**ALCF DEVELOPER SESSION**

# INTEL PERFORMANCE PROFILING TOOLS
# ON AURORA

**JAEHYUK KWACK**
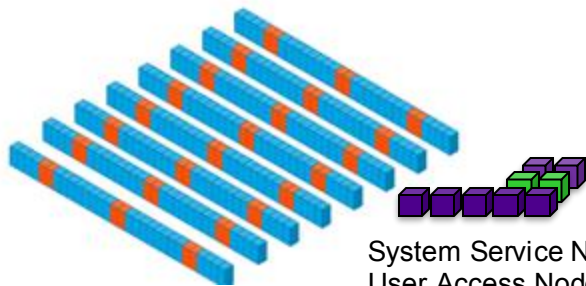( jkwack@anl.gov )
ALCF Perf. Engr. Group

9/25/2024

# AURORA OVERVIEW

# AURORA HIGH-LEVEL SYSTEM OVERVIEW



**COMPUTE RACK**
64 Compute blades
32 Switch blades
GPU: 49.1 TB HBM
CPU: 8.2 TB HBM, 64 TB DDR5
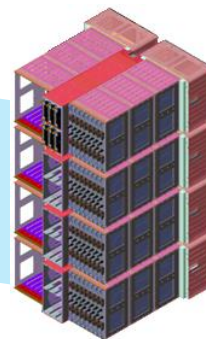
System Service Nodes (SSNs)
User Access Nodes (UANs)
DAOS Nodes (DNs)
Gateway Nodes (GNs)
  IOF service, scalable library loading
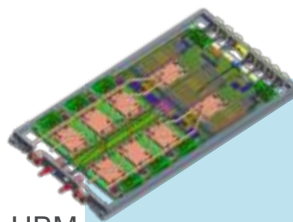  DAOS <-> Lustre data mover

**AURORA SYSTEM**
166 Compute racks
10,624 Nodes
GPU: 8.16 PB HBM
CPU: 1.36 PB HBM, 10.9 PB DDR5
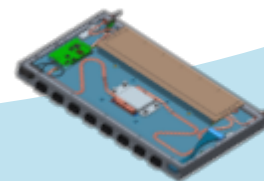DAOS: 64 racks, 1024 nodes
      230 PB (usable), 31 TB/s

**COMPUTE BLADE**
2x Intel Xeon Max Series w HBM
6x Intel Data Center GPU Max Series
GPU: 768 GB HBM
CPU: 128 GB HBM, 1024 GB DDR5

**SWITCH BLADE**
1 Slingshot switch
64 ports
Dragonfly topology

# AURORA EXASCALE COMPUTE BLADE

## NODE CHARACTERISTICS

| | |
|---|---|
| 6 | GPU - Intel Data Center GPU Max Series (#) |
| 2 | CPU - Intel Xeon CPU Max Series (#) |
| 768 | GPU HBM Memory (GB) |
| 19.66 | Peak GPU HBM BW (TB/s) |
| 128 | CPU HBM Memory (GB) |
| 2.87 | Peak CPU HBM BW (TB/s) |
| 1024 | CPU DDR5 Memory (GB) |
| 0.56 | Peak CPU DDR5 BW (TB/s) |
| ≧ 130 | Peak Node DP FLOPS (TF) |
| 200 | Max Fabric Injection (GB/s) |
| 8 | NICs (#) |

