効率の並列化をやっただけ かも知れない)ほとんどの場合数倍程度
- if文や条件文(dependency)を並列化するためには、 Mask, List vector, Gather and scatterの中から 最も適切なものを最も適切な方法で使う必要がある
- 実測する場合には<u>粒度(1個のCPUのメモリ占有率)</u> を揃える必要がある
- 5. 最善のテストは最大数のCPUを用いて、粒度を意識して 絶対効率を実測すること

今後解明すべきこと

 Core数を1-10万個と更に増やした場合もscalability は保たれるか?(どうしたら予測できるか?)

- 2. 通信は本当にボトルネックになるのか?この限界を 明確に示すにはどうすればよいのか?
- 3. Cacheのヒット率を更に上げるにはどうするのか?
- Flat MPIとHybrid (プロセス並列+自動並列)は Core数が増えた時、どちらが高効率なのか?
- 5. スパコンの種類によって高効率のプログラムは 異なるのか(FX1型とT2K型)?







ه ا

•

2010Year

50



FX1とHX600計算速度

コン ピュータ	プログラ ム	グリッド数	総コア 数	プロセス 並列数	自動 並列数	CPU時間 (sec)	計算速度 (Gflops)	効率 (%)
FX1	1Da	1024**3	256	64	4	152.0	9	0.4
FX1	1Da	1024**3	128	128	1	232.2	6	0.5
FX1	1Da	1124*1024**2	128	128	1	28.4	52	4.1
FX1	1Da	602*402*402	256	64	4	0.417	294	11.5
FX1	1Da	902*402*402	256	64	4	0.704	261	10.2
FX1	1Da	1124*1024**2	128	128	1	8.23	180	14.1
FX1	3Da	1024**3	128	128	1	7.32	185	14.4
HX600	1Da	512**3	64	64	1	4.06	42	6.5
HX600	1Da	1024**3	256	64	4	9.81	138	5.4

1Da:1次元領域分割+f(i,j,k,m) m: MHD方程式の8個成分 2Da:2次元領域分割+f(i,j,k,m)

3Da:3次元領域分割+f(i,j,k,m)3Db:3次元領域分割+f(m,i,j,k) 1次元領域分割も配列数を適切に選べばかなりの計算速度が得られる。

55

FX1の計算速度:1次元コード:f(i1)

・1次元配列f(i1)MHDコードはメモリを最小化したプログラム ・領域分割法を利用できないので新しいユーザー(プログラム) 分割法の開発が必要

コンピュータ	グリッド数	総コア	プロセス	自動	CPU時間	計算速度	効率		
		数	並列数	並列数	(sec)	(Gflops)	(%)		
FX1	62*32*32	4	1	4	0.0227	3.5	8.8		
FX1	602*302*302	4	1	4	14.4	4.8	12.0		
FX1	802*402*402	4	1	4	32.6	5.0	12.5		
FX1	1002*502*502	4	1	4	65.7	4.8	12.1		
FX1	50*50*26	4	1	4	0.0238	3.4	8.6		
FX1	402*402*202	4	1	4	8.78	4.7	11.7		
FX1	602*602*302	4	1	4	28.9	4.8	11.9		
FX1	802*802*402	4	1	4	66.7	4.9	12.2		
$f(1): 11=1+n1^{(1)}+n2^{(K-1)}+n3^{(m-1)}, n1=nx2, n2=n1*ny2, n3=n2*nz2$									

56