

GETTING READY TO AURORA EXA-SCALE SUPERCOMPUTER USING INTEL ADVISOR ROOFLINE ON INTEL CPUS AND GPUS

JAEHYUK KWACK
ALCF Perf. Engr. Group

ZAKHAR MATVEEV
Intel

OVERVIEW OF AURORA TESTBED SYSTEMS

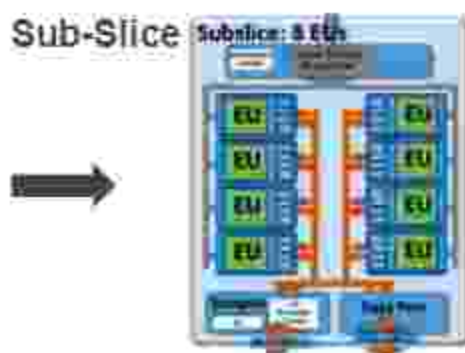
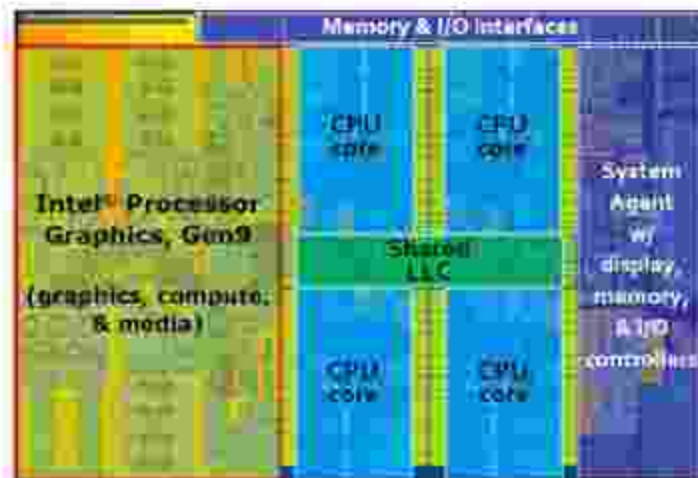
AURORA: A HIGH-LEVEL VIEW

- One of coming exa-scale systems supported by US DOE
 - 21 science/engineering project teams are actively working for exa-scale systems via ECP (Exascale Computing Project)
- Intel-HPE machine arriving at Argonne
 - Sustained Performance > **1 ExaFlops**
- Intel Xeon processor and Intel X^e GPUs
 - 2 Xeons (Sapphire Rapids)
 - 6 GPUs (Ponte Vecchio [PVC])
- Greater than 10 PB of total memory
- HPE/Cray XE platform and HPE Slingshot network
- Filesystem
 - Distributed Asynchronous Object Store (DAOS)
 - ≥ 230 PB of storage capacity
 - Bandwidth of > 25 TB/s
 - Lustre
 - 150 PB of storage capacity
 - Bandwidth of ~ 1TB/s

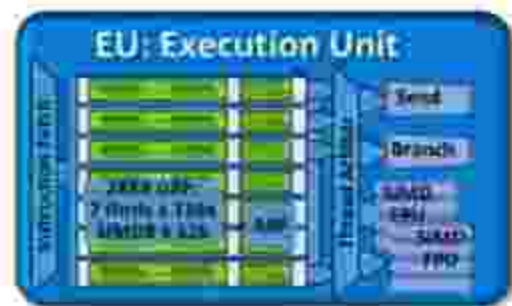


INTEL GEN9 GPU ON AURORA TESTBED

- Gen9: Double precision peak performance: 100-300 GF
 - Low by design due to power and space limits
 - Integrated GPU
- Hardware hierarchies
 - A GPU tile has multiple slices
 - A slice has multiple Sub-Slices
 - A sub-slice has multiple EUs



Execution Unit (EU)