Argonne
NATIONAL LABORATORY

# PROFILING TOOLS ON AURORA

Argonne
NATIONAL LABORATORY

# PROFILING TOOLS ON AURORA
## A list of popular profiling tools on Aurora

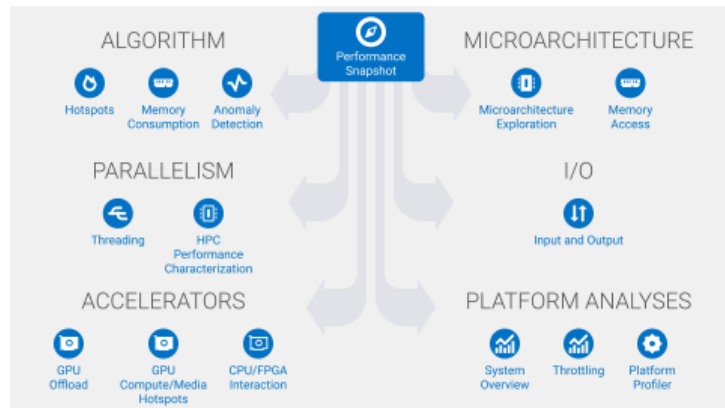- Intel VTune

- Intel Advisor

- Intel APS (Application Performance Snapshot)

- Intel xpu-smi

- Intel unitrace


- THAPI from Argonne (https://github.com/argonne-lcf/THAPI)

- HPCToolkit from Rice University (https://hpctoolkit.org/)

- TAU from University of Oregon (https://www.paratools.com/tau)

# INTEL VTUNE
## Profile GPU Performance



- Multi-GPU systems analysis

- GPU Offload cost profiling
  - CPU vs GPU boundness
  - Offload overhead & host-to-device traffic, GPU compute vs data transfer
  - GPU utilization and software queues

- GPU Hotspots analysis
  - XVE (Xe Vector Engine) and memory efficiency metrics, GPU occupancy limiting factors
  - Memory hierarchy diagram and throughput analysis

- Source level in-kernel profiling (need to build with "*-fdebug-info-for-profiling -gline-tables-only*")
  - Dynamic instruction count
  - Basic Block execution latency
  - Memory latency
  - HW-assisted stall sampling
  - Memory Access Analysis Tool (MAAT)

Argonne
NATIONAL LABORATORY

# INTEL VTUNE

## Analysis types for Intel GPUs

■ HPC Performance Characterization
- $ <mpi launcher> <mpi param-s> vtune –collect hpc-performance -r <result_dir> <my app> <app param-s>
- Provides a different aspect of application performance
- High level hardware information, CPU cores utilization, GPU stacks utilization including XVE HW metrics and top offload regions, CPU-side memory metrics, and CPU instruction statics

■ GPU offload
- $ <mpi launcher> <mpi param-s> vtune –collect gpu-offload -r <result_dir> <my app> <app param-s>
- Serves studies of an application offload implementation and assesses its efficiency
- Traces Level-zero and OpenCL API functions in oneAPI software stack; detects long latency host functions; shows time spent in data allocation and transfer function as well as kernel device time

# INTEL VTUNE

## Analysis types for Intel GPUs

- GPU Compute/Media Hotspots
  - `$ <mpi launcher> <mpi param-s> vtune –collect gpu-hotspots –r <result_dir> <my app> <app param-s>`
  - The most accurate analysis in tracing kernels on GPU
  - Allows to analyze the most time-consuming GPU kernels, characterize GPU usage based on GPU hardware metrics, identify performance issues caused by memory latency or inefficient kernel algorithms, and analyze GPU instruction frequency per certain instruction types.

# VTune server for pre-collected results on Aurora via SSH terminal

▪ Step1: Open a new terminal and log into Aurora login node (no X11 forwarding required)

```
$ ssh <username>@bastion.alcf.anl.gov
$ ssh <username>@login.aurora.alcf.anl.gov
```

▪ Step2: Start VTune server on an Aurora login node after loading oneapi module and setting corresponding environmental variables for VTune

```
$ module load oneapi
$ vtune-backend --data-directory=<location of precollected VTune results>
```

▪ Step3: Open a new terminal with SSH port forwarding enabled (need 2 hops):

```
$ ssh -L 127.0.0.1:<port printed by vtune-backend>:127.0.0.1:<port printed by vtune-backend>
<username>@bastion.alcf.anl.gov
$ ssh -L 127.0.0.1:<port printed by vtune-backend>:127.0.0.1:<port printed by vtune-backend>
<username>@login.aurora.alcf.anl.gov
```

▪ Step4: Open the URL printed by VTune server in your local web browser (e.g., firefox on your laptop)

# INTEL APS
## VTune Profiler's Application Performance Snapshot

- Provides an aggregated view of an application at scale, designed for large MPI workloads.

