MILC Staggered CG Performance on Intel KNL

Carleton DeTar, Douglas Doerfler, Steven Gottlieb, Ashish Jha, Bálint Joó, Dhiraj Kalamkar, Ruizi Li*, Doug Toussaint



* : Presenter

Outline

- Staggered QPhiX library
- Benchmarks and performance results
- Conclusions and outlook

Staggered QPhiX library

- Developed from standard open source QPhiX library for Wilson quarks. Joint work by Intel, Jlab, and MILC collaboration
- Portable to various IA ISA's: SSE, AVX2, KNC, AVX512
- Supports MPI, OpenMP
- Current implementation on CG, supports both single and double precision
- Extending to gauge force

Staggered QPhiX library

• Data structure AoSoA :

*KS_Color_Vector[3][2][VECLEN], *Gauge[8][3][3][2][VECLEN], etc.

Data layout: EO checkerboarded, 4D or 3D block for SP or DP, with size 2 along each dimension to fill up the inner VECLEN array, i.e., data at (x, y, z, t), (x+Nx/2, y, z, t), (x, y +Ny/2, z, t), (x+Nx/2, y+Ny/2, z, t), ... are stored consecutively.

Hardware configuration

- Intel Endeavor Cluster:
 - Intel Xeon Phi[™] processor 7210 & 7250 (KNL): 64 & 68 cores @ 1.3 & 1.4 GHz, 8 or 16 GB MCDRAM, 6x16 GB DDR4 @ 2.1 & 2.4 GHz and 115 GB/s peak BW
 - Intel Broadwell(BDW) multi-core processor: Intel Xeon Dual Socket processor E5-2697 v4, 18 Cores/Socket, 36 Cores @ 2.3 GHz, 128GB DDR4 @ 2.4 GHz
- NERSC Cori Cluster:
 - Intel Haswell(HSW) multi-core processor, based on Cray's XC40 architecture: 2 sockets, 16 cores/socket @ 2.3 GHz, 128 GB DDR4 @ 2.1 GHz, 1.92 Pflops (theoretical peak)

Benchmarks

- Multi-mass CG using 9 or 11 masses
- Benchmarks :
 - QPhiX Dslash one KNL node
 - MILC+QPhiX vs. baseline MILC code, one and multi-node

QPhiX HISQ Dslash Performance



One KNL 7250: 64 cores; 2 threads/core except L = 8

Further improve: Compressed gauge; Software prefetches

7

MILC+QPhiX performance

- QPhiX staggered Dslash bandwidth on single KNL:
 - Increases w. increasing lattice volume(L > 12);
 - (Model BW) Hits 80% peak read bandwidth w. hardware prefetches.
- The next two slides compare MILC+QPhiX and Baseline MILC code on:
 - KNL 7250 one node;
 - KNL 7210 and Broadwell multi-node.





MILC+QPhiX vs. baseline MILC, single KNL 7250 weak scaling, up to 64 cores w. OpenMP Lattice volume 8³ x 24 per core

- a. All data in MCDRAM.
- b. MCDRAM used as cache, all data in DDR4 memory.
- c. DDR4 mode (no MCDRAM use).



1 node = 2 sockets, 16 cores/socket



MILC+QPhiX vs. baseline MILC w. MPI+OpenMP (baseline MILC code w. MPI on BDW), multi-node weak scaling, up to 16 nodes & 4D communications Lattice volume 24³ x 60 per node

Conclusions and outlook

- Conclusions:
 - Staggered QPhiX improves multi-mass CG performance by a factor of 1.5 ~ 2.5 in DP.
 - Benefits from MCDRAM, hyperthreading.
- Outlook:
 - Explore other NUMA modes: hemisphere, quadrant.
 - Omni-path network and various communication strategies for faster across-chip communications
 - Optimize other routines in the code.

Backup slides

• QPhiX staggered Dslash BW (DP)





MILC staggered multi-mass Conjugate Gradient(CG)

• Algorithm: shifted polynomials in Krylov space, e.g. *B. Jegerlehner, hep-lat/9612014*

$$(M - \sigma_i)a_i = b$$

where α_i is within a set of vector solutions, each with a bare quark mass σ_i , and b is the source color vector.

Update α_i with Dslash for smallest σ_j , and the rest α_j , $j \neq i$ with local linear algorithm.

Computational requirements

 $Flops = (1205 + 15 \times masses) \times iters \times V$

 $Bytes = ((171 + 12 \times masses) \times iters + 9$ $\times masses) \times V \times size of (complex)$

Where *masses*, *iters*, *V* are number of quark masses, CG iterations, and total lattice volume. (*masses* = 9 or 11 in our tests)

Intel Xeon Phi Knights Landing (KNL) architecture

HBM Modes



Image quoted from http://www.anandtech.com/show/9794/a-few-notes-on-intels-knights-landing-and-mcdram-modes-from-sc15

Intel Xeon Phi Knights Landing (KNL) architecture



Knights Landing Overview



Chip: 36 Tiles interconnected by 2D Mesh Tile: 2 Cores + 2 VPU/core + 1 MB L2

Memory: MCDRAM: 16 GB on-package; High BW DDR4: 6 channels @ 2400 up to 384GB IO: 36 lanes PCIe Gen3. 4 lanes of DMI for chipset Node: 1-Socket only Fabric: Omni-Path on-package (not shown)

Vector Peak Perf: 3+TF DP and 6+TF SP Flops Scalar Perf: ~3x over Knights Corner Streams Triad (GB/s): MCDRAM : 400+; DDR: 90+

Source Intel: All products, computer systems, dates and figures specified are preliminary based on current expectations, and are subject to change without notice. KNL data are preliminary based on current expectations and are subject to change without notice. IBinary Compatible with Intel Xeon processors using Haswell learners on States and are subject to change numbers are based on STREAN-like memory access pattern where MCERAL users and are memory. Results have been estimated based on internal Intel analysis and an entitient of the memory purposes only. Any ofference in system