ADVISOR ON INTEL CPUS AND GPUS

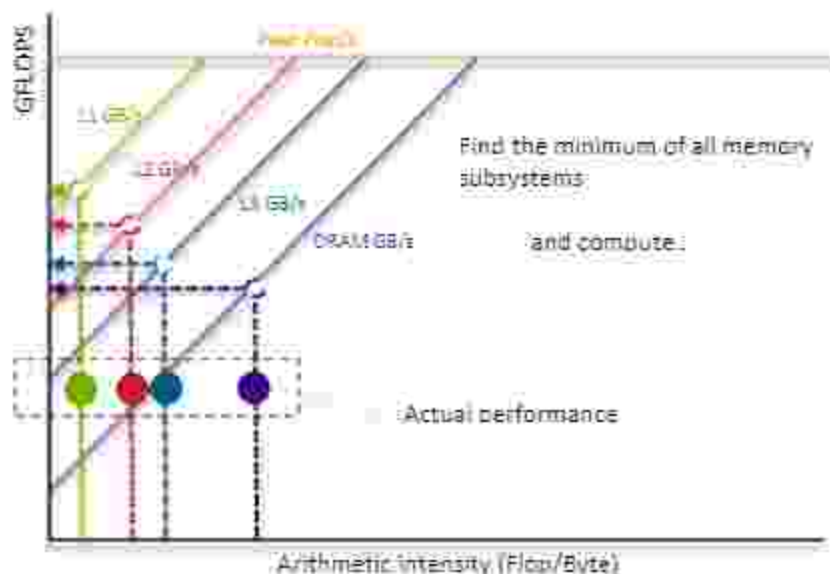
AUTOMATED ROOFLINE ANALYSIS

Via Intel Advisor

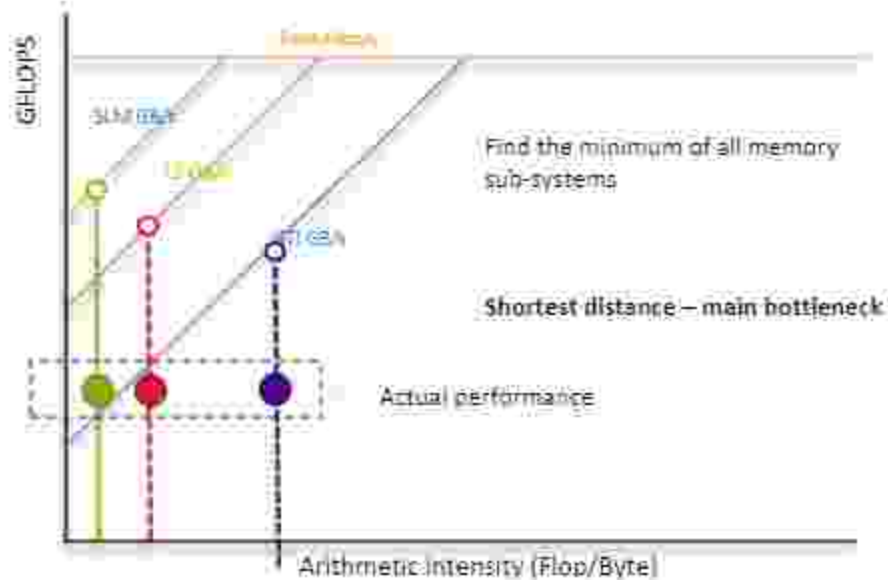
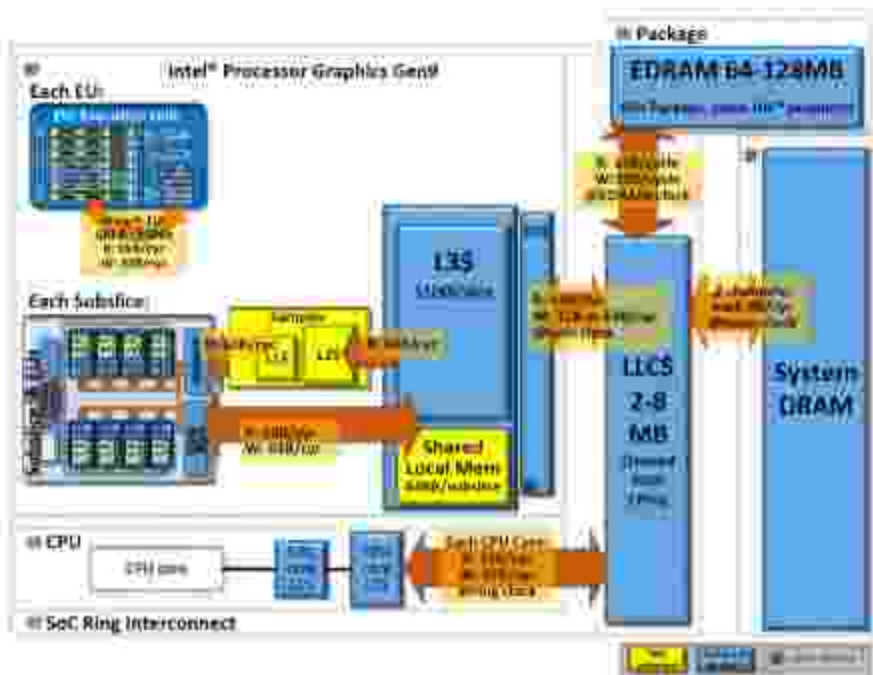


- Intel Advisor is a part of the Intel oneAPI Base Toolkit
- Performance headroom against hardware limitations
- Identify which optimizations will payoff the most
- Detect and prioritize bottlenecks
- Visualize optimization progress
- GPU roofline analysis
 - GPU recommendations
 - See the CPU and GPU performance side-by-side
 - GPU application performance characterization