- Captures performance aspects of compute intensive applications
  - MPI and OpenMP usage and imbalance,
  - CPU and GPU utilization
  - CPU stalls due to memory accesses, vectorization, I/O, and memory footprint

- Command lines
  - `$ <mpi launcher> <mpi param-s> aps -r <result_dir> <my app> [<app param-s>]`
  - `$ aps-report <result_dir>`

# INTEL ADVISOR
## Overview

- A design and analysis tool for developing performant code
  - Performant CPU Code: Design your application for efficient threading, vectorization, and memory use.
  - Efficient GPU Offload: Identify parts of the code that can be profitably offloaded. Optimize the code for compute and memory.
  - Flow Graph Design and Analysis: Create, visualize, and analyze task and dependency computation for heterogeneous algorithms.

Summary of CPU/GPU analysis





GPU Roofline
w/ memory hierarchy

Argonne
NATIONAL LABORATORY

# INTEL ADVISOR

## Advisor Roofline analysis

- Advisor version on Aurora
  ```
  $ advisor --version
  Intel(R) Advisor 2024.2.1 (build 615624) Command Line Tool
  Copyright (C) 2009-2024 Intel Corporation. All rights reserved.
  ```

- Step1: Setting the environments
  ```
  $ module load oneapi
  $ export PRJ=<your_project_dir>
  ```

- Step 2-a: Collecting the GPU Roofline data on a single GPU (Survey analysis and Trip Count with FLOP analysis)
  ```
  $ advisor --collect=roofline --profile-gpu --project-dir=$PRJ -- <your_executable> <your_arguments>
  ```

- Step 2-b: Collecting the GPU Roofline data on one of MPI ranks(Survey analysis and Trip Count with FLOP analysis)
  ```
  $ mpiexec -n 1 gpu_tile_compact.sh advisor --collect=survey --profile-gpu --project-dir=$PRJ --
  <your_executable> <your_arguments> : -n 1 gpu_tile_compact.sh <your_executable> <your_arguments>
  $ mpirun -n 1 gpu_tile_compact.sh advisor --collect=tripcounts --profile-gpu --flop --no-trip-counts --
  project-dir=$PRJ -- <your_executable> <your_arguments> : -n 1 gpu_tile_compact.sh <your_executable>
  <your_arguments>
  ```

- Step 3-a: Generate a GPU Roofline report, and then review the HTML report
  ```
  $ advisor --report=all --project-dir=$PRJ --report-output=${PRJ}/roofline_all.html
  ```

- Step 3-b: Download the project folder to your laptop and open it with Advisor Client
  - https://www.intel.com/content/www/us/en/developer/articles/tool/oneapi-standalone-components.html#advisor

# THAPI: TRACING HETEROGENEOUS APIs

## A lightweight tool for tracing and sampling

- Overview
  - THAPI is a portable, programming model-centric tracing framework for heterogeneous systems.
    - OpenCL, L0, Cuda, HIP, OMPT, MPI
  - Two Components:
    - Tracing Events
      - LTTng based tracing
    - Parsing of the trace
      - Babeltrace2 library and tools

- Device Sampling
  - Ability to sample device telemetry with API tracing
    - Holistic view of system performance and behavior
    - Help understand H/W behavior in application context
    - Power/energy optimization
    - Resource management
    - Improves Debugging
  - Power, Fabric and Memory Traffic, Engine Activities
  - Timeline Visualization
    - Perfetto

# ADVISOR & VTUNE EXAMPLES ON AURORA

# WARMING-UP

Advisor on Aurora

- $ advisor -version
- $ advisor -help
- $ advisor -help collect

VTune on Aurora

- $ vtune -version
- $ vtune -help
- $ vtune -help collect
- $ vtune -help collect gpu-hotspots
- Vtune server for post-processing via ssh forwarding (i.e., $ vtune-backend)

