2019 Intel® eXtreme Performance Users Group (IXPUG) meeting

### Massively scalable computing method to tackle large eigenvalue problems for nanoelectronics modeling

Hoon Ryu, Ph.D. (E: <u>elec1020@kisti.re.kr</u>)

Korea Institute of Science and Technology Information (KISTI) KISTI Intel® Parallel Computing Center



## **Two Big Features of Advanced Device Designs**



**Miniaturization + Functionalization** 



- Size miniaturization of CMOS transistor
  - → Quantum / atomistic effects start to play!
  - $\rightarrow$  <u>New structures</u> to increase the TR density!

Intel eXtreme Performance Users Group (IXPUG) meeting / 2019 Jan.

#### Diversification in device functions

- → Devices dedicated to specific functions: Sensors, LEDs
- → No CMOS transistors: Need to explore the feasibility of <u>new materials and structures</u> beyond CMOS



## **Electronic Structure Calculations**

### **Prediction of electron motions in nanoscale materials**



• The state of electron motions in an electrostatic field created by the stationary nuclei.

 $\rightarrow$  Prediction of electron motions in nanoscale materials and devices

- Physics, Chemistry, Materials Science, Electrical Engineering
  - $\rightarrow$  Huge customers in the society of computational science



### **Electronic Structure Calculations**

### In a perspective of "numerical analysis"



- Two PDE-coupled Loop: Schrödinger Equation and Poisson Equation
- Both equations involve system matrices (Hamiltonian and Poisson)
  - $\rightarrow$  DOFs of those matrices are proportional to the # of grids in the simulation domains



- (Stationary) Schrödinger Equations

   → Normal Eigenvalue Problem
   HΨ = EΨ

   Poisson Equations
  - $\rightarrow$  Linear System Problem  $-\nabla(\epsilon \nabla V) = \rho \rightarrow Ax = b$

### **Electronic Structure Calculations**

### In a perspective of "numerical analysis"

- Two PDE-coupled Loop: <u>Schrödinger Equation</u> and <u>Poisson Equation</u>
- Both equation involve system matrices (Hamiltonian and Poisson)
  - $\rightarrow$  DOFs of those matrices are proportional to the # of grids in the simulation domains



• Schrödinger Equations • Normal Eigenvalue Problem  $(H)\Psi = E\Psi$ • Poisson Equations • Linear System Problem  $-\nabla(\epsilon\nabla V) = \rho \rightarrow Ax = b$ How large are these system matrices? Why do we need to handle those?

Intel eXtreme Performance Users Group (IXPUG) meeting / 2019 Jan.

5



# **Needs for "Large" Electronic Structures**



Electron motion happens in cores, but we need more

- 1. Quantum Simulations of "Realizable" Nanoscale Materials and Devices
  - $\rightarrow$  Needs to handle large-scale atomic systems (~ A few tens of nms)



# **Development Strategy: DD, Matrix Handling**

System matrices for Schrödinger and Poisson equations



Intel eXtreme Performance Users Group (IXPUG) meeting / 2019 Jan.



7

## **Development Strategy: Numerical Algorithms**

Schrödinger equations





## **Development Strategy: Numerical Algorithms**

PCC KISTI



**Poisson equations** 

#### Poisson Eqs. w/ CG Algorithm

- → A Problem of Solving Linear Systems
- Conv. Guaranteed: Symmetric & Positive Definite
- Poisson is always S & PD.

 $-\nabla(\varepsilon\nabla V) = \rho$ 

Steps for Iteration: Purely Scalable

Algebraic Ops.

We want to solve 
$$Ax = b$$
. First compute  $r_0 = b - Ax_0$ ,  $p_0 = r_0$   
loop for  $(j=1; j \le K; j++)$   
 $a_j \in \langle r_j \bullet r_j \rangle / \langle Ap_j \bullet p_j \rangle$ ;  
 $x_{j+1} \in x_j + a_j p_j$ ;  
 $r_{j+1} \in r_j - a_j Ap_j$ ;  
if  $(||r_{j+1}||/||r_0|| < e)$   
declare  $r_{j+1}$  is the solution of  $Ax = b$  and break the loop  
 $c_j \in \langle r_{j+1} \bullet r_{j+1} \rangle / \langle r_j \bullet r_j \rangle$ ;  
 $p_{j+1} \in r_{j+1} + c_j p_j$ ;  
end loop