- Memory Level Roofline (MLR) on CPU

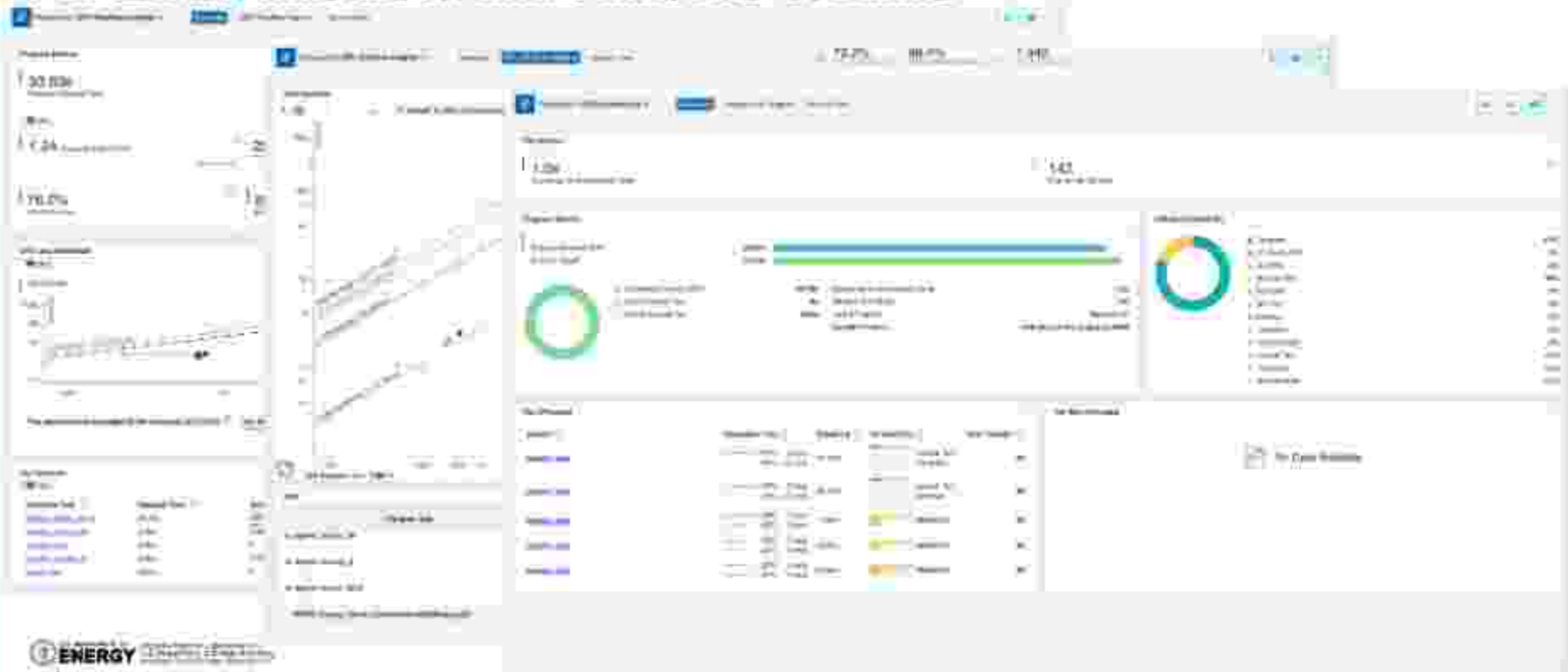


MEMORY LEVEL ROOFLINES (MLR) ON GPU



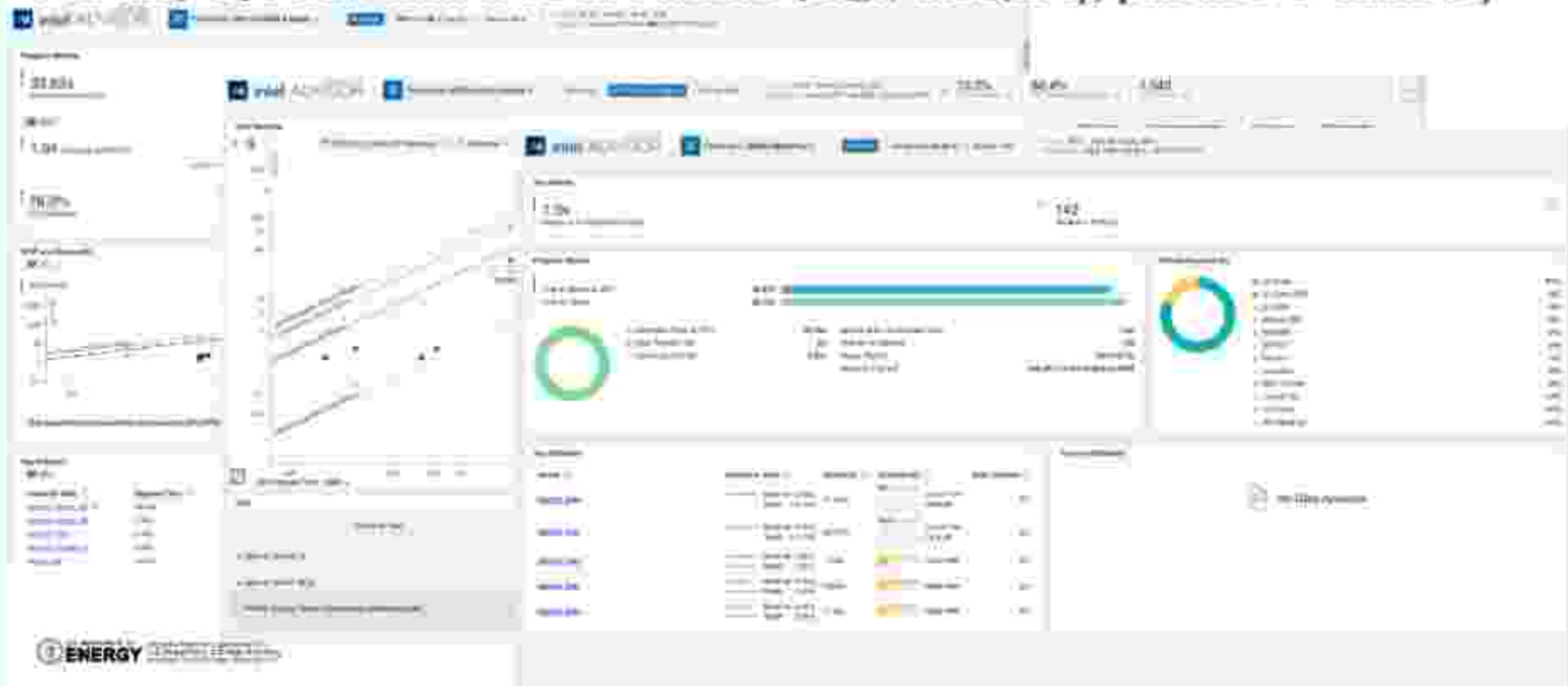
ADVISOR GUI INTERFACE

For Linux, Mac(x86), and Windows systems



“FAT” STAND-ALONE HTML REPORT

For any system with web browsers (e.g., Mac(M1), phones & tablets)



