

NEC SX-Aurora と格子QCD計算

T. Aoyama (KEK)

June 17th, 2020

7th HPC-Phys workshop (online)

Outline

- Overview of NEC SX-Aurora TSUBASA system
- Bridge++ Project for General-purpose Code Set of Lattice Gauge Theory Simulations
- Porting and Optimizing Bridge++ to SX-Aurora
- Summary

SX-Aurora TSUBASA: Overview

- **Newest product of NEC SX series**
 - Processors on PCIe card form factor, equipped in Xeon servers via PCIe Gen3.
- **Vector Engine (VE)**
 - Vector processor with 8 cores.
 - 64 vector registers of 16 kbit each.
 - HBM2 memory 6ch provides 1.2 TB/s bandwidth.



	SX-Aurora	Xeon CascadeLake	NVIDIA V100
# cores	8	28	5120
DP Performance	2.45 TFlops	2.42 TFlops	7.8 TFlops
Memory Capacity	48/24 GB	up to 1 TB	32/16 GB
Memory Bandwidth	1.2 TB/s	140 GB/s	900 GB/s
Memory Type	HBM2	DDR4	HBM2
Cache	16 MB shared	38.5 MB	6 MB
B/F	0.5	0.06	0.12

- **Vector Host (VH)**
 - Xeon server accommodates 8 VEs.
 - VHs interconnected by Infiniband EDR: up to 64 VEs in a Rack.

SX-Aurora TSUBASA

- **Programming Model**

- “VE execution model”
 - program runs on VE, as if on an ordinary node.
 - system calls e.g. I/O are offloaded to VH underneath.
 - cf. GPUs: offload tasks from host program to devices.
- Ordinary C/C++/Fortran programs run just by recompilation.
 - directives to control in detail.

- **Software Environment**

- NEC C/C++/Fortran Compilers with auto vectorization. OpenMP supported. MPI library provided.
- BLAS, LAPACK, and other optimized mathematical libraries.
- Profiler and Debugger.
- LLVM-based compiler also being developed.
 - intrinsics available. recent release supports auto-vectorization.
 - cf. LLVM-VE-RV project on github.

Bridge++: A Lattice QCD Code Set

- **Overview**

- General-purpose code set for simulations of Lattice Gauge Theory.
- Object-oriented design using C++.
- Development policy:
 - Readable, Extendable, Portable, and High-performance.



- **History**

- Launched in 2009, first public release in 2012.
- Latest version: 1.5.4 (March 2020).
- Adopted in research works, acknowledged in 48 papers.

- **Members**

- Y. Akahoshi, S. Aoki, Y. Namekawa (YITP), T. Aoyama, H. Matsufuru (KEK), I. Kanamori (RIKEN), K. Kanaya, Y. Taniguchi (Tsukuba), H. Nemura (RCNP), and contributors.

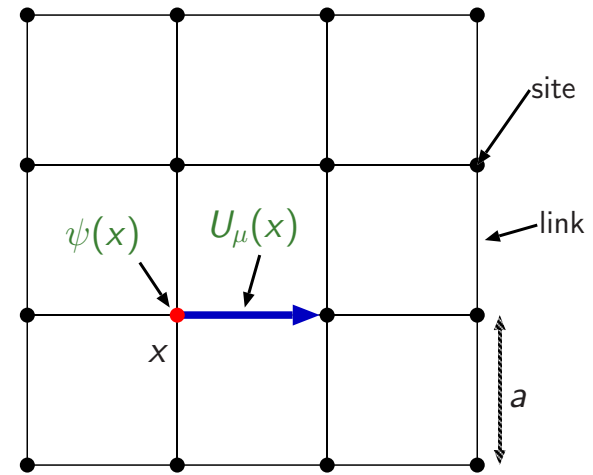
Bridge++: A Lattice QCD Code Set

- **Development of Bridge++**
 - Provides a tool set as a library for User applications, including:
 - Various Fermion and Gauge actions
 - Linear Solver algorithms, Linear algebraic operations
 - Simulation algorithms, Random numbers, I/O manipulations
- **Target platforms**
 - **“core” library:**
 - scalar processors, multicore cluster systems
 - OpenMP + MPI for parallelization
 - **Extensions:**
 - system-specific implementations
 - GPUs w/ OpenCL, OpenACC
 - Manycore processors and SIMD instructions, e.g. KNL, Intel Skylake
 - **Vector processors**