APS on Aurora

- $ aps –version
- $ aps -help

# VTUNE HPC-PERFORMANCE ANALYSIS

```
$ mpiexec -n 12 -ppn 12 gpu_tile_compact.sh vtune -collect hpc-performance -r VTune_hpc-performance
./Comp_GeoSeries_omp_mpicxx_DP 2048 1000
```

# VTUNE GPU-OFFLOAD ANALYSIS

```
$ mpiexec -n 12 -ppn 12 gpu_tile_compact.sh vtune -collect gpu-offload -r VTune_gpu-offload
./Comp_GeoSeries_omp_mpicxx_DP 2048 1000
```

# VTUNE GPU-HOTSPOTS ANALYSIS

```
$ mpiexec -n 12 -ppn 12 gpu_tile_compact.sh vtune -collect gpu-hotspots -r VTune_gpu-hotspots
./Comp_GeoSeries_omp_mpicxx_DP 2048 1000
```

# VTUNE INSTRUCTION COUNT ANALYSIS

```
$ mpiexec -n 12 -ppn 12 gpu_tile_compact.sh vtune -collect gpu-hotspots -knob characterization-mode=instruction-count -r VTune_inst-count ./Comp_GeoSeries_omp_mpicxx_DP 2048 1000
```

# VTUNE SOURCE ANALYSIS

```
$ mpiexec -n 12 -ppn 12 gpu_tile_compact.sh vtune -collect gpu-hotspots -knob profiling-mode=source-analysis -r
VTune_source ./Comp_GeoSeries_omp_mpicxx_DP 2048 1000
```

# VTUNE MEMORY LATENCY ANALYSIS

```
$ mpiexec -n 12 -ppn 12 gpu_tile_compact.sh vtune -collect gpu-hotspots -knob profiling-mode=source-analysis -
knob source-analysis=mem-latency -r VTune_mem-latency ./Comp_GeoSeries_omp_mpicxx_DP 2048 1000
```

# APS ANALYSIS

```
$ mpiexec -n 192 gpu_tile_compact.sh aps -r APS_NMPI192 ./amr_wind
../test/test_files/abl_godunov/abl_godunov.inp
```

```
$ aps-report --metrics=?
APS_NMPI192_MaxGrid256_MaxStep13_jobid810828

| Available Metrics:

|--------------------

Elapsed Time

MPI Time

MPI Time

MPI Imbalance

MPI Hotspot 1 - MPI_Waitall

MPI Hotspot 1 - MPI_Waitall

MPI Hotspot 2 - MPI_Allreduce

MPI Hotspot 2 - MPI_Allreduce

MPI Hotspot 3 - MPI_Init

MPI Hotspot 3 - MPI_Init

MPI Hotspot 4 - MPI_Isend

MPI Hotspot 4 - MPI_Isend

MPI Hotspot 5 - MPI_Comm_create

MPI Hotspot 5 - MPI_Comm_create

Disk I/O Bound

Disk I/O Bound

Disk read

Disk write

Resident Memory Usage per Rank

Resident Memory Usage per Node
```

```
Virtual Memory Usage per Rank

Virtual Memory Usage per Node

Instructions Per Cycle Rate

Average CPU Frequency

Physical Core Utilization

Average Physical Core Utilization

Memory Stalls

Cache Stalls

DRAM Stalls

Average DRAM Bandwidth

DRAM Bandwidth Peak

DRAM Bandwidth Average

DRAM Bandwidth Bound

SP GFLOPS

DP GFLOPS

Vectorization

SP FLOPs

DP FLOPs

Non-FP

FP Arith/Mem Rd Instr. Ratio

FP Arith/Mem Wr Instr. Ratio

GPU Accumulated Time

GPU Stack Utilization
```

```
XVE State: Active

XVE State: Stalled

XVE State: Idle

GPU Occupancy

GPU Inbound PCIe Read

GPU Inbound PCIe Write

GPU Outbound PCIe Read

GPU Outbound PCIe Write

Network Controller Inbound PCIe Read

Network Controller Inbound PCIe Write

Network Controller Outbound PCIe Read

Network Controller Outbound PCIe Write

Inbound PCIe Read Per Device

Inbound PCIe Write Per Device

Outbound PCIe Read Per Device

Outbound PCIe Write Per Device

GPU Accumulated Time Per Device

GPU Stack Utilization Per Device

XVE State: Active Per Device

XVE State: Stalled Per Device

XVE State: Idle Per Device

GPU Occupancy Per Device
```