### **Performance Bottleneck?**

### **Matrix-vector multiplier: Sparse matrices**

#### Vector Dot-Product (VVDot)



#### **Main Concerns for Performance**

- Collective Communication
  - → May not be the main bottleneck as we only need to collect a single value from a single MPI process
- Matrix-vector Multiplier
  - → Communication happens, but would not be a critical problem as it only happens between adjacent ranks
  - $\rightarrow$  Data locality affects vectorization efficiency







### Single-node Performance: Whole domain in MCDRAM

#### With no HBM With HBM (a) (b) **Description of BMT Target and Test Mode** • 10 CB states in 16x43x43(nm<sup>3</sup>) [100] Si:P quantum dot COMM COMM MVMUL → Material candidates for Si Quantum Info. Processors VVDOT VVDOT MEMOP MEMOP (Nature Nanotech. 9, 430) OTHERS OTHERS 3 → 15.36Mx15.36M Hamiltonian Matrix (~11GB) Speed-up Speed-up with HBM (a.u.) with HBM secs) • Xeon Phi 7210: 64 cores Wall-time (x10<sup>3</sup> s MCDRAM control w/ numactrl; Quad Mode H. Rvu, Intel® HPC Developer Conference (2017) **Results** With no MCDRAM 25 $2^{6}$ 2<sup>7</sup> Number of cores $\rightarrow$ No clear speed-up beyond 64 cores With MCDRAM $\rightarrow$ Up to ~4x speed-up w.r.t. the case w/ no MCDRAM $\rightarrow$ Intra-node scalability up to 256 cores 0 **Points of Questions** $2^{4}$ 25 $2^{6}$ 2<sup>8</sup> 2<sup>4</sup> 25 $2^{6}$ $2^7$ 2<sup>8</sup> 27 Number of cores Number of cores How is the performance compared to the one under $2^4$ cores = (1 MPI proc(s), 16 threads), $2^7$ cores = (2 MPI proc(s), 64 threads) other computing environments? (GPU, CPU-only etc..) $2^5$ cores = (2 MPI proc(s), 16 threads), $2^8$ cores = (4 MPI proc(s), 64 threads) → In terms of speed and energy consumption $2^{6}$ cores = (2 MPI proc(s), 32 threads)

### Speed in various computing platforms

#### **Description of BMT Target and Test Mode**

- 10 CB states in 16x43x43(nm<sup>3</sup>) [100] Si:P quantum dot
   → Material candidates for Si Quantum Info. Processors
  - (Nature Nanotech. **9**, 430)
  - → 15.36Mx15.36M Hamiltonian Matrix (~11GB)
- Specs of Other Platforms
  - → Xeon(V4): 24 cores of Broadwell (BW) 2.50GHz
  - → Xeon(V4)+KNC: 24 cores BW + 2 KNC 7120 cards
  - $\rightarrow$  Xeon(V4)+P100: 24 cores BW + 2 P100 cards
  - $\rightarrow$  KNL(HBM): the one described so far

#### Results

- KNL slightly beats Xeon(V4)+P100
  - $\rightarrow$  Copy-time (CPIN): a critical bottleneck of PCI-E devices
  - → P100 shows better kernel speed, but the overall benefit reduces due to data-transfer between host and devices
  - $\rightarrow$  CPIN would even increase if we consider periodic BCs
- Another critical figure of merit: Energy-efficiency









#### **Description of BMT Target and Test Mode**

- 10 CB states in 16x43x43(nm<sup>3</sup>) [100] Si:P quantum dot
   → Hamiltonian DOF: 15.36Mx15.36M (~11GB)
- Description of Device Categories
  - $\rightarrow$  Xeon(V4): 24 cores of Broadwell (BW) 2.50GHz
  - → Xeon(V4)+KNC: 24 cores BW + 2 KNC 7120 cards
  - $\rightarrow$  Xeon(V4)+P100: 24 cores BW + 2 P100 cards
  - $\rightarrow$  KNL(HBM): the one described so far

#### **Power Measurement**

- w/ RAPL (Running Ave. Power Limit) API
- Host (CPU+Memory), PCI-E Devices





When problem sizes exceed > 16GB?