Members: Kanamori, Matsufuru, Namekawa, Aoyama

Hotspot in Lattice QCD

- QCD (Quantum ChromoDynamics) describes strong interaction among quarks and gluons.
- Numerical simulation of Lattice QCD formulated on discrete spacetime provides powerful tool.
 - Fermions (quarks): $\psi(x)$ on lattice site
 - Gauge fields (gluons): $U_\mu(x)$ on links
- Hotspot: inversion of “Fermion matrix”



$$D_{x,y} = \delta_{x,y} - \kappa \sum_{\mu} \left[(1 - \gamma_{\mu}) U_{\mu}(x) \delta_{x+\hat{\mu},y} + (1 + \gamma_{\mu}) U_{\mu}^{\dagger}(y) \delta_{x-\hat{\mu},y} \right]$$

$U_{\mu}(x)$: 3×3 complex matrix, γ_{μ} : 4×4 matrix, $\kappa = 1/2(4 + m)$

- Discretized differential operator with color and spin d.o.f
 → large sparse matrix of typically $\mathcal{O}(10^7) \times \mathcal{O}(10^7)$.
- Iterative linear equation solver algorithms employed
 → matrix-vector multiplication many times.

Porting and Optimizing Bridge++

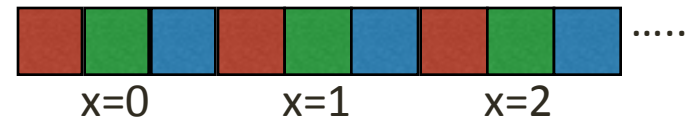
- **Strategy**

- Vectorization along site-loop.
 - rely on compiler's auto-vectorization.
 - promote loop unrolling for color/spin d.o.f. using constants.

- **Switching Data Layout**

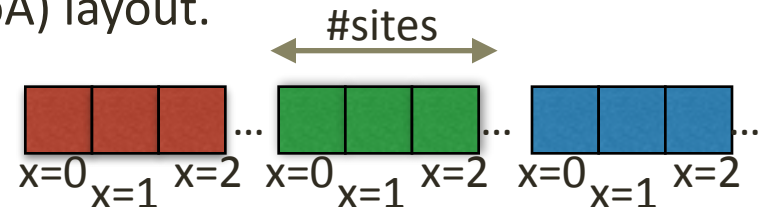
- “core” library: Array of Structure (AoS) layout.

- site d.o.f. packed innermost.



- Vector library: Structure of Array (SoA) layout.

- contiguous w.r.t. site loop index.



	core library (AoS)	Vector library (SoA)
Wilson mult on 1 core	1.07 GFlops	41.7 GFlops

Porting and Optimizing Bridge++

- **Parallelization within VE**

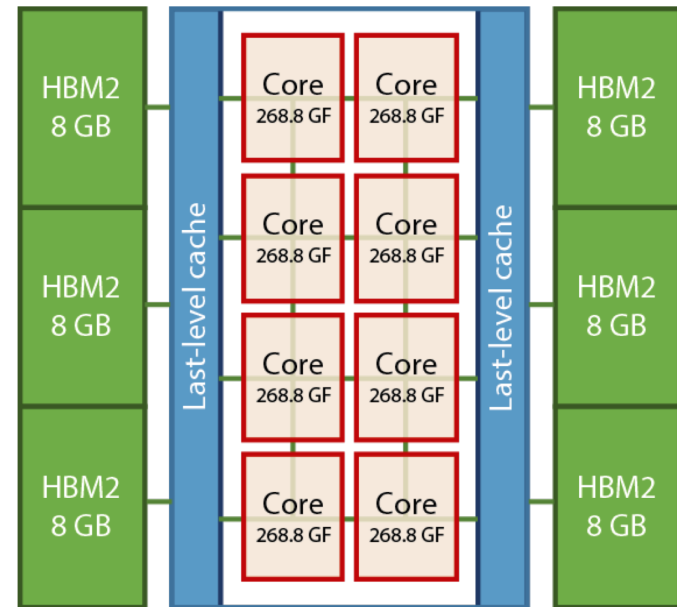
- Using 8 cores in a VE via flat-MPI.
- Performance:

	Single core 16x16x8x8	1 VE (8 cores) 16x16x16x32 / [1,1,2,4]
Wilson mult	41.7 GFlops	81.8 GFlops
peak performance	308 GFlops	2.42 TFlops
theoretical bandwidth	409 GB/s	1.2 TB/s

- B/F of Wilson mult ~ 2.2
 - Expected performance from memory bandwidth:
 - 180 GFlops (single core), 540 GFlops (1 VE)
- Profiling:
 - Vector instruction ratio: $\sim 99.90\%$
 - Average Vector Length: 256.0 \rightarrow seems well vectorized.

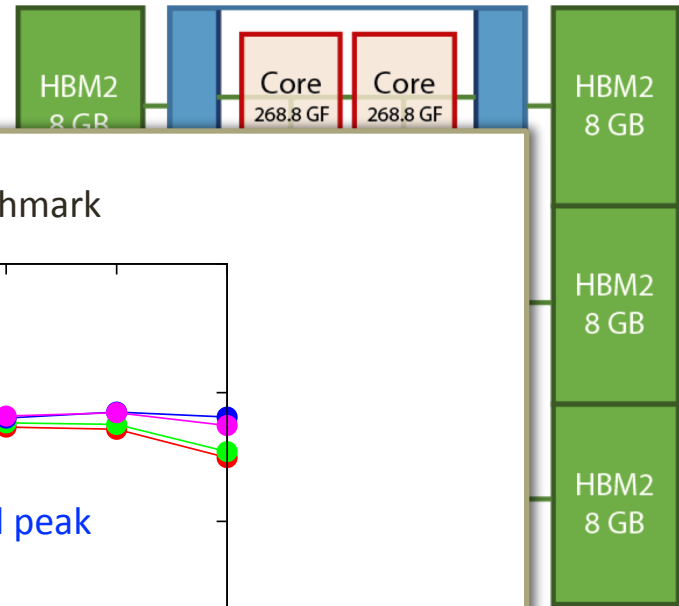
Further optimization

- Memory subsystem overview
 - 6 HBM2 modules connected.
 - 8 channels on each module.
 - 128bits interface x 8ch
x 1.6GHz x 6HBM2 = 1.2 TB/s
 - 128 byte-cells within module.
 - classified in 32 banks
 - access to contiguous 128 bytes will be most effective
 - VE equips 16 MB LLC shared by 8 cores.
 - connected by NOC (network on chip).
 - bandwidth between LLC and core = 409.6 GB/s.
- Bank conflict occurs by simultaneous access with 192KB strides
 - 128bytes x 6HBM2 x 8ch x 32banks (round-robin) = 192KB