ARGONNE ADVISOR ROOFLINE USE-CASES

ARGONNE ADVISOR ROOFLINE USE-CASES

Presented at roofline performance analysis tutorials (@SC, ECPAM)

Application	ANL PoC	Programming models	Field of Science
AMR-Wind	Jaehyuk Kwack	SYCL/DPG++	Wind farm simulation (CFD/ES)
GAMESS/RI-MP2	Colleen Bertoni	OpenMP Target	Electronic structure method
NekBench	Kris Rowe	DGCA with OpenCL backend	Computational Fluid Dynamics
TestSNAP	Yasaman Ghadar	OpenMP Target	Material Science / Chemistry
LAMMPS kernel	Chris Knight	OpenMP Target, OpenCL, SYCL/DPG++	Material Science / Chemistry
XSbench/RSBench	John Tramm	OpenMP Target	Monte Carlo neutron transport
SW4Lite	Brian Hornarding	RAJA with SYCL/DPG++ backend	Seismic modeling

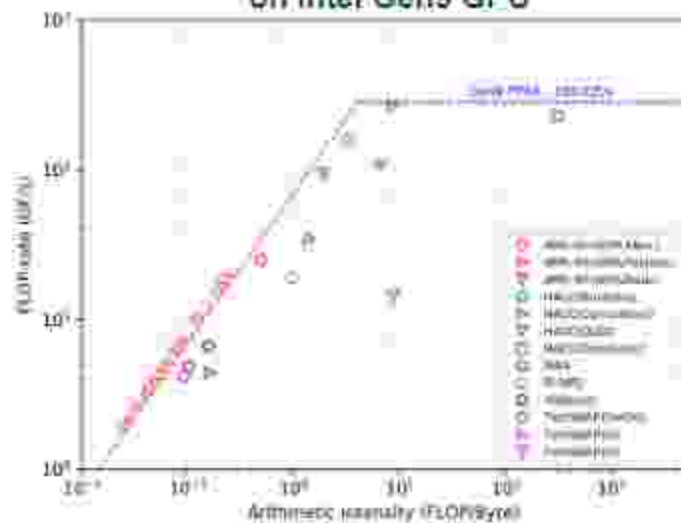
SC21 P3HPC WORKSHOP (2021)

Performance, Portability & Productivity in HPC @ SC21 on November 14th, 2021

- Title: Evaluation of Performance Portable Programming Models and Framework across AMD, Intel and NVIDIA GPUs using Applications and Mini-Apps
- Authors: J. Kwack, J. Tramm, C. Bertoni, Y. Ghadar, B. Hornerding, E. Rangel, C. Knight, and S. Parker
- Employed ECP applications and mini-apps

Application	Kernel	Application	Kernel
AMR-Wind	MLABec	SW4	curvilinearitsg
	MLPoisson	GAMNESS RI-MP2	RIMP2@omp
	MLNode	XSBench	XSBench\$omp
HACC	BarExtras	TestSNAP	FusedDeIDr
	Corrections		U1
	DuDt		Y1
	Geometry		

EDRAM-based Advisor roofline results on Intel Gen9 GPU



PORTING

GAMESS RI-MP2 MINI-APP

FROM INTEL CPU TO INTEL GPU

GAMESS AND RI-MP2 MINI-APP

- GAMESS is the General Atomic and Molecular Electronic Structure System
 - General-purpose electronic structure code (many methods and capabilities)
 - ~1 million lines of Fortran
 - Optional C/C++ GPU-accelerated libraries/applications in GAMESS
- Scientific problem of interest
 - FMO/RI-MP2 calculations towards accurate simulations of catalysis reactions inside a mesoporous silica nanoparticle
- RI-MP2 mini-app
 - Computes RI-MP2 perturbative correction to the Hartree-Fock energy
 - RI-MP2 mini-app examples

Coronene



H₂O clusters



Fullerene



INTRO OF RI-MP2 MINI-APP

An ECP Proxy Application

- Computes RI-MP2 perturbative correction to the Hartree-Fock energy

$$E^{(2)} = \sum_{i \leq j} (2 - \delta_{ij}) \sum_{ab}^{vir} \frac{(ia|jb)[2(ja|jb) - (ib|ja)]}{\epsilon_i + \epsilon_j - (\epsilon_a + \epsilon_b)}$$

- Simplifies expression with the RI approximation letting it be written in terms of matrix multiplication

$$(ia|jb) = \sum_{\mu} B_{ia\mu}^i B_{b\mu}^j$$

- Combine i,j, loops over occ orbitals, and then reduce over a,b, virtual orbitals to compute the final correlation energy

```
for(int JACT=0; JACT<NACT; JACT++){
  for(int IACT=0; IACT<=JACT; IACT++){
    // A kernel to compute QVW via Matrix Multiplication
    QVW[0:NVIR][0:NVIR] = B32[IACT][0:NVIR][0:NAUXBASD] *
                        B32[JACT][0:NAUXBASD][0:NVIR];

    // Accumulating E2 using QVW[:::][i1:::][i2:::][i3:::] and eab[:::][i1][i2]
    for(int IB=0; IB<NVIR; IB++){
      for(int IA=0; IA<NVIR; IA++){
        T2jab = QVW[IB+NVIR][IA] /
                ( e13[IACT][IACT] - eab[IB][IA] );
        Qc = QVW[IB][IA] + QVW[IB][IA];
        E2_c += T2jab * (Q_c - QVW[IA][IB]);
      }
    }
    E2 += E2_c * E2C;
  }
}
```