#### Matrix-vector multiplier

```
for (unsigned int i = 0; i < nSize; i++) {</pre>
   double real sum = 0.0;
   double imaginary sum = 0.0:
   const unsigned int nSubStart = pMatrixRow[i];
   const unsigned int nSubEnd = pMatrixRow[i + 1]:
   for (unsigned int j = nSubStart; j < nSubEnd; j++) { 1, index</pre>
        const unsigned int nColIndex = pMatrixColumn[i]:
                                                         2. Matrix
        const double m real = pMatrixReal[j];
                                                            element
        const double m imaginary = pMatrixImaginary[i];
        const double v real = pVectorReal[nColIndex];
        const double v imaginary = pVectorImaginary[nColIndex];
        real_sum += m_real * v_real - m_imaginary * v_imaginary;
        imaginary sum += m real * v imaginary + m imaginary * v real;
   }
   pResultReal[i] = real sum;
   pResultImaginary[i] = imaginary_sum;
```

#### Data to be saved in memory

- index: column index of matrix nonzero elements (indirect index)
- · matrix and vector: matrix nonzero elements and vector elements

#### Intel eXtreme Performance Users Group (IXPUG) meeting / 2019 Jan.

#### Available options for MCDRAM utilization

- Cache mode: use MCDRAM like L3 cache
- Preferred mode: First fill MCDRAM then go to DRAM
   → numactl -preferred=1 ...
- Library memkind: use dynamic allocations in code
   → hbw malloc(), hbw free()







### When problem sizes exceed > 16GB?



#### Which component would be most affected by the enhanced bandwidth of MCDRAM?

• MVMul is tested with 11GB Hamiltonian matrix

Good for us – matrix is built upon the definition of geometry!

• Matrix nonzero elements drive the most remarkable performance improvement when combined w/ MCDRAM

### **Vectorization efficiency**

#### Matrix-vector multiplier: Revisit

```
for (unsigned int i = 0; i < nSize; i++) {</pre>
   double real sum = 0.0;
   double imaginary sum = 0.0:
   const unsigned int nSubStart = pMatrixRow[i];
   const unsigned int nSubEnd = pMatrixRow[i + 1];
   for (unsigned int j = nSubStart; j < nSubEnd; j++) {</pre>
        const unsigned int nColIndex = pMatrixColumn[j];
        const double m real = pMatrixReal[i]:
        const double m imaginary = pMatrixImaginary[i];
        const double v real = pVectorReal[nColIndex];
        const double v imaginary = pVectorImaginary[nColIndex];
        real sum += m real * v real - m imaginary * v imaginary;
        imaginary sum += m real * v imaginary + m imaginary * v real;
   }
   pResultReal[i] = real sum;
   pResultImaginary[i] = imaginary_sum;
```



• Efficiency of vectorization would not be super excellent

→ Vector elements should be "gathered" onto register before processing vectorization for matrix-vector multiplier

### **Vectorization efficiency**

#### Matrix-vector multiplier: Revisit

for (unsigned int i = 0; i < nSize; i++) {</pre> double real sum = 0.0; double imaginary sum = 0.0: const unsigned int nSubStart = pMatrixRow[i]; const unsigned int nSubEnd = pMatrixRow[i + 1]; for (unsigned int j = nSubStart; j < nSubEnd; j++) {</pre> const unsigned int nColIndex = pMatrixColumn[j]; const double m real = pMatrixReal[i]: const double m imaginary = pMatrixImaginary[i]; const double v real = pVectorReal[nColIndex]; const double v imaginary = pVectorImaginary[nColIndex]; real sum += m real \* v real - m imaginary \* v imaginary; imaginary sum += m real \* v imaginary + m imaginary \* v real; } pResultReal[i] = real sum; pResultImaginary[i] = imaginary\_sum;