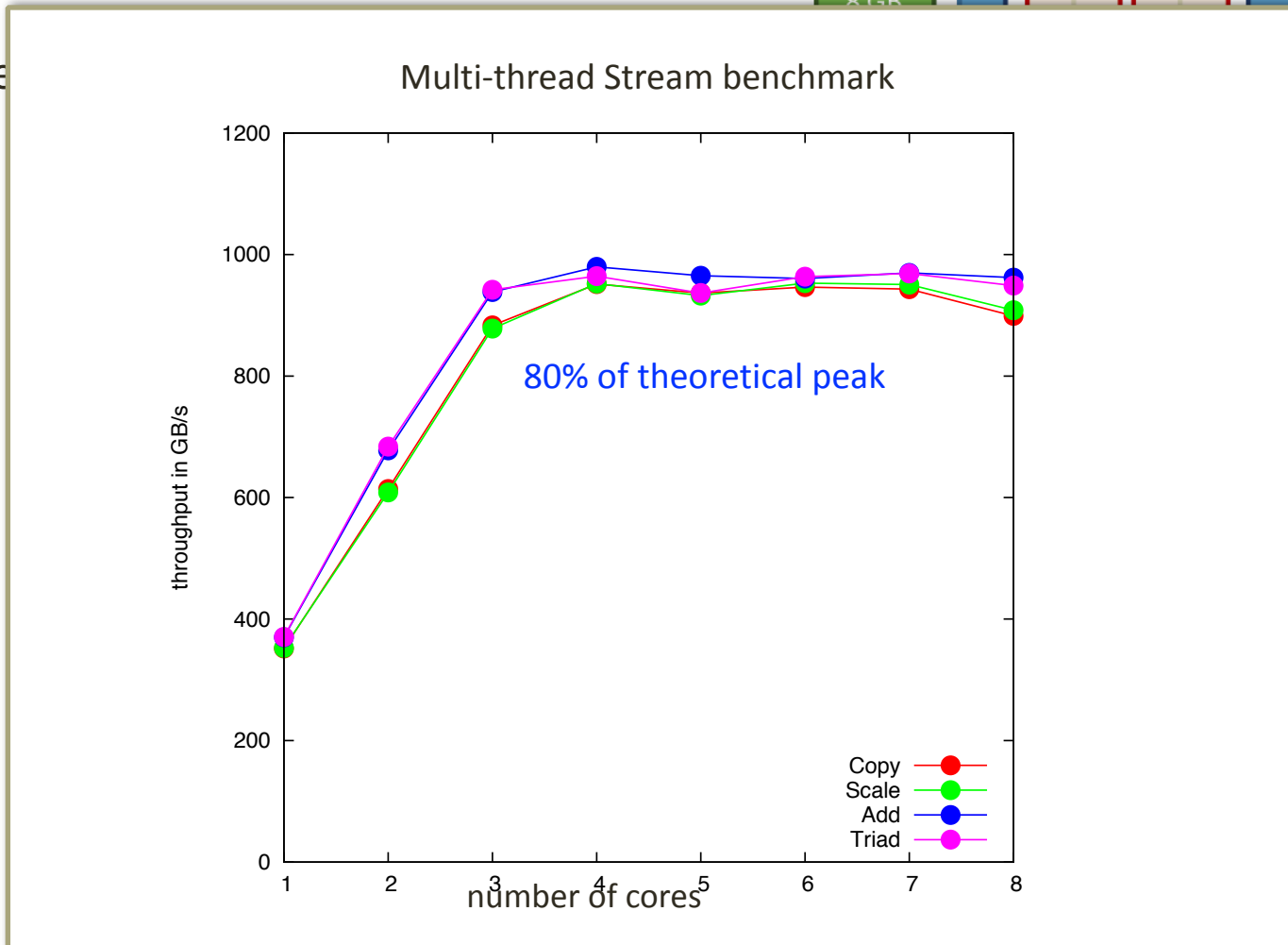


Block diagram of a vector processor

Further optimization



- Me



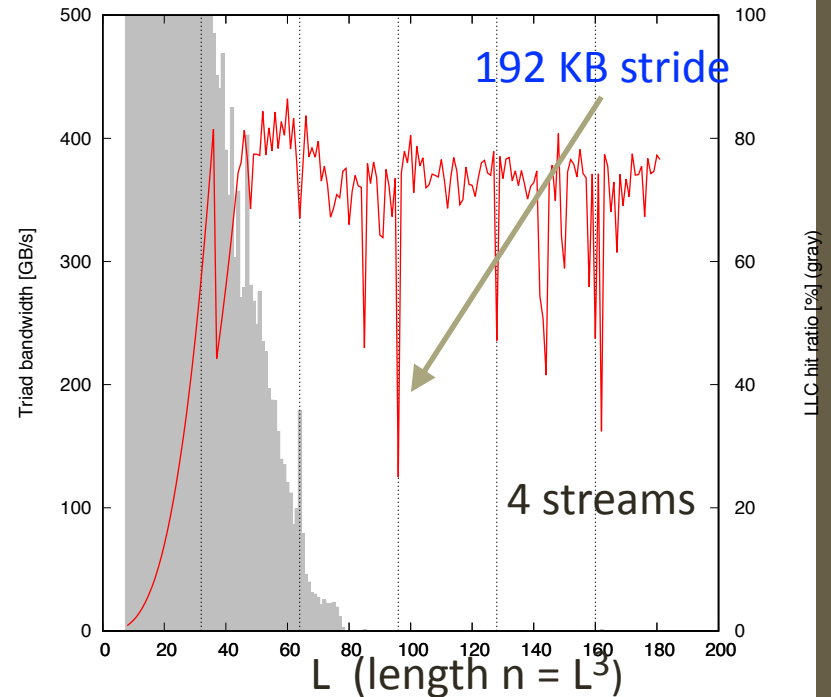
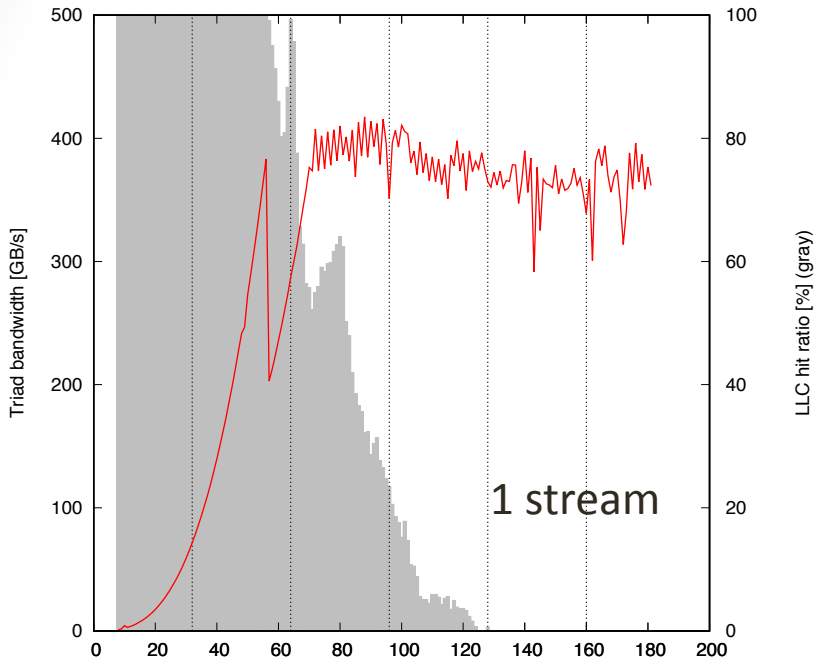
processor

- 128bytes x 6HBM2 x 8ch x 32banks (round-robin) = 192KB

Further optimization

- STREAM benchmark with multiple streams.

```
for (i=0; i<n; ++i) {  
    a0[i] = b0[i] + v * c0[i];  
    a1[i] = b1[i] + v * c1[i];  
    a2[i] = b2[i] + v * c2[i];  
    . . . .
```



