

Crusher Tips & Tricks

Tom Papatheodore

HPC Engineer

System Acceptance & User Environment Group

Oak Ridge Leadership Computing Facility (OLCF)

ECP BOF Days – May 10, 2022

ORNL is managed by UT-Battelle LLC for the US Department of Energy





GPU-to-GRU Performance Expectations

- Is your application GPU-bound or data-transfer-bound?
 - If GPU-bound, is it compute-bound or memory-bound?
 - Consider Compute performance, HBM, and HBM bandwidth in table
 - Also consider register usage, occupancy, launch configuration parameters
 - This is where kernel profiling can be helpful
 - If data-transfer-bound, consider CPU-to-GPU bandwidth in table.

GPU-to-GPU Comparison

	Summit	Crusher (per GCD)	Crusher / Summit
Compute performance	~7.8 TF	~26 TF	~3.3X
HBM	16 GB	64 GB	4X
HBM bandwidth	0.9 TB/s	1.6 TB/s	~1.8X
CPU-to-GPU bandwidth	50 GB/s	36 GB/s	~0.7X
L1 Cache	Up to 128 KB	16 KB	0.125X – 0.5X
L2 Cache	6 MB	8 MB	1.3X



Slurm Tips – Flags and Completed Jobs

-u flag gives unbuffered output, which can be helpful when debugging.

-1 flag prepends the task ID to lines of stdout.

To show completed jobs in a specific time period, specify a start (-S) and end (-E) time

The default time window depends on other options (see man sacct)

ition Accoun	t AllocCPUS	State E	lxitCode
batch stf01	6 128	COMPLETED	0:0
batch stf01	6 128	COMPLETED	0:0
batch stf01	.6 128	COMPLETED	0:0
batch stf01	.6 256	COMPLETED	0:0
batch stf01	6 256	COMPLETED	0:0
batch stf01	.6 256	CANCELLED+	0:0
batch stf01	.6 256	COMPLETED	0:0
batch stf01	.6 128	COMPLETED	0:0
batch stf01	.6 128	FAILED	6:0
	batchstf01batchstf01batchstf01batchstf01batchstf01batchstf01batchstf01batchstf01	batchstf016128batchstf016256batchstf016256batchstf016256batchstf016128batchstf016128batchstf016128	batchstf016128COMPLETEDbatchstf016256COMPLETEDbatchstf016256COMPLETEDbatchstf016256CANCELLED+batchstf016256COMPLETEDbatchstf016128COMPLETEDbatchstf016128FAILED

CAK RIDGE

Then you can drill

into more details

about each job

with the - j flag

and customized

output.

Slurm Tips – Capture Job Information



108121.extern 2022-05-05T17:35:03

108121.0 2022-05-05T17:35:04

CAK RIDGE

Open slide master to edit

00:02:34

00:02:32

Slurm Tips – Test CPU/GPU Binding Before Running

[crusher: ~]\$ module load rocm
[crusher: ~]\$ module load craype-accel-amd-gfx90a

[crusher: ~]\$ git clone https://code.ornl.gov/olcf/hello_jobstep.git

[crusher: ~]\$ cd hello_jobstep/

[crusher: ~/hello_jobstep]\$ make
crusher
CC -std=c++11 -fopenmp --rocm-path=/opt/rocm-4.5.0 -x hip -D__HIP_ARCH_GFX90A__=1 --offload-arch=gfx90a -I/opt/rocm4.5.0/include -c hello_jobstep.cpp
CC -fopenmp --rocm-path=/opt/rocm-4.5.0 -L/opt/rocm-4.5.0/lib -lamdhip64 hello_jobstep.o -o hello_jobstep

[crusher: ~/hello_jobstep]\$ OMP_NUM_THREADS=2 srun -N1 -n8 -c8 --gpus-per-node=8 --gpu-bind=closest ./hello_jobstep | sort