Argonne NATIONAL LABORATORY

```
jkwack@aurora-uan-0009:/lus/flare/projects/Aurora_deployment/jkwack/ExaWind_flare/amr-wind/JK_build_2024.07.30.002_w_mpi> aps-report
APS_NMPI192_MaxGrid256_MaxStep13_jobid810828 --metrics="MPI Time, GPU Occupancy"

| Metric Table%

|------------------------------------------------------------------------------------------

Metric Name                          Node Name        Rank    Metric Value    Outlier Type

MPI Time, s                          x4305c3s2b0n0     135       30.0779          None

MPI Time, s                          x4305c3s3b0n0      72       30.037           None

MPI Time, s                          x4305c4s3b0n0      79       30.0068          None

MPI Time, s                          x4305c3s2b0n0      71       30.004           None

MPI Time, s                          x4305c4s1b0n0     109       29.9483          None

MPI Time, s                          x4305c2s4b0n0     130       29.946           None

MPI Time, s                          x4305c3s2b0n0     103       29.9379          None

MPI Time, s                          x4305c3s2b0n0     167       29.9315          None

MPI Time, s                          x4305c2s7b0n0      53       29.9133          None

MPI Time, s                          x4305c3s3b0n0     136       29.9061          None

MPI Time, s                          x4305c3s1b0n0      70       29.901           None

……

GPU Occupancy, % of Peak Value       x4305c4s2b0n0     N/A       54.4             None

GPU Occupancy, % of Peak Value       x4305c4s1b0n0     N/A       54.4             None

GPU Occupancy, % of Peak Value       x4305c4s0b0n0     N/A       54.4             None

GPU Occupancy, % of Peak Value       x4305c3s6b0n0     N/A       54.4             None

GPU Occupancy, % of Peak Value       x4305c4s3b0n0     N/A       54.3             None

GPU Occupancy, % of Peak Value       x4305c3s3b0n0     N/A       54.3             None

GPU Occupancy, % of Peak Value       x4305c3s5b0n0     N/A       54.2             None
```

# ADVISOR ROOFLINE ANALYSIS

```
$ mpiexec -n 1 -ppn 12 gpu_tile_compact.sh advisor --collect=survey --profile-gpu --project-dir=Advisor_results --
./Comp_GeoSeries_omp_mpicxx_DP 2048 1000 : -n 11 -ppn 12 gpu_tile_compact.sh ./Comp_GeoSeries_omp_mpicxx_DP 2048 1000
```

```
$ mpiexec -n 1 -ppn 12 gpu_tile_compact.sh advisor --collect=tripcounts --profile-gpu --flop --no-trip-counts --project-
dir=Advisor_results -- ./Comp_GeoSeries_omp_mpicxx_DP 2048 1000 : -n 11 -ppn 12 gpu_tile_compact.sh
./Comp_GeoSeries_omp_mpicxx_DP 2048 1000
```

```
$ advisor --report=all --project-dir=Advisor_results --report-output=Advisor_results/roofline_all.html
```