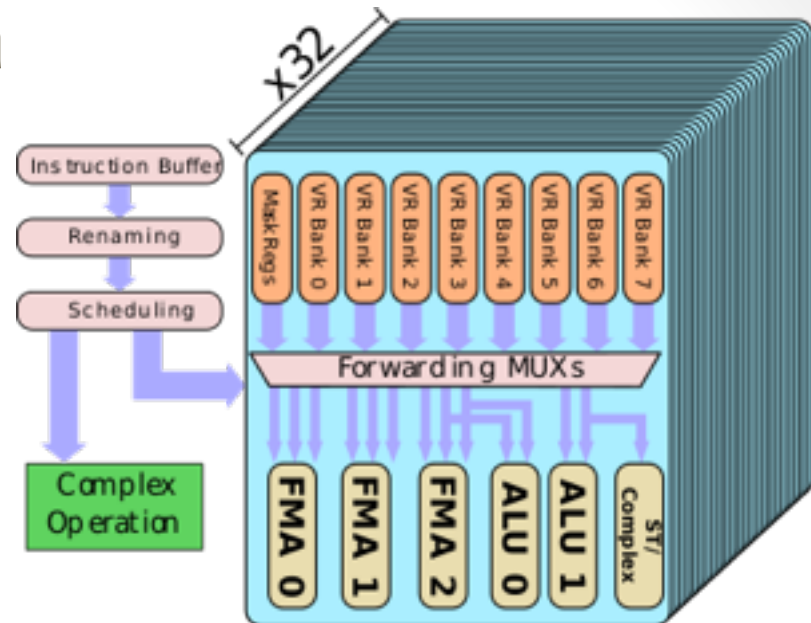
- Insert padding of appropriate size to avoid bank conflicts.
- Padding size to be chosen by examining memory access pattern.

x 1.5 for Wilson mult.

cf. Memory First: A performance tuning strategy focusing on memory access patterns, N. Ebata, R. Egawa, Y. Isobe, R. Takaki, H. Takizawa, Research Poster at ISC2019.

Further optimization

- Vector processing unit overview
 - 64 vector registers, vector length 256 elem. of 8B.
 - 32 vector pipelines (VPP).
 - 1 vector instruction execute 256 arith ops with 8 clock cycles.
 - 6 execution pipes:
 - 3 FMAs, 2 ALUs, 1 DIV/SQRT.



from WikiChip

- Invoke vector instructions: two-fold loops as an idiom.
 - Outer loop (maybe further parallelized over threads)
 - Inner vectorized loop (long enough i.e. ≥ 256 to fill vector registers)
 - apply blocking to inner loop in unit of vector length (VL=256), and specify compiler directives for optimization.

$(x,y,z,t) \rightarrow (x,y)$ and (z,t) for Wilson mult

Further optimization

- Wilson mult performance evolution:

	Single core	1 VE (8 cores)
core library (AoS)	1.07 GFlops	
Vector library (SoA)	41.7 GFlops	81.8 GFlops
insert padding		133 GFlops
two-fold loop, blocking, vector register directives		~200 GFlops
pack/unpack revised	70.1 GFlops	271.5 GFlops
peak performance	308 GFlops	2.42 TFlops

- About 20% (single core), or 10% (1 VE) of peak performance obtained.
- Compared to the expected performance from memory bandwidth, further improvement may be possible.

Summary

- Extension of Bridge++, a general-purpose code set for lattice simulations, to vector processors is being carried out.
- Data layout significantly affects the performance. Data should be aligned contiguously along site loop index (SoA format), efficient for load into vector registers.
- Further elaborations on memory access pattern, e.g. including padding to bank conflict improve performance significantly.
- An idiom to write vector loops will be two-fold loops, vectorized inner loop and outer loop. Further optimization to be applied e.g. loop blocking and compiler directives.
- Related work:
 - “Lattice QCD on a novel vector architecture”, B. Huth, N. Meyer, T. Wittig, LATTICE2019, arXiv:2001.07557[cs.DC].

素粒子原子核宇宙 シミュレーションプログラム



- 素粒子・原子核・宇宙物理学分野の大規模数値シミュレーションに基づく理論的研究を推進
 - これまでの KEK 「大型シミュレーション研究」 に続く共同利用プログラム
 - 計算基礎科学連携拠点 (JICFuS) によるサポートの下で共同利用を実施し、計算資源を提供
 - 大規模実験プロジェクトとの連携：実験データを精密に予言する理論計算により実験の成果の最大化を図る
- 研究課題を公募
 - 研究テーマを設定しグループを組織、代表者が申請
 - 公募要項については、KEK共同利用案内を参照
<https://www2.kek.jp/uskek/apply/pna-sp.html>
 - 一般利用の課題は随時受付け
- シミュレーションプログラムの情報は：
<http://research.kek.jp/group/pna-sp/>

KEK シミュレーションプログラム

