

A performance evaluation of CCS QCD Benchmark on the COMA (Intel® Xeon Phi™ , KNC) system

Taisuke Boku^{a,b}, Ken-Ichi Ishikawa^{*c,d}, Yoshinobu Kuramashi^{a,e}, Lawrence Meadows^f, Michael D'Mello^f, Maurice Troute^f, Ravi Vemuri^f

^aGraduate School of Systems and Information Engineering, University of Tsukuba, Tsukuba, Ibaraki 305-8573, Japan

^bCenter for Computational Sciences, University of Tsukuba, Tsukuba, Ibaraki 305-8577, Japan

^cGraduate School of Science, Hiroshima University, Higashi-Hiroshima, Hiroshima 739-8526, Japan

^dCore of Research for the Energetic Universe, Hiroshima University, Higashi-Hiroshima, Hiroshima 739-8526, Japan

^eFaculty of Pure and Applied Sciences, University of Tsukuba, Tsukuba, Ibaraki 305-8571, Japan

^fIntel Corporation, USA

E-mail: ishikawa@theo.phys.sci.hiroshima-u.ac.jp

The most computationally demanding part of Lattice QCD simulations is solving quark propagators. Quark propagators are typically obtained with a linear equation solver utilizing HPC machines. The CCS QCD Benchmark is a benchmark program solving the Wilson-Clover quark propagator, and is developed at the Center for Computational Sciences (CCS), University of Tsukuba. We optimized the benchmark program for a Intel® Xeon Phi™ (Knights Corner, KNC) system named “COMA (PACS-IX)” at CCS Tsukuba under the Intel Parallel Computing Center program. A single precision BiCGStab solver with the overlapped Restricted Additive Schwarz (RAS) preconditioner was implemented using SIMD intrinsics, OpenMP and MPI in the offload mode. With the reverse-offloading technique, we could reduce the communication and offloading overheads. We observed a performance of ~ 200 GFlops sustained for the Wilson-Clover hopping matrix multiplication on the lattice sizes larger than $24^3 \times 32$ on a single card of the COMA system. A good weak scaling performance was observed on the local lattice sizes larger than $24^3 \times 32$.

34th annual International Symposium on Lattice Field Theory

24-30 July 2016

University of Southampton, UK

*Speaker.

1. Introduction

The success of Lattice QCD simulations owes much to the development of numerical algorithms and optimization for the quark solver, and evolution of HPC machines. We have developed a quark solver benchmark program called “CCS QCD Benchmark” (CCS-QCD) [1], which solves the Wilson-Clover quark propagator, at the Center for Computational Sciences (CCS), University of Tsukuba. This is designed to be as simple as possible and is written in plain Fortran 90 so that new algorithms or new HPC architectures can be evaluated quickly with this benchmark program.

A new architecture system equipped with the Intel[®] Xeon Phi[™] (Knights Corner, KNC) co-processor cards has been installed at CCS in 2014. The name of the system is “COMA (PACS-IX)” [2]. This is the ninth system of the PACS/PAX series [3]. The Intel[®] Xeon Phi[™] (Knights Corner, KNC) co-processor cards is based on the Intel Many Integrated Core[™](MIC) architecture, and has many physical cores compatible to x86-64 on a chip. Although the programming model is common to x86-64 based systems, it requires some tuning tips to fully extract the many core performance. The communication among the co-processor cards requires HOST-HOST and HOST-KNC communication like GPGPU computing. There have been a lot of studies on the QCD program for KNC systems in the past few years [4, 5]. We optimize the CCS-QCD program for the COMA system to extract the best performance. The basic strategy to optimize the CCS-QCD program for the KNC system, such as prefetching, threading, SIMD-vectorization *etc.*, is almost the same as those studied in Refs. [4, 5]. This year a first result using the next generation Intel[®] Xeon Phi[™] system (Knights Landing, KNL) has been presented in this conference [6]. In this talk, we especially focus on the parallel performance of the CCS-QCD on the COMA system.