MPI 000 - OMP 000 - HWT 001 - Node crusher166 - RT GPU ID 0 - GPU ID 4 - Bus ID d1 MPI 000 - OMP 001 - HWT 004 - Node crusher166 - RT GPU ID 0 - GPU ID 4 - Bus ID d1 MPI 001 - OMP 000 - HWT 008 - Node crusher166 - RT GPU ID 0 - GPU ID 5 - Bus ID d6 MPI 001 - OMP 001 - HWT 012 - Node crusher166 - RT GPU ID 0 - GPU ID 5 - Bus ID d6 MPI 002 - OMP 000 - HWT 018 - Node crusher166 - RT GPU ID 0 - GPU ID 2 - Bus ID c9 MPI 002 - OMP 001 - HWT 021 - Node crusher166 - RT GPU ID 0 - GPU ID 2 - Bus ID c9 MPI 003 - OMP 000 - HWT 025 - Node crusher166 - RT GPU ID 0 - GPU ID 3 - Bus ID ce MPI 003 - OMP 001 - HWT 028 - Node crusher166 - RT GPU ID 0 - GPU ID 3 - Bus ID ce MPI 004 - OMP 000 - HWT 032 - Node crusher166 - RT GPU ID 0 - GPU ID 6 - Bus ID d9 MPI 004 - OMP 001 - HWT 037 - Node crusher166 - RT GPU ID 0 - GPU ID 6 - Bus ID d9 MPI 005 - OMP 000 - HWT 040 - Node crusher166 - RT GPU ID 0 - GPU ID 7 - Bus ID de MPI 005 - OMP 001 - HWT 044 - Node crusher166 - RT GPU ID 0 - GPU ID 7 - Bus ID de MPI 006 - OMP 000 - HWT 049 - Node crusher166 - RT GPU ID 0 - GPU ID 0 - Bus ID c1 MPI 006 - OMP 001 - HWT 052 - Node crusher166 - RT GPU ID 0 - GPU ID 0 - Bus ID c1 MPI 007 - OMP 000 - HWT 059 - Node crusher166 - RT GPU ID 0 - GPU ID 1 - Bus ID c6 MPI 007 - OMP 001 - HWT 060 - Node crusher166 - RT GPU ID 0 - GPU ID 1 - Bus ID c6

CAK RIDGE National Laboratory

. . .

https://code.ornl.gov/olcf/hello_jobstep



MPI ranks should target the GPU associated with their L3 cache region



edit



Open slide master to edit



Open slide master to edit

\$ OMP NUM THREADS=1 srun -N1 -n8 -c1 --qpus-per-node=8 ./hello jobstep | sort - OMP 000 - HWT 000 - Node crusher124 - RT GPU ID 0,1,2,3,4,5,6,7 - GPU ID 0,1,2,3,4,5,6,7 - Bus ID c1,c6,c9,ce,d1,d6,d9,de 000 MPI - RT GPU ID 0,1,2,3,4,5,6,7 - GPU ID 0,1,2,3,4,5,6,7 - Bus ID c1,c6,c9,ce,d1,d6,d9,de OMP 000 008 - Node crusher124 MPI 001 -НМТ - Node crusher124 - RT GPU ID 0,1,2,3,4,5,6,7 - GPU ID 0,1,2,3,4,5,6,7 - Bus ID c1,c6,c9,ce,d1,d6,d9,de MPI 002 _ OMP 000 HWT 016 _ - RT GPU ID 0,1,2,3,4,5,6,7 - GPU ID 0,1,2,3,4,5,6,7 - Bus ID c1,c6,c9,ce,d1,d6,d9,de - Node crusher124 MPI 003 - OMP 000 - HWT 024 RT GPU ID 0,1,2,3,4,5,6,7 - GPU ID 0,1,2,3,4,5,6,7 - Bus ID c1,c6,c9,ce,d1,d6,d9,de MPI 004 OMP 000 032 - Node crusher124 -- RT GPU ID 0,1,2,3,4,5,6,7 - GPU ID 0,1,2,3,4,5 - Node crusher124 - OMP 000 040 6,7 - Bus ID c1,c6,c9,ce,d1,d6,d9,de MPI 005 - нмт - RT GPU ID 0,1,2,3,4,5,6,7 - GPU ID 0,1,2,3,4,5,6,7 - Bus ID c1,c6,c9,ce,d1,d6,d9,de - Node crusher124 MPT 006 - OMP 000 - HWT 048 - RT GPU ID 0,1,2,3,4,5,6,7 - GPU ID 0,1,2,3,4,5,6,7 - Bus ID c1,c6,c9,ce,d1,d6,d9,de 007 - OMP 000 - HWT 056 - Node crusher124 MPI



10

Incorrect mapping: can result in all 8 MPI ranks targeting the same GPU.