RI-MP2 CODE VERSIONS

- We constructed four RI-MP2 variants to explore with w25.rand.
- These variants add progressive/incremental levels of optimization as follow:
 - **V0-CPU:** an initial version,
 - **V1-CPU:** QVV computations use **MKL DGEMM**,
 - **V3-GPU:** **offloaded V1-CPU to GPU using OpenMP Target offloading** model with optimal values for **num_teams**, and **threads_limit**
 - **V5-GPU:** QVV computation uses **less MKL DGEMM calls by restructuring an outer loop**; For E2, **IAC T loop** is distributed **subslices w/ "target teams distribute"**, and then **IB loop** uses **EUs w/ "parallel for"** at each subslice, with optimal values for **num_teams**, and **threads_limit**

V0-CPU

- QVV: Computing QVV
- E2:E2 accumulation

```
double *QVV;  
double E2_local=0.0E0;  
QVV = new double[NVIR*NVIR];  
  
for(int IACT=0; IACT<NACT; IACT++){  
    for(int IACT=0; IACT==IACT; IACT++){
```

QVV

```
// Compute QVV  
for(int j = 0; j < NVIR; j++)  
    for(int i = 0; i < NVIR; i++)  
        QVV[i,j] += B[i,IACT]*B[i,IACT];  
}
```

E2

```
// Accumulate E2  
double E2 = 0.0;  
for(int IBA=0; IBA<NVIR; IBA++){  
    for(int IBB=0; IBB<NVIR; IBB++){  
        double E2ab = QVV[IBA, IBB] * B[IBA, IACT] * B[IBB, IACT];  
        double Q_c = 2*QVV[IBA, IBB];  
        E2_c += E2ab * Q_c - QVV[IBA, IBB];  
    } // loop for IBA and IBB  
    double FAC = IACT==IACT ? 1.0E0 : 2.0E0;  
    E2_local += FAC*E2_c;  
}
```

```
E2 = *E2 + E2_local;  
delete[] QVV;
```

- QVV L2 cache bound
 - 31.33 GF/s
 - L2 A= 0.5



Data traffic (Bytes) ratios:

L2 L3: DRAM = 4.2 2.2 1

- E2 L2 cache bound
 - 5.13 GF/s
 - L2 A= 0.091



Data traffic (Bytes) ratios:

L2 L3: DRAM = 4.3 1.4 1

V1-CPU

- **QVV: Computing QVV using MKL DGEMM**
- **E2:E2 accumulation**

```
double *QV;
double E2_local=0.0;
QV = new double[NVIR*NVIR];

for(int JACT=0; JACT<NACT; JACT++)
for(int IACT=0; IACT<NACT; IACT++)

// Compute QV
int n=NACT;
int k=NACT*NVIR;
double msa = 1.0;
double mso = 0.0;
#pragma omp for collapse(2)
for(int JACT=0; JACT<NACT; JACT++)
for(int IACT=0; IACT<NACT; IACT++)
    QV[IACT*NVIR+JACT] = msa;
// End of QV
```

```
E2
// Accumulate E2
double E2 = 0.0;
for(int IB=0; IB<NVIR; IB++)
for(int IA=0; IA<NVIR; IA++)
    double sumab = QV[IB*NVIR+IA] * E2;
    double sumc = QV[IB*NVIR+IA];
    E2 = sumab + sumc;
// Loop for IA and IB
double E2 = E2 + E2_local;
E2_local = E2;
// Loop for IACT and JACT

*E2 = *E2 + E2_local;
delete[] QV;
```