There are three running modes for a typical KNC system in executing a MPI program; “native”, “symmetric”, and “offload modes”. We employ the “offload mode” for the CCS-QCD to utilize the single precision acceleration to the solver algorithm, where the single precision solver is added and involved as the preconditioner to the double precision solver. The directive based programming is available as the Language Extensions for Offload (LEO) in the Intel[®] compiler. The total amount of the code modification on the original code is minimized by offloading the single precision solver to the accelerator. The performance of the whole program relies on the performance of the single precision solver added.

To have the best performance, together with the tuning and the optimization for the computational part on the co-processor, the MPI communication among the co-processor cards in the offload mode must be considered as the MPI functions cannot be used in offload regions. Typically the MPI functions and manipulating data are handled by the host CPU code in the offload mode. The data transfer from a host CPU to the co-processor on the host and vice versa can be done only at the beginning or end of the offload region using the directive. The MPI communication splits the entire solver code into many parts of the offload regions. The offloading overhead could be a bottleneck of the performance.

To get rid of the limitation in the offload mode and reduce the offloading overhead, we implement a proxy server code running on the host CPU which handles the request of the MPI-communication from the offload region. The communication between the host proxy and the offloaded code on the KNC is done via the Symmetric Communications InterFace (SCIF) [7]. With the proxy code, the entire code of the single precision solver can be packed in a single offload

region and the MPI requests are reversely offloaded to the host proxy. This strategy is called reverse-offloading. The proxy and reverse-offloading have been introduced in Ref. [4]. This strategy also enables us to apply the communication-computation overlapping. The single precision solver is programmed in this way. The tuned code for the COMA system is also available at [1].

This paper is organized as follows. In the next section, we mainly describe the details of the reverse-offloading and the communication-computation overlapping. In section 3 we show the performance of the code and summarize this paper.

2. Tuning the CCS-QCD for the COMA system

The COMA system is composed of 393 computational nodes equipped with two CPUs (Intel® Xeon E5-2680v2) and two Xeon Phi™ 7110P co-processor (KNC) cards. All the nodes are connected by full-bisection bandwidth of Fat-Tree network of InfiniBand FDR. The theoretical peak performance is 1.001 PFlops including 157.2 TFlops of CPUs. Making a full use of the co-processors (84% of the system peak) is inevitable to get the best performance of the system.

The CCS QCD Benchmark (CCS-QCD) implements the BiCGStab solver algorithm for the even-odd site preconditioned Wilson-Clover quark matrix in double precision. The code is written in Fortran 90 Language and parallelized in the X , Y , and Z directions using MPI. We replace the double precision BiCGStab solver algorithm to the double precision flexible BiCGStab solver algorithm [8]. We add a single precision BiCGStab solver, to be offloaded to the co-processor, as the preconditioner to the flexible BiCGStab. Hereafter, we refer to the single precision BiCGStab solver as the solver for simplicity, and we focus on the solver performance.

The single precision solver is further preconditioned with the even-odd site and the overlapped Restricted Additive Schwarz (RAS) preconditioning [9]. The solver is written in the C/C++ language to make use of the SIMD intrinsic functions of the Intel C/C++ compiler. The SIMD length is 16 for float and we embed four time slices of the spinor and gauge link fields into a SIMD vector (`__m512`); four time slices of a two-component spinor at a color index, and four time slices of first two column elements of a SU(3) matrix at a row-color index. The so-called SU(3)-reconstruction technique is employed. The many cores of the Xeon Phi™ are organized by OpenMP parallel threading. Loop blocking on the lattice sites and explicit prefetching is employed to enhance the use of the cache locality in a physical core. We use one co-processor within a MPI process. Thus two MPI processes are located on the COMA node. The co-processor identifier (0 or 1) is assigned to the MPI process according to the even/odd-ness of the MPI RANK.

2.1 Reverse-offloading

The reverse offloading [4] is implemented as follows. The solver code to be offloaded is essentially the same as that in the native mode except for the MPI communication part. Instead of calling the MPI APIs, the solver code calls wrapper functions similar to the MPI, in which the communication requests from the KNC are translated into the MPI-requests, and they are reversely offloaded to the proxy on the host CPU via SCIF. We refer to this communication API as SCIF simply.