# THAPI / IPROF ON AURORA
## A lightweight tool for tracing

```
$ module load thapi
$ mpirun -n 12 -ppn 12 gpu_tile_compact.sh iprof ./Comp_GeoSeries_omp_mpicxx_DP 2048 1000
```

```
THAPI: Trace location: /home/jkwack/thapi-traces/thapi_aggreg--2024-09-25T06:27:15-05:00
BACKEND_MPI | 1 Hostnames | 12 Processes | 12 Threads |


        Name |     Time | Time(%) | Calls |  Average |      Min |     Max |
    MPI_Init |   16.76s |  98.13% |    12 |    1.40s | 434.72ms |   2.48s |
MPI_Finalize | 219.50ms |   1.29% |    12 |  18.29ms |  17.48ms | 22.41ms |
  MPI_Reduce |  99.94ms |   0.59% |    96 |   1.04ms |    908ns | 19.77ms |
MPI_Comm_rank |  13.49us |   0.00% |    12 |   1.12us |    504ns |  3.27us |
MPI_Comm_size |   7.60us |   0.00% |    12 | 633.67ns |    485ns |   733ns |
       Total |   17.08s | 100.00% |   144 |
......
```

```
BACKEND_OMP | 1 Hostnames | 12 Processes | 12 Threads |


                                  Name |     Time | Time(%) | Calls | Average |     Min |     Max |
                   ompt_target_exit_data | 468.50ms |  39.62% |    12 | 39.04ms | 31.71ms | 45.23ms |
   ompt_target_data_transfer_from_device | 460.39ms |  38.93% |    12 | 38.37ms | 31.12ms | 44.28ms |
                             ompt_target | 191.37ms |  16.18% |   132 |  1.45ms |  1.30ms |  9.28ms |
                  ompt_target_enter_data |  30.59ms |   2.59% |    12 |  2.55ms |  1.84ms |  3.58ms |
     ompt_target_data_transfer_to_device |  20.24ms |   1.71% |    12 |  1.69ms |  1.19ms |  2.46ms |
                 ompt_target_submit_emi |   9.36ms |   0.79% |   132 | 70.90us |  7.16us |  1.13ms |
                  ompt_target_data_alloc |   1.33ms |   0.11% |    24 | 55.28us | 24.86us | 93.15us |
                 ompt_target_data_delete | 775.82us |   0.07% |    24 | 32.33us |  5.08us | 67.70us |
                                   Total |    1.18s | 100.00% |   360 |
```

```
BACKEND_ZE | 1 Hostnames | 12 Processes | 12 Threads |


                              Name |    Time | Time(%) | Calls | Average |     Min |      Max | Error |
                      zeModuleCreate |   2.72s |  77.65% |   132 | 20.61ms | 103.89us | 224.80ms |     0 |
       zeCommandListAppendMemoryCopy | 483.38ms |  13.80% |   180 |  2.69ms |  9.00us |  44.28ms |     0 |
              zeEventHostSynchronize | 197.31ms |   5.63% |   312 | 632.41us |    108ns |   9.25ms |     0 |
                        zeEventCreate |  28.86ms |   0.82% | 49920 | 578.08ns |    223ns | 146.23us |     0 |
          zeCommandListCreateImmediate |  22.23ms |   0.63% |    24 | 926.29us | 55.87us |   2.90ms |     0 |
                       zeModuleDestroy |  10.31ms |   0.29% |   132 | 78.08us |  5.86us | 478.02us |     0 |
                        zeEventDestroy |   9.04ms |   0.26% | 49920 | 181.14ns |    108ns |  23.10us |     0 |
            zeContextMakeMemoryResident |   7.99ms |   0.23% |    84 | 95.15us |  5.18us | 610.76us |     0 |
          zeCommandListAppendLaunchKernel |   7.28ms |   0.21% |   132 | 55.12us |  6.57us | 636.24us |     0 |
                    zeCommandQueueCreate |   3.23ms |   0.09% |    12 | 268.84us | 232.46us | 299.16us |     0 |
                        zeMemAllocDevice |   2.55ms |   0.07% |    84 | 30.41us | 13.05us |  69.45us |     0 |
   zeDriverGetExtensionFunctionAddress |   2.00ms |   0.06% |   132 | 16.69us |    289ns | 233.68us |    12 |