IPCC KISTI

Assembly

Block 2:
leal (%rcx,%rdi,1), %r15d
vpadddy (%rll,%rl5,4), %ymmO, %ymml
kxnorw %kO, %kO, %kl
vpxord %zmm9, %zmm9, %zmm9
vpxord %zmmll, %zmmll, %zmmll
kxnorw %k0, %k0, %k2
vmovupsz (%rax,%r15,8), %zmm10
vmovupsz (%rl0,%rl5,8), %zmml2
add \$0x8, %edi
vgatherdpdz (%r9,%ymm1,8), %k2, %zmmll
vgatherdpdz (%r14,%ymm1,8), %k1, %zmm9
vmulpd %zmml1, %zmm10, %zmm8
vmulpd %zmm10, %zmm9, %zmm13
vfmsub231pd %zmm12, %zmm9, %zmm8
vfmadd231pd %zmm12, %zmm11, %zmm13
vaddpd %zmm4, %zmm8, %zmm4
vaddpd %zmm2, %zmm13, %zmm2
cmp %r8d, %edi
jb 0x417920 <block 2=""></block>

• Efficiency of vectorization would not be super excellent

→ Vector elements should be "gathered" onto register before processing vectorization for matrix-vector multiplier

### **Vectorization efficiency**

#### Matrix-vector multiplier: Revisit

```
for (unsigned int i = 0; i < nSize; i++) {</pre>
   double real sum = 0.0;
   double imaginary sum = 0.0:
   const unsigned int nSubStart = pMatrixRow[i];
    const unsigned int nSubEnd = pMatrixRow[i + 1];
   for (unsigned int j = nSubStart; j < nSubEnd; j++) {</pre>
        const unsigned int nColIndex = pMatrixColumn[j];
        const double m real = pMatrixReal[i]:
        const double m imaginary = pMatrixImaginary[i];
        const double v real = pVectorReal[nColIndex];
        const double v imaginary = pVectorImaginary[nColIndex];
        real sum += m real * v real - m imaginary * v imaginary;
        imaginary sum += m real * v imaginary + m imaginary * v real;
   }
    pResultReal[i] = real sum;
    pResultImaginary[i] = imaginary_sum;
```

Efficiency of vectorization would not be super excellent
 → Vector elements should be "gathered" onto register before proc

Intel eXtreme Performance Users Group (IXPUG) meeting / 2019 Jan.

#### Default (-O3) Vector length 2 Normalized vectorization overhead 1.020 Vector cost : 26.0 Estimated potential speedup: 1.70 AVX2 (-AVX2) Vector length 2 Normalized vectorization overhead 1.020 Vector cost : 24.5 Estimated potential speedup: 1.490 MIC-AVX512 (-xMIC-AVX512) Vector length 8 Normalized vectorization overhead 1.104 ier Vector cost : 8.370 Estimated potential speedup: 3.930



## **Extremely large-scale problems**

### In NURION computing resource

#### **NURION System Overview**



- 132 SKL (Xeon 6148) nodes / 8,305 KNL (Xeon Phi 7250) nodes
- Ranked at 13<sup>th</sup> in Top500.org as of 2018. Nov.
  - → Rpeak 25.7pFLOPS, Rmax, 13.9pFLOPS. https://www.top500.org/system/179421







## **Extremely large-scale problems**

### In NURION computing resource

#### **Description of BMT Target**



- $\rightarrow$  contains 400 million (0.4 billion) atoms,
- $\rightarrow$  Hamiltonian matrix DOF = 4 billion x 4 billion

#### **Computing Environment**

- Intel® Xeon Phi 7250 (NURION)
  - → 1.4GHz/68 cores, 96GB DRAM, 16GB MCDRAM (/node) → OPA (100GB)

#### Other Information for Code Compile and Runs

- Intel® Parallel Studio 17.0.5
- Instruction set for vectorization: MIC-AVX512
- MCDRAM allocation: numactl –preferred=1
- OPA fabric. 4 MPI processes / 17 threads per node

#### Memory Placement Policy Control

→ export I\_MPI\_HBW\_POLICY = hbw\_bind, hbw\_preferred, hbw\_bind (HBW memory for RMA operations and for Intel® MPI Library first. If HBW memory is not available, use local DDR) RMA: Remote Memory Access (for MPI communications)



```
export NUMACTL="numactl --preferred=1"
export I_MPI_HBW_POLICY=hbw_bind,hbw_preferred,hbw_bind
export I_MPI_FABRICS=ofi
```

ulimit -s unlimited cd \$PBS\_O\_WORKDIR cat \$PBS\_NODEFILE time mpirun \$NUMACTL ...