To launch the solver to the co-processor and the proxy server on the host CPU simultaneously, asynchronous offloading is used. Figure 1 shows the code snippet launching the single precision

```

1 extern ``C`` void
2 assign_inv_mult_offl_ras_solver( const float      *kappa,
3                                 const float      *stol,
4                                 int               *iter,
5                                 const mic_wqf_eo *swe,
6                                 mic_wqf_eo *sve)
7 {
8     int offload_signal = 0xF;
9     const float skappal = *kappa;
10    const float stoll = *stol;
11    int iterl = *iter;
12
13    ////////////////////////////////////////////////////////////////////
14    // Asynchronous offloading of native code
15    // proxy server is running parallel to the native code
16    ////////////////////////////////////////////////////////////////////
17    static mic_wqf_eo *swel, *svel;
18    if (swel == 0)
19    {
20        #pragma offload target(mic:mic_targetid) \
21        nocopy(swel[0:1] : alloc_if(1) free_if(0) preallocated targetptr) \
22        nocopy(svel[0:1] : alloc_if(1) free_if(0) preallocated targetptr)
23    {
24        swel = (mic_wqf_eo *)_mm_malloc(sizeof(mic_wqf_eo), 64);
25        svel = (mic_wqf_eo *)_mm_malloc(sizeof(mic_wqf_eo), 64);
26    }
27    }
28
29    #pragma offload target(mic:mic_targetid) \
30    nocopy(scif mic) \
31    in( swe [0:1] : into(swel[0:1]) alloc_if(0) free_if(0) targetptr) \
32    out(svel[0:1] : into(sve [0:1]) alloc_if(0) free_if(0) targetptr) \
33    in(skappal,stoll) \
34    inout(iterl) \
35    signal(offload_signal)
36    {
37        assign_inv_mult_eoprec_wd_bicgstab_mic(&skappal,&stoll,&iterl,svel,swel);
38    }
39
40
41
42    ////////////////////////////////////////////////////////////////////
43    // process proxy
44    // receive requests from the native code located
45    // in the above offload region.
46    ////////////////////////////////////////////////////////////////////
47    process_cmds();
48
49    ////////////////////////////////////////////////////////////////////
50    // wait for finishing of offloading
51    ////////////////////////////////////////////////////////////////////
52    #pragma offload_wait target(mic:mic_targetid) wait(offload_signal)
53    ....

```

Figure 1: Offloading part in the host CPU code.

solver and run the proxy asynchronously, where “#pragma offload*” are the directive for offloading. The schematic diagram of the program behavior is shown in Figure 2.

The input spinor vector memory is allocated on the co-processor in the lines 19–27, where the pointers `swel` and `svel` are kept on the co-processor across the offload region. In the lines 29–40, the single precision solver is called on the co-processor after the data transmission (copy in) of the source vector to `swel` of the co-processor memory from `swe` of the host memory (the start point of the blue and red horizontal arrows in Figure 2). This offload region is asynchronous as indicated by the directive `signal(offload_signal)` and the control is non-blocking over the host CPU. Then, without waiting for the termination of the offload region, the host CPU calls `process_cmds` at the line 47, which is the proxy function, where the MPI requests from the solver `assign_inv_mult..._mic` located at the line 38 are processed (the blue and red horizontal arrows and the communication arrows between them in Figure 2). When the solver converges, the solver sends a termination signal to the proxy `process_cmds`. After receiving the termination command in the proxy, the control of CPU moves to the line 52, where CPU waits for the termination signal from the offload region asynchronously emitted in the lines 29–40. This includes the completion of the data transmission (copy out) of the solution vector to `sve` of the host memory from `svel` of the co-processor memory (the end point of the blue and red horizontal arrows in Figure 2).

2.2 Communication-computation overlapping

Using the reverse-offloading and the SCIF, the co-processor can use the communication functions similar to the non-blocking MPI and DMA transmission functions. With this function, we

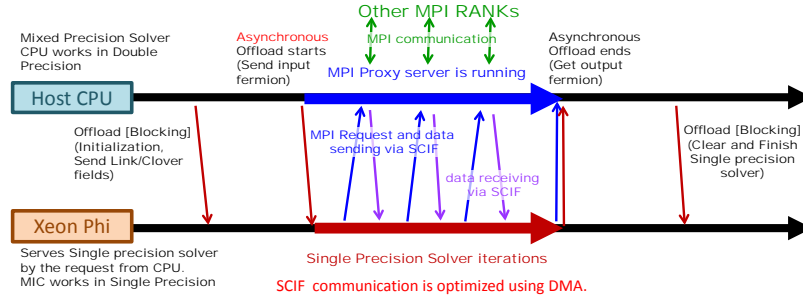


Figure 2: Reverse offloading in the solver.