Open slide master to edit

Debugging Issues on a Compute Node

\$ salloc -A stf016 -t 60 -N 1

salloc: Pending job allocation 108355
salloc: job 108355 queued and waiting for resources
salloc: job 108355 has been allocated resources
salloc: Granted job allocation 108355

Salloc S make Curren bipco

•••

1

Current system: crusher hipcc --amdgpu-target=gfx90a -c p2p.cpp hipcc --amdgpu-target=gfx90a p2p.o -o p2p

\$ srun -n8 -c8 --gpus-per-node=8 --gpu-bind=closest ./p2p --correct

Once on a compute node, you can use rocm-smi, gdb, gstack, etc.

Get process IDs from

ps -ef | grep <username>, top, etc.

	\$ squeue grep tpape108355batch interpretation	athe eract tpapa	athe R	13:58	1 crusher001	
inal 2	<pre>\$ ssh crusher001 \$ hostname crusher001</pre>					
Term	\$ rocm-smishowpid	s DCm System ===== KF1	Management D Processes	Interface		
	KFD process informat.	ion:				
	PID PROCESS NAME	GPU(s)	VRAM USED	SDMA US	ED CU OCCUPANCY	
	65592 p2p	8	2197880832	2 0	0	
	<pre>\$ rocm-smishowmemuse ====================================</pre>					
	GPU[0] : GPU memory use (%): 1					
	GPU[0] : Memory Activity: 363340274					
	GPU[1] : GPU memory use (%): 2					
	GPU[1]	2U[1] : Memory Activity: 363699843				
	GPU[2]	: GPU memory use (%): U				
	GPU[2]	· CPU memory use (%) · A				
	GPU[3]	· Memory Activity· 363537625				
	GPU[4]	U[4] : GPU memory use (%): 0				
	GPU[4]	: Memory Activity: 363009653				
	GPU[5]	: GPU memory use (%): 0				
	GPU[5]	: Memo	ry Activity	: 362536691		
	GPU[6]	: GPU 1	memory use	(%): 0		
	GPU[6]	: Memo	ry Activity	: 362110208		
	GPU[7]	: GPU 1	memory use	(응): 0		
	GPU[7]	: Memo	ry Activity	: 361753816	,	





CAK RIDGE

National Laboratory

12

Making sure you actually ran on a GPU:

\$ srun ... rocprof --stats/a.out

rocprof output with MPI:

Creates 1 output file per process

\$ cat rocprof_wrapper.sh #!/bin/bash rocprof --stats --hip-trace -o my_output_\${SLURM_PROCID}.csv "\$@" \$ srun/rocprof wrapper.sh ./a.out

Creates output file from specific process only

```
$ cat rocprof_wrapper_single.sh
#!/bin/bash
if [ ${SLURM_PROCID} -eq 0 ];then
    rocprof --stats --hip-trace -o my_output_${SLURM_PROCID}.csv "$@"
else
    "$@"
fi
$ srun ... ./rocprof_wrapper_single.sh ./a.out'
```

Viewing trace files

- chrome://tracing
- <u>https://ui.perfetto.dev/</u>

GPU HW Atomics

HW atomics can be much faster than sw-enabled atomics.

Importantly, users must explicitly request HW atomics via the <u>-munsafe-fp-atomics</u> flag. Otherwise, atomic operations will be performed via CAS loops.

As the name of the flag implies, this might be unsafe to do in some situations, and this is determined by the "granularity" of memory allocated by the user. The "granularity" in this case refers to the coherency between the memory being used inside a kernel and the memory in the rest of the system.

- With coarse-grained memory, coherence with the rest of the system (CPU and GPUs) is obtained at synchronization points. (e.g., hipDeviceSynchronize)
- Fine-grained memory allows CPU and GPU (and multiple GPUs) to synchronize while the GPU kernel is running.

Only coarse-grained memory can use HW atomics. If HW atomics are requested on fine-grained memory, they will (silently) produce a no-op (i.e., give incorrect results).



Managed Memory – XNACK

Works by using either zero-copy memory access or page migration on page fault. To understand which will be used, users need to know about XNACK.

XNACK allows AMD GPUs to migrate memory pages (allocated in a unified virtual address space) between the CPU and GPU upon page-fault.

Importantly, support for XNACK must be enabled at both compile-time and run-time.

- A code that is compiled with xnack support **and** runs in an environment with XNACK enabled, will migrate pages between the CPU and GPU on page-fault.
- A code that is not compiled with xnack support and does not run in an environment with XNACK enabled, will still function as expected with managed memory, but the CPU will access the GPU's memory in a zero-copy fashion (CPU reads directly from GPU memory), and vice versa.