#### **Snapshot of a PBS script**





## **Extremely large-scale problems**



### In NURION computing resource





## Summary

### **KISTI Intel® Parallel Computing Center**

- Introduction to Code Functionality
- Main Numerical Problems and Strategy of Development
- Performance (speed and energy consumption) in a single KNL node
  - $\rightarrow$  Benefits against the case of CPU + 2xP100 GPU devices
- Performance in extremely huge computing environment
  - → Strong scalability up to 2,500 KNL nodes in NURION system
- (Appendix) Strategy of Performance Improvement towards PCI-E devices
- (Appendix) List of Related Publications

### Thanks for your attention!!

## Appendix: Strategy for offload-computing

Asynchronous Offload (for Xeon(V4) + KNC, Xeon(V4) + GPU)

### The real bottleneck of computing: Overcome with asynchronous offload [H. Ryu et al., Comp. Phys. Commun. (2016) (http://dx.doi.org/10.1016/j.cpc.2016.08.015)

- Vector dot-product is not expensive: All-reduce, but small communication loads
- Vector communication is not a big deal: only communicates between adjacent layers
- Sparse-matrix-vector multiplication is a big deal: Host and PCI-E device shares computing load



Intel eXtreme Performance Users Group (IXPUG) meeting / 2019 Jan.

of Supercomp,

# **Appendix: Strategy for offload-computing**

**Data-transfer and Kernel Functions for GPU Computing** 

### Data-transfer between host and GPU Devices

- 3x increased bandwidth with pinned memory
- Overlap of computation and data-transfer with <u>asynchronous streams</u>

### Speed-up of GPU Kernel Function (MVMul)

• Treating several rows at one time with WARPs Data-Access with a thread-base (no WARPs)



<u>H. Ryu et al., J. Comp. Elec. (2018)</u> (http://dx.doi.org/10.1007/s10825-018-1138-4)

#### [Synchronous Data Transfer with Pageable Memory]







## **List of Related Publications**



#### **Journals and Conference Proceedings**

#### **Journal Articles / Book Chapters**

- [1] H. Ryu, O. Kwon, Journal of Computational Electronics (2018) https://doi.org/10.1007/s10825-018-1138-4
  - → "Fast, Energy-efficient Electronic Structure Simulations for Multi-million Atomic Systems with GPU Devices",
- [2] S. Choi, W. Kim, M. Yeam, H. Ryu, International Journal of Quantum Chemistry (2018) https://doi.org/10.1002/qua.25622
  - → "On the achievement of high fidelity and scalability for large-scale diagonalizations in grid-based DFT simulations"
- [3] O. Kwon, H. Ryu, A Book Chapter in "High Performance Parallel Computing", InTechOpen (2018) <u>https://doi.org/10.5772/intechopen.80997</u>
  - → "Acceleration of Large-scale Electronic Structure Simulations with Heterogeneous Parallel Computing"
- [4] H. Ryu, Y. Jeong, J. Kang, K. Cho, Computer Physics Communications (2016) https://doi.org/10.1016/j.cpc.2016.08.015
  - → "Time-efficient simulations of tight-binding electronic structures with Intel Xeon Phi<sup>TM</sup> many-core processors

#### **Conference Proceedings / Presentations**

- [1] H. Ryu, K.-H. Hong (2018), Proceedings of IEEE SISPAD, <u>https://doi.org/10.1109/SISPAD.2018.8551719</u>
  - → "Optical Properties of Organic Perovskite Materials for Finite Nanostructures"
- [2] J. Kang, O. Kwon, J. Jeong, K. Lim, H. Ryu (2018), Proceedings of HPCS, https://doi.org/10.1109/HPCS.2018.00063
  - → "Performance Evaluation of Scientific Applications on Intel Xeon Phi Knights Landing Clusters"
- [3] H. Ryu, O. Kwon (2017), Top 20 Research Posters in GPU Technology Conference, http://www.gputechconf.com/resources/poster-gallery/2017/computational-physics
  - → "Q-AND: Fast, Energy-efficient Computing of Electronic Structures for Multi-million Atomic Structures with GPGPU Devices"
- [4] H. Ryu, Y. Jeong (2016), Proceedings of IEEE CLUSTER, https://doi.org/10.1109/CLUSTER.2016.76
  - → "Enhancing Performance of Large-scale Electronic Structure Calculations with Many-core Computing"