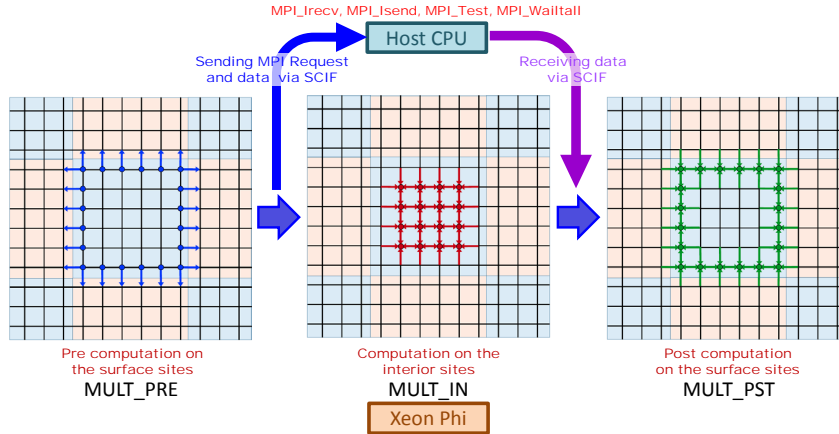


Figure 3: Task separation in the hopping matrix multiplication.

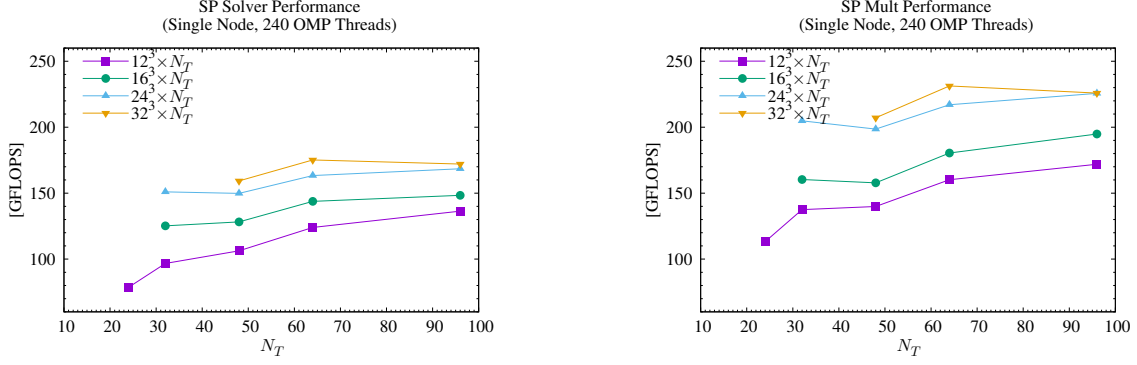
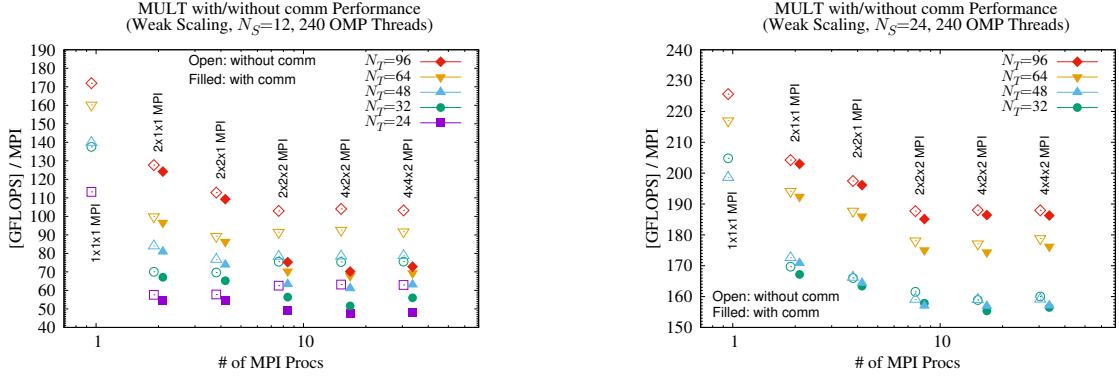
implement the communication-computation overlapping in the Wilson-Clover Dirac hopping matrix multiplication (MULT).