XNACK is not enabled by default at runtime, so users unaware of this might find lower-thanexpected performance (using managed memory) and not realize why.

xnack is enabled at compile-time by default, but to use page-migration XNACK needs to be enabled at run-time with HSA_XNACK=1



Helpful Resources



Crusher Quick-Start Guide:

https://docs.olcf.ornl.gov/systems/crusher_quick_start_guide.html

- Known Issues: <u>https://docs.olcf.ornl.gov/systems/crusher_quick_start_guide.html#known-issues</u>
- SW Compatibility: <u>https://docs.olcf.ornl.gov/systems/crusher_quick_start_guide.html#determining-the-compatibility-of-cray-mpich-and-rocm</u>

ECP Frontier COE Confluence:

https://confluence.exascaleproject.org/display/FCOE/Frontier+Center+of+Excellence

• Forum to ask questions of the community – including vendors.

Crusher Office Hours (every Monday):

https://www.olcf.ornl.gov/crusher-office-hours/

• Must sign up, 5 teams per office hour

\$ man intro						
intro_blacs	intro_blas3	intro_craype-api	intro_irt	intro_libsci	intro_PAPI	intro_scalapack
intro_blas1	intro_cblas	intro_dsmml	intro_lapack	intro_mpi	intro_perftools	—
intro_blas2	intro_craypat	intro_hugepages	intro_lapacke	intro_papi	intro_pmi	





Slurm Tips – Test CPU/GPU Binding Before Running

[crusher: ~]\$ module load rocm

[crusher: ~]\$ OMP_NUM_THREADS=1 srun -1 -N1 -n8 -c8 --gpus-per-node=8 --gpu-bind=closest /bin/bash -c 'echo \$(hostname)
\$(grep Cpus_allowed_list /proc/self/status) GPUS: \$ROCR_VISIBLE_DEVICES' | sort

0: crusher076 Cpus_allowed_list: 0-7 GPUS: 4 1: crusher076 Cpus_allowed_list: 8-15 GPUS: 5 2: crusher076 Cpus_allowed_list: 16-23 GPUS: 2 3: crusher076 Cpus_allowed_list: 24-31 GPUS: 3 4: crusher076 Cpus_allowed_list: 32-39 GPUS: 6 5: crusher076 Cpus_allowed_list: 40-47 GPUS: 7 6: crusher076 Cpus_allowed_list: 48-55 GPUS: 0 7: crusher076 Cpus_allowed_list: 56-63 GPUS: 1

This command gives some broad details about the CPUs and GPUs your job step's processes have access to.



Clang, Clang Wrappers, and Cray's Compiler Wrappers

Vendor	Programming Environment	Compiler Module	Language	Compiler Wrapper	Compiler
Cray PrgEnv-cray		cce	С	cc	craycc
	PrgEnv-cray		C++	CC	craycxx Or crayCC
		Fortran	ftn	crayftn - Crai	
AMD PrgEnv-amd		С	cc	amdclang	
	PrgEnv-amd	rocm	C++	CC	amdclang++
		Fortran	ftn	amdflang	
GCC PrgEnv-gnu		gcc	С	cc	\${GCC_PATH}/bin/gcc
	PrgEnv-gnu		C++	CC	<pre>\${GCC_PATH}/bin/g++</pre>
			Fortran	ftn	<pre>\${GCC_PATH}/bin/gfortran</pre>

NOTE: It is highly recommended to use the Cray compiler wrappers whenever possible, but if for some reason you don't want to use the wrappers...

- Cray compilers, use craycc, crayCC, and crayftn (don't use Cray's bare Clang compilers)
- AMD compilers, use amdclang, amdclang++, and amdflang (don't use AMD's bare Clang compilers)

Another Way to Check if GPU is Being Used

Making sure you actually ran on a GPU (Cray OpenMP Offload only!)

```
$ CRAY_ACC_DEBUG=1 srun -A stf016 -t5 -N1 -n1 -c8 --gpus-per-node=8 --gpu-bind=closest ./vAdd_ompGPU
```

srun: job 108498 queued and waiting for resources
srun: job 108498 has been allocated resources

ACC: Transfer 3 items (to acc 805306368 bytes, to host 0 bytes) from vAdd_ompGPU.cpp:30 ACC: Execute kernel __omp_offloading_69_158925e9_main_130_cce\$noloop\$form from vAdd_ompGPU.cpp:30 ACC: Transfer 3 items (to acc 0 bytes, to host 268435456 bytes) from vAdd ompGPU.cpp:30

= 1.0000000000000000
= 0.00000000000000000000000000000000000
= 268435456
= 0.915606