- QVV DP FMA bound
 - 61.96 GF/s
 - L2 A/= 1.34

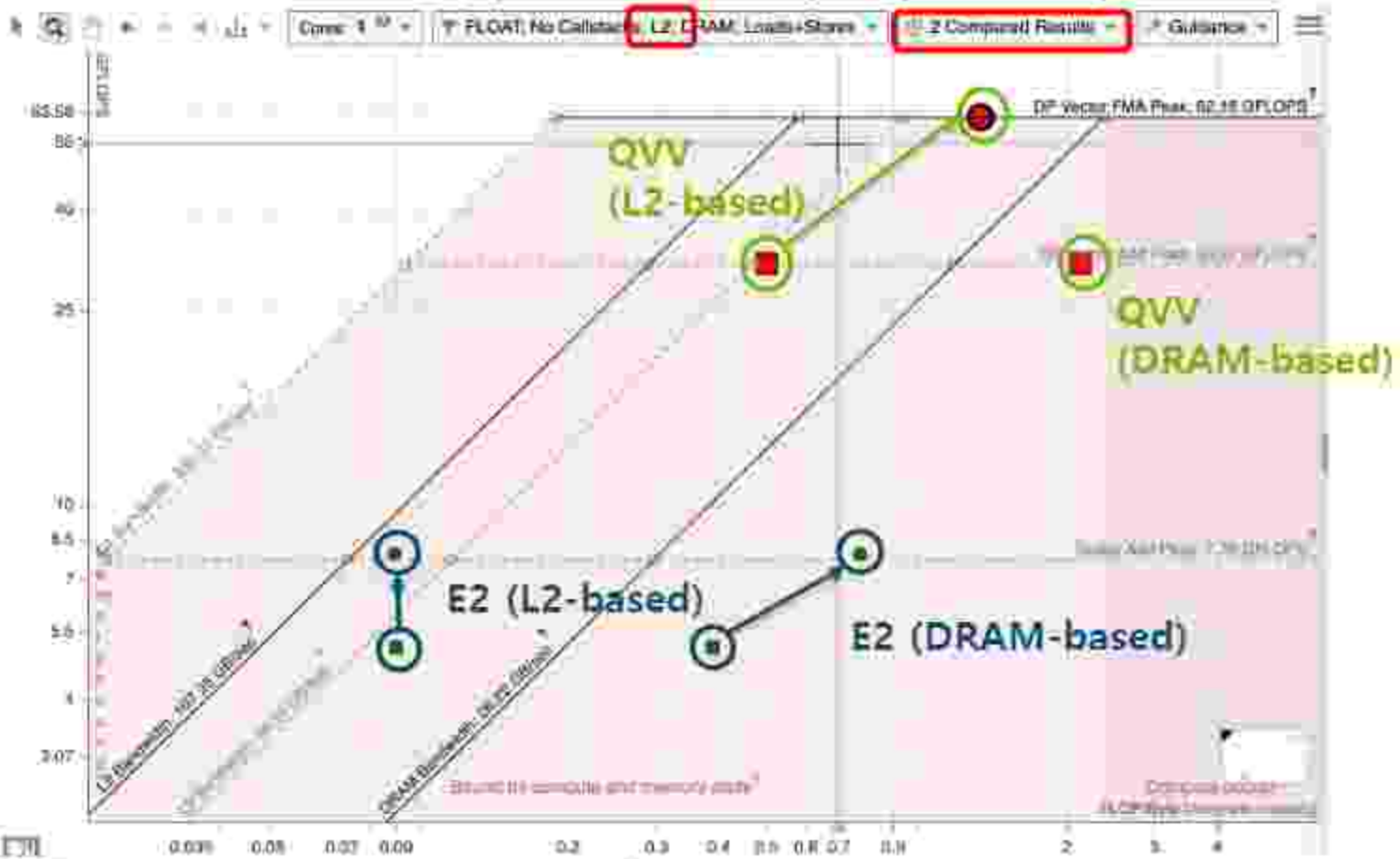


- E2 L2 cache bound
 - 7.95 GF/s
 - L2 A/= 0.090



V0-CPU to V1-CPU

- QVV computation is compute-bound (by MKL) w/ much lower memory pressure; it consequently increases performance of E2.



V3-GPU: initial CPU-like version

```
double *QVV;
double E2_local=0.0;
int dimm=0;
QVV = new double(NVIR*NVIR);
double *B32I; *B32J;

#pragma omp target enter data map(alloc:QVV(0:NVIR*NVIR))
device(dimm)
#pragma omp target enter data
map(to: a1(0:NACT*NACT), esb(0:NVIR*NVIR), B32(0:B32size))
device(dimm)
for(int JACT=0; JACT<NACT; JACT++)
for(int IACT=0; IACT<NACT; IACT++){

    // Compute QVV
    int n=NVIR;
    int k=NUMTHREAD;
    double rrr = 1.0;
    double drr = 0.0;
    B32I = B32I+IACT,0,0;
    B32J = B32J+IACT,0,0;
    #pragma omp target variant dispatch
    use_device_ptr(B32I, B32J, B32, QVV); device(dimm)
    dgemv('T', n, n, rrr, drr, B32I, 0, B32J, 0, &QVV[IACT], n);

    // Accumulate E2
    double E2 = 0.0;
    #pragma omp target teams distribute parallel for
    reduction(+:E2) rrr teams(80) thread_limit(72) device(dimm)
}
```

```
double E2=0; E2<NVIR; E2++{
for(int I=0; I<NVIR; I++){
double rrrab = QVV[I*IACT+IACT]-esb[I*IACT];
double Q = 0+QVV[I*IACT];
E2 += rrrab * Q; Q = QVV[IACT*I];
} // loop for IA and IB
} // omp target teams distribute parallel for
double E2 = E2+Q; E2 = E2+Q; E2 = E2+Q;
E2_local += E2*E2;
} // loop for IACT and JACT
```

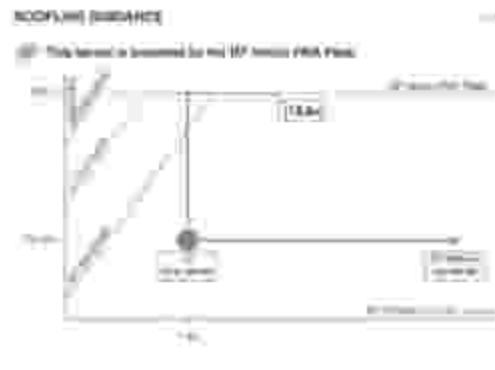
```
E2 = E2 + E2_local;
delete[] QVV;
```

```
#pragma omp target exit data map(release:QVV(0:NVIR*NVIR))
device(dimm)
#pragma omp target exit data
map(release: a1(0:NACT*NACT), esb(0:NVIR*NVIR), B32(0:B32size))
device(dimm)
```