The implementation of the communication-computation overlapping for the the Wilson-Clover Dirac hopping matrix is usually organized by splitting the task into two pieces; (i) the computation within the node, and (ii) the computation using the data from other nodes. In this case the conditional branch statements are required to distinguish the site location and the direction of the hopping operation. To avoid the use of the conditional branch, as the KNC has a rather less performance on the conditional branch, we split the task in slightly different manner as depicted in Figure 3. The lattice sites in a process are split into the interior sites and the surface sites. In the “MULT_PRE” the spin-projection and data packing to the send buffer are carried out. After submitting the MPI request to SCIF in no-blocking way, the hopping computation in the interior sites continues in “MULT_IN”. In this region there is no conditional branches on the hopping directions. The communication request is processed on the host CPU during the computation. After receiving the data from the host CPU, the stencil computation on the surface sites follows in “MULT_PST”. With the task separation a good performance is expected in “MULT_IN”.

3. Results and Summary

We benchmark the tuned code on the COMA system. The performance of the hopping multiplication (MULT) is measured varying the local lattice size and the parallelism.

The single node performance, which is the baseline for the parallel computation, is shown in Figure 4. In this case there is no task splitting in the MULT computation and periodic boundary


Figure 4: Solver (left) and MULT (right) performance with single process.

Figure 5: Performance of MULT in weak-scaling test ($N_S = 12$ (left) and $N_S = 24$ (right)).

condition is imposed. Basic optimizations used for KNC systems are applied and we obtain ~ 200 GFlops for sufficiently larger lattice sizes (> 40000 sites) in a process.

The parallel performance is tested with the weak-scaling benchmarking. We test two local lattice sizes for the spatial size $N_S = 12$ and 24 , and vary the temporal local lattice size. The parallelisms tested are $1 \times 1 \times 1$ (baseline), $2 \times 1 \times 1$, $2 \times 2 \times 1$, $2 \times 2 \times 2$, $4 \times 2 \times 2$, and $4 \times 4 \times 2$ for the weak-scaling.

Figure 5 shows the performance of MULT. Filled symbols include the communication time of waiting for receiving data before `MULT_PRE`, open ones are the timing without the communication time. For the results with small local volume of $N_S = 12$ (left panel), a large gap between the open and filled symbols appears for the cases with three-dimensional parallelism. On the other hand, the performance degradation (open vs filled) is small for sufficiently larger local volume (right panel). This indicates the communication time is hidden behind the computation time. We need $N_S = 24$ for the communication hiding. We observe ~ 190 GFlops for the MULT performance at the local lattice size $24^3 \times 96$. The weak-scaling behavior is good as seen from the cases with three-dimensional parallel partitioning.

We have implemented the single precision solver code for Intel[®] Xeon Phi[™] (KNC) and applied it into the CCS QCD Benchmark. By using the reverse-offloading technique and the SCIF interface, we have achieved ~ 190 GFlops for the Wilson-Clover hopping multiplication with a $24^3 \times 96$ local lattice size on the COMA system. The next generation Intel[®] Xeon Phi[™] (KNL) system named ‘‘Oakforest-PACS’’ will appear soon in the Joint Center for Advanced High Performance Computing (JCAHPC) in Japan [10]. It could be interesting to optimize the CCS QCD Benchmark to the ‘‘Oakforest-PACS’’ system.

Acknowledgments

This work was supported by the Intel Parallel Computing Center at the Center for Computational Sciences (CCS), University of Tsukuba [11].

