

### **Omnitrace By Example**

Bob Robey, <u>Gina Sitaraman</u>, Ian Bogle, Giacomo Capodaglio, Asitav Mishra

AMD @ ORNL 29 May 2024

> AMD together we advance\_

#### Acknowledgements

- Jonathan Madsen
- David Galiffi
- Nicholas Curtis
- and the rest of the DCGPU team

#### Agenda

- Omnitrace for Application Profiling and Tracing
- A Simple Ghost Exchange MPI Example Suite
- Orig: CPU implementation
- Ver1: OpenMP<sup>®</sup> offload port with Managed Memory
- Ver2: Add roctx ranges
- Ver3: Allocate MPI buffers on device (Under construction)
- Ver4: Allocate all buffers once
- Ver5: Convert from 2D to 1D indexing
- Ver6: Add explicit data map directives

#### **Omnitrace for Application Profiling and Tracing**

- Get high level view of entire application run
- Holistic view of CPU, GPU, and system activity
- Sampling and binary instrumentation modes
- Visualize in Perfetto



4

[Public]

#### **MPI Ghost Exchange Example Suite**

- Many applications need to exchange ghost cells with adjacent processes
- This example suite, developed by Bob Robey, implements the exchange for a regular cartesian grid
- We start with the examples in Chapter 8 of Parallel and High Performance Computing, Manning Publications
  git clone <u>https://EssentialsOfParallelComputing/Chapter8</u> https://github.com/essentialsofparallelcomputing/Chapter8

cd GhostExchange

- These examples include various versions ranging from simple methods to those using MPI Datatypes and MPI Cartesian topology capabilities
- For GPU-Aware MPI, the versions using MPI Datatypes have not been optimized in most MPI implementations
- We will use the simpler methods. We will pack the column data into to a buffer and send the buffer. We can send row data directly. If corner data is needed, synchronization is required
- From CPU code, we port to GPU and optimize incrementally using Omnitrace to guide the process
- Uses OpenMP target offload mechanism for offloading compute to AMD GPUs
- Repo: <u>https://github.com/amd/HPCTrainingExamples/tree/main/MPI-examples/GhostExchange/GhostExchange\_ArrayAssign</u>

#### **MPI Ghost Exchange Example – How does it work?**

- A rectangular domain is partitioned into a 2D computational grid, distributed among MPI processes
- An initial solution in specified on a cell-wise basis, then advanced with a 5-point stencil averaging operator
- Halo cells are located along the boundary, and around MPI domains (*ghost cells*) when doing parallel runs
- Boundary conditions are of *outflow* type, enforced prior to ghost halo exchanges
- Example of 2-step halo exchange with 3x3 grid of processes, each owning a 4x4 subset of the mesh:





#### **MPI Ghost Exchange Examples – How to run it?**

- Parameters:
  - -x nprocx -y nprocy -i imax -j jmax -h nhalo -t (0 or 1) -c (0 or 1) -I maxIter

nprocx = number of processes in x dimension

nprocy = number of processes in y dimension

imax = number of mesh cells in x dimension

jmax = number of mesh cells in y dimension

nhalo = number of halo layers

maxIter = maximum number of iterations

- -t = enable/disable sync before MPI calls to accurately time MPI overhead
- -c = include/exclude corner cell updates
- Example run on Frontier with 4 ranks:

srun -N1 -n4 -c7 ./GhostExchange -x 2 -y 2 -i 20000 -j 20000 -h 2 -t -c -I 100

#### **Getting Started with Omnitrace - Configuring Omnitrace Runtime**

• First, create a default configuration file

```
omnitrace-avail -G ~/.omnitrace.cfg
export OMNITRACE_CONFIG_FILE=~/.omnitrace.cfg
```

· Contains settings to control Omnitrace runtime behavior, modify settings as desired



Refer to documentation for more omnitrace-avail capabilities: https://rocm.github.io/omnitrace/runtime.html

#### **Running Omnitrace on Ghost Exchange Examples**

- Set up your environment on Frontier
  - module load cce/17.0.0
    module load rocm/5.7.0
    module load omnitrace/1.11.2
    module load craype-accel-amd-gfx90a
    module load cmake/3.23.2
- Build the code

```
mkdir build; cd build; cmake ..; make -j8
```

Instrument the binary

```
omnitrace-instrument -o ./GhostExchange.inst -- ./GhostExchange
```

• Profile the instrumented binary

```
srun -N1 -n4 -c7 --gpu-bind=closest -A <proj> -t 05:00 omnitrace-run --
./GhostExchange.inst -x 2 -y 2 -i 20000 -j 20000 -h 2 -t -c -I 100
```

#### Understanding output from omnitrace-instrument

[omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/available.json'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/available.txt'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/instrumented.json'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/instrumented.txt'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/instrumented.txt'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/excluded.json'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/excluded.txt'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/excluded.txt'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/excluded.txt'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/overlapping.json'... Done [omnitrace][exe] Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/overlapping.json'... Done

[[ssitaram@login05.frontier build]\$ cat omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/available.txt

StartAddress	AddressRange	#Instructions	Ratio	Linkage V	isibility	Module		Function	FunctionSignature
0x21162c	9	3	3.00	local	hidden	/sysdeps/x86_64/crti.S		_fini	_fini
0x211614	23	7	3.29	local	hidden	/sysdeps/x86_64/crti.S		_init	_init
0x20250e	45	13	3.46	global	default	/sysdeps/x86_64/start.S		_start	_start
0x210af0	1963	485	4.05	unknown	unknown	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>//PI-examples/GhostExchange/GhostEx</pre>	Cartesian_print	Cartesian_print
0x20c4d0	7069	1509	4.68	unknown	unknown	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>IPI-examples/GhostExchange/GhostEx</pre>	boundarycondition_update	boundarycondition_update
0x20e070	10868	2403	4.52	unknown	unknown	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>IPI-examples/GhostExchange/GhostEx</pre>	ghostcell_update	ghostcell_update
0x209d00	10184	2198	4.63	unknown	unknown	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>IPI-examples/GhostExchange/GhostEx</pre>	haloupdate_test	haloupdate_test
0x2025e0	29900	6532	4.58	global	default	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>IPI-examples/GhostExchange/GhostEx</pre>	main	main
0x209ab0	591	157	3.76	unknown	unknown	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>IPI-examples/GhostExchange/GhostEx</pre>	parse_input_args	parse_input_args
0x211300	303	88	3.44	global	default	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>IPI-examples/GhostExchange/GhostEx</pre>	malloc2D	malloc2D
0x211430	15	4	3.75	global	default	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>IPI-examples/GhostExchange/GhostEx</pre>	malloc2D_free	malloc2D_free
0x2112a0	13	3	4.33	global	default	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>IPI-examples/GhostExchange/GhostEx</pre>	cpu_timer_start	cpu_timer_start
0x2112b0	71	19	3.74	global	default	/ccs/home/ssitaram/git/HPCTrainingExamples/M	<pre>IPI-examples/GhostExchange/GhostEx</pre>	cpu_timer_stop	cpu_timer_stop
0x20257f	85	23	3.70	local	default	GhostExchange		do_fini	do_fini
0x20253b	68	19	3.58	local	default	GhostExchange		do_init	do_init
0x21143f	320	88	3.64	global	default	GhostExchange		<pre>no_mmap_for_malloc</pre