                        zeKernelCreate |   1.83ms |   0.05% |  1752 |  1.04us |    684ns |  12.13us |     0 |
......
```

```
Device profiling | 1 Hostnames | 12 Processes | 12 Threads | 12 Devices | 12 Subdevices |


                                 Name |     Time | Time(%) | Calls |   Average |    Min |     Max |
zeCommandListAppendMemoryCopy(D2M) | 203.59ms |  51.65% |    12 |  16.97ms | 6.03ms | 30.29ms |
                      Comp_Geo_l29 | 172.75ms |  43.83% |   132 |   1.31ms | 1.29ms |  1.40ms |
zeCommandListAppendMemoryCopy(M2D) |  17.58ms |   4.46% |    96 | 183.11us |   80ns |  2.10ms |
zeCommandListAppendMemoryCopy(S2M) | 217.76us |   0.06% |    48 |   4.54us | 1.28us | 16.48us |
zeCommandListAppendMemoryCopy(M2M) |  24.40us |   0.01% |    12 |   2.03us | 1.36us |  2.80us |
zeCommandListAppendMemoryCopy(M2S) |    960ns |   0.00% |    12 |  80.00ns |   80ns |    80ns |
                             Total | 394.16ms | 100.00% |   312 |


Explicit memory traffic (BACKEND_MPI) | 1 Hostnames | 12 Processes | 12 Threads |


      Name | Byte | Byte(%) | Calls | Average | Min | Max |
MPI_Reduce | 768B | 100.00% |    96 |  8.00B |  8B |  8B |
     Total | 768B | 100.00% |    96 |


Explicit memory traffic (BACKEND_OMP) | 1 Hostnames | 12 Processes | 12 Threads |


                                 Name |     Byte | Byte(%) | Calls | Average |     Min |     Max |
              ompt_target_data_alloc | 805.31MB |  50.00% |    24 | 33.55MB | 33.55MB | 33.55MB |
  ompt_target_data_transfer_to_device | 402.65MB |  25.00% |    12 | 33.55MB | 33.55MB | 33.55MB |
ompt_target_data_transfer_from_device | 402.65MB |  25.00% |    12 | 33.55MB | 33.55MB | 33.55MB |
             ompt_target_data_delete |       0B |   0.00% |    24 |   0.00B |      0B |      0B |
                               Total |   1.61GB | 100.00% |    72 |


……
```

```
Explicit memory traffic (BACKEND_ZE) | 1 Hostnames | 12 Processes | 12 Threads |


                              Name |     Byte | Byte(%) | Calls | Average |      Min |      Max |
          zeContextMakeMemoryResident | 845.41MB |  51.21% |    84 | 10.06MB |      8B | 33.55MB |
  zeCommandListAppendMemoryCopy(M2D) | 402.71MB |  24.40% |    96 |  4.19MB |      4B | 33.55MB |
  zeCommandListAppendMemoryCopy(D2M) | 402.65MB |  24.39% |    12 | 33.55MB | 33.55MB | 33.55MB |
  zeCommandListAppendMemoryCopy(S2M) |   2.40kB |   0.00% |    48 | 50.00B |      8B |    144B |
  zeCommandListAppendMemoryCopy(M2S) |     768B |   0.00% |    12 | 64.00B |     64B |     64B |
  zeCommandListAppendMemoryCopy(M2M) |     684B |   0.00% |    12 | 57.00B |     57B |     57B |
                             Total |   1.65GB | 100.00% |   264 |
```

# KNOWN ISSUES
## w/ the latest SDK (oneapi/eng-compiler/2024.07.30.002 )

▪ Issue 1:
- – Symptom: got the following error message at the end
  - • Exception: 0xb, Segmentation fault
  - • Module: libswip.so
- – Workaround: Before running Advisor, VTune, and APS, do the following:
  - • export SWIP_NULL_SOCKET=1
- – W/ 2025.0 version, the issue will be gone

▪ Issue 2:
- – Symptom: got the following error message after finishing collection
  - • vtune: Error: Cannot stop collection of GPU events
- – Workaround: it is a false alarm. You can ignore it.
- – Intel is investigating this issue now.

▪ Any other issues: please send an email to support@alcf.anl.gov or jkwack@anl.gov (JaeHyuk Kwack)

Argonne
NATIONAL LABORATORY

# THANKS!

Argonne
NATIONAL LABORATORY