References

- [1] CCS QCD Benchmark program, Center for Computational Sciences (CCS), University of Tsukuba, <https://www.ccs.tsukuba.ac.jp/eng/research-activities/published-codes/qcd/>.
- [2] COMA System, Center for Computational Sciences (CCS), University of Tsukuba, <http://www.ccs.tsukuba.ac.jp/eng/research-activities/supercomputers/>.
- [3] History of PACS/PAX Series, Center for Computational Sciences (CCS), University of Tsukuba, <http://www.ccs.tsukuba.ac.jp/eng/research-activities/projects/ha-pacs/history/>.
- [4] B. Joó *et al.*, in *Supercomputing*, ser. Lecture Notes in Computer Science, J. M. Kunkel, T. Ludwig, and H. W. Meuer, Eds., Springer Berlin Heidelberg, 2013, vol. 7905, pp. 40–54, http://dx.doi.org/10.1007/978-3-642-38750-0_4.
- [5] S. Heybrock, B. Joó, D. D. Kalamkar, M. Smelyanskiy, K. Vaidyanathan, T. Wettig and P. Dubey, doi:10.1109/SC.2014.11 [arXiv:1412.2629 [hep-lat]]; R. Li and S. Gottlieb, PoS LATTICE **2014** (2015) 034 [arXiv:1411.2087 [hep-lat]]; H. Jeong *et al.*, PoS LATTICE **2013** (2014) 423 [arXiv:1311.0590 [physics.comp-ph]]; D. Barthou *et al.*, J. Phys. Conf. Ser. **510** (2014) 012005 [arXiv:1401.2039 [hep-lat]]; S. Mukherjee, O. Kaczmarek, C. Schmidt, P. Steinbrecher and M. Wagner, PoS LATTICE **2014** (2015) 044 [arXiv:1409.1510 [cs.DC]]; O. Kaczmarek, C. Schmidt, P. Steinbrecher and M. Wagner, arXiv:1411.4439 [physics.comp-ph]; Y. C. Jang *et al.* [SWME Collaboration], PoS LATTICE **2014** (2014) 035 [arXiv:1411.2223 [hep-lat]]; P. Arts *et al.*, PoS LATTICE **2014** (2015) 021 [arXiv:1502.04025 [cs.DC]]; P. Boyle, A. Yamaguchi, G. Cossu and A. Portelli, PoS LATTICE **2015** (2015) 023 [arXiv:1512.03487 [hep-lat]]; M. Schröck, S. Simula and A. Strelchenko, PoS LATTICE **2015** (2016) 030 [arXiv:1510.08879 [hep-lat]]; S. Heybrock, M. Rottmann, P. Georg and T. Wettig, PoS LATTICE **2015** (2016) 036 [arXiv:1512.04506 [physics.comp-ph]]; H. Kobayashi, Y. Nakamura, S. Takeda and Y. Kuramashi, PoS LATTICE **2015** (2016) 029.
- [6] C. DeTar, D. Doerfler, S. Gottlieb, A. Jha, D. Kalamkar, R. Li and D. Toussaint, arXiv:1611.00728 [hep-lat].
- [7] “Symmetric Communications Interface (SCIF) For Intel® Xeon Phi™ Product Family Users Guide”, Intel Corporation, http://registrationcenter-download.intel.com/akdlm/irc_nas/9669/scif_userguide.pdf.
- [8] H. Tadano and T. Sakurai, LSSC’07, Lec. Notes Com-put. Sci. 4818, 721 (2008); S. Aoki *et al.* [PACS-CS Collaboration], Phys. Rev. D **79** (2009) 034503 doi:10.1103/PhysRevD.79.034503 [arXiv:0807.1661 [hep-lat]].
- [9] Y. Osaki and K. I. Ishikawa, PoS LATTICE **2010** (2010) 036 [arXiv:1011.3318 [hep-lat]].
- [10] Joint Center for Advanced High Performance Computing (JCAHPC), <http://jcahpc.jp/eng/index.html>.
- [11] Intel® Parallel Computing Center at Center for Computational Sciences (CCS), University of Tsukuba, Intel Corporation, <https://software.intel.com/en-us/articles/intel-parallel-computing-center-at-center-for-computational-sciences-university-of-tsukuba>.