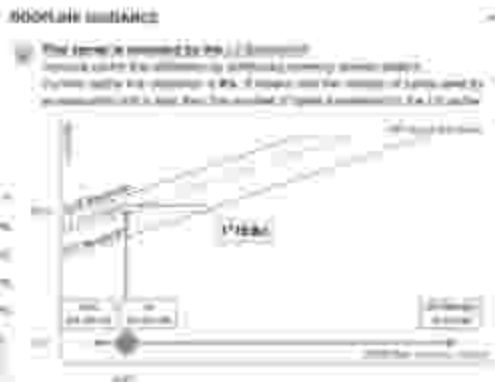
- Computing QVV using MKL DGEMM w/ GPU offloading
- E2 accumulation w/ GPU offloading

V3-GPU: initial CPU-like version

- QV: DP FMA bound
 - 127.25 GF/s
 - L3 AI= 1.14



- E2: L3 GPU cache bound
 - 1.41 GF/s
 - L3 AI= 0.071



V5-GPU: restructured loops for better GPU performance

```
double *QVV;
double E2_local=0.0E0;
int dmm=0;
QVV = new double[NVIR*NACT*NVIR];
double *B32;

#pragma omp target enter data map(alloc:QVV(0:NVIR*NACT*NVIR))
device(dmm)
#pragma omp target enter data
map(to:es1(0:NACT*NACT), esb(0:NVIR*NVIR), B32(0:BS2size))
device(dmm)
for(int IACT=0; IACT<NACT; IACT++){

    // Compute QVV
    int es=NVIR*(IACT+1);
    int es=NVIR;
    int es=NVIR;
    double one = 1.0;
    double zero = 0.0;
    B32 = BS2(IACT,0,0);
    #pragma omp target variant dispatch use_device_ptr(B32,B32,QVV)
    device(dmm)
    dgemv('T', 'N', es, es, es, one, B32, es, B32, es, zero, QVV, es);

    // Accumulate E2
    #pragma omp target teams distribute reduction(* E2_local)
    teams(10) thread_limit(7) hws(1) from/
    for(int IACT= IACT; IACT<NACT; IACT++){
        double E2_c=0.0;
        #pragma omp parallel for reduction(+:E2_c)

```

```
for(int IB=0; IB<NVIR; IB++)
for(int IS=0; IS<NVIR; IS++)
    double E2ab = QVV[IB,IS,IACT] + (es-IB)*E2_c;
es=IB,IS;
    double E2_c = 0*QVV[IB,IS,IACT];
    E2_c += E2ab * (E2_c = QVV[IA,IACT,IB]);
} // loop for IB and IS
double EAC = (EACT-EACT) * (1.0E0) + (E2E1);
E2_local += E2_c*E2_c;
} // loop for IACT
} // loop for IACT
```

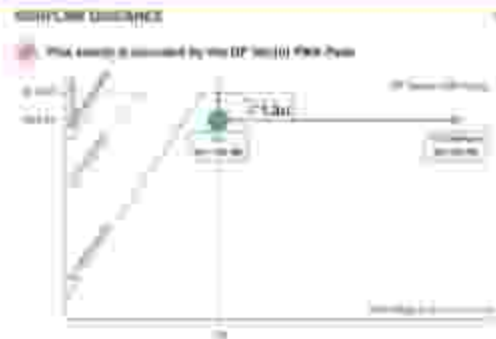
```
*E2 = *E2 + E2_local;
delete[] QVV;
```

```
#pragma omp target exit data map(release:QVV(0:NVIR*NACT*NVIR))
device(dmm)
#pragma omp target exit data
map(release:es1(0:NACT*NACT), esb(0:NVIR*NVIR), B32(0:BS2size))
device(dmm)
```

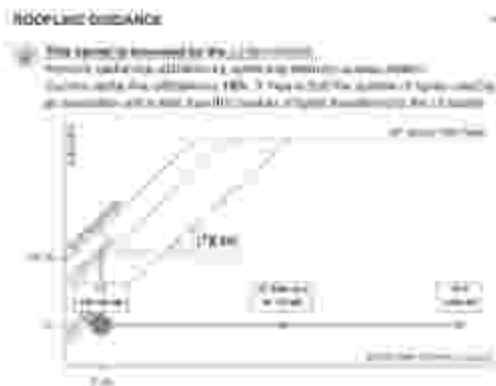
- Computing QVV with less MKL DGEMM offloading calls on GPU by restructuring an outer loop
- E2 accumulation w/ hierarchial GPU offloading (IACT-loop to Subslices (SS), and IB-loop to Execution Units (EUs))

V5-GPU: restructured loops for better GPU performance

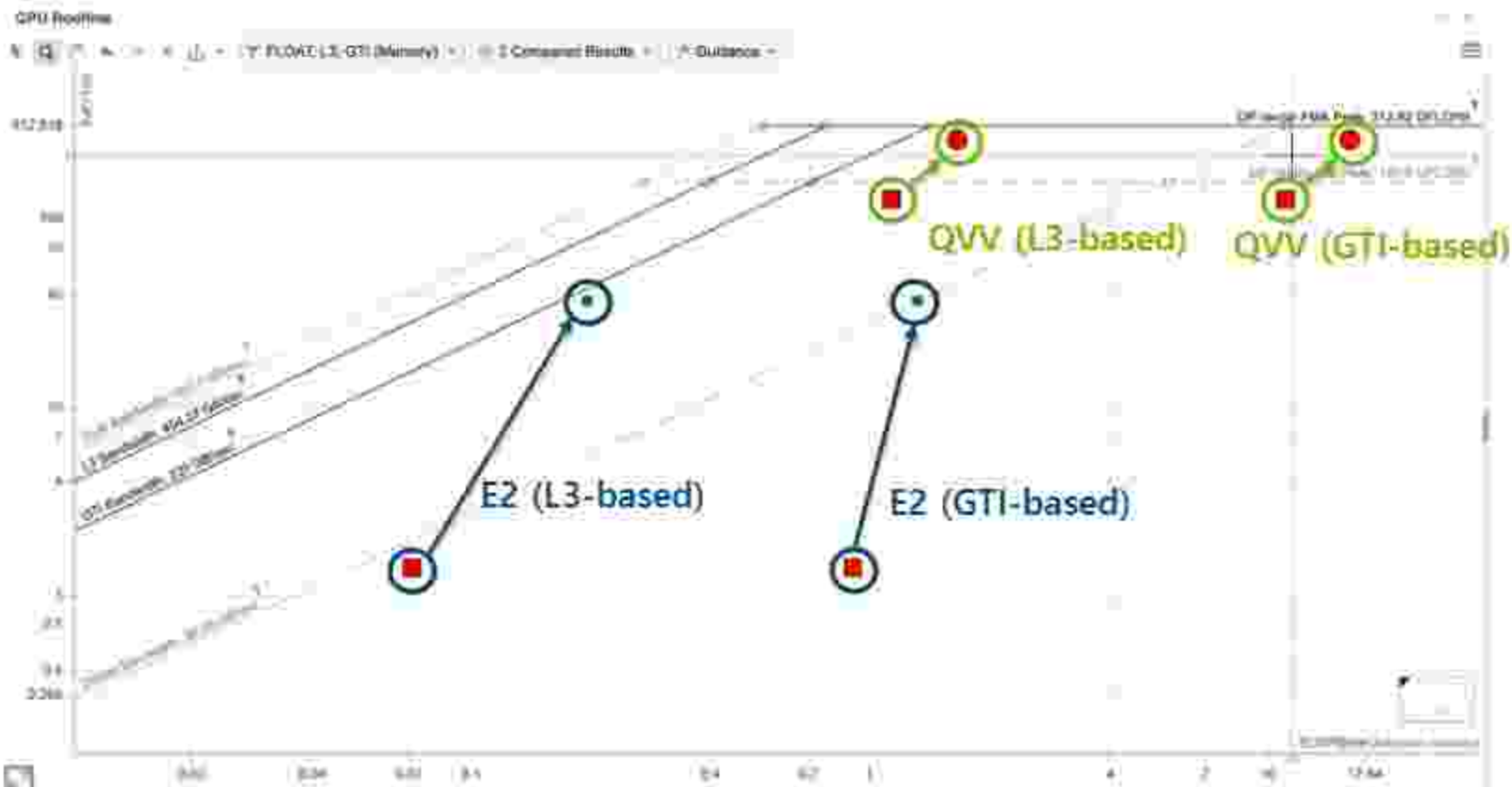
- QVV: DP FMA bound
 - 262.84 GF/s
 - L3 AI= 1.66



- E2: L3 GPU cache bound
 - 37.0 GF/s
 - L3 AI= .019



Comparison from V3-GPU to V5-GPU



SUMMARY & CONCLUDING REMARKS

PERFORMANCE RESULTS

Tested Input: w25 rand (random data with structure of 25 H₂O clusters)

Employed Compute Node Intel(R) Xeon(R) CPU E3-1585 v5

- CPU: Intel Xeon Skylake @ 3.5GHz (4 cores)
- GPU: Intel Gen9 GT4e @ 1.15GHz (72 EUs)

H₂O clusters



Selected versions	QW		E2		Overall	
	Target	Time(s)	Target	Time(s)	Time(s)	Speedup over V00
V0-CPU	CPU	158.68	CPU	1.34	158.02	1.0x
V1-CPU	CPU	86.41	CPU	0.92	87.33	1.8x
V3-GPU	GPU	32.38	GPU	23.06	55.44	2.9x
V5-GPU	GPU	21.49	GPU	0.86	22.35	7.1x

CONCLUDING REMARKS

- Argonne application developers and other US DOE collaborators have been actively porting their applications to Aurora testbed systems with Intel oneAPI toolkits for the coming Aurora Exa-scale supercomputer at Argonne National Laboratory.
- As an example, we present four versions of GAMESS RI-MP2 mini-app on Intel CPU and GPUs.
- Intel Advisor successfully provides the detailed performance data via roofline analysis features on both Intel CPUs and GPUs.
- Call-to-actions
 - If you are working on ECP projects and interested in Aurora testbeds with Intel oneAPI toolkit, you can get the access via <https://www.jlsc.anl.gov/accessing-jlsc-resources/>.
 - If you are interested in Intel GPUs with Intel oneAPI toolkits for your own applications, you may try the Intel DevCloud system via the following link
 - <https://www.intel.com/content/www/us/en/developer/tools/devcloud/overview.html>
 - Enjoy!

ACKNOWLEDGEMENT

- This work was supported by
 - the Argonne Leadership Computing Facility, which is a DOE Office of Science User Facility supported under Contract DE-AC02-06CH11357,
 - and by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of two U.S. Department of Energy organizations (Office of Science and the National Nuclear Security Administration).
- We also gratefully acknowledge the computing resources provided and operated by the Joint Laboratory for System Evaluation (JLSE) at Argonne National Laboratory.
- We would like to thank **Colleen Bertoni** at Argonne National Laboratory for discussing GAMESS RI-MP2 miniapp, **Kirill Rogozhin**, **Yulia Efimenko**, and **Ekaterina Guseva** for **Advisor** and **Louise Huot** for **MKL** for providing technical supports for issues, reviewing data and providing feedback for this study.

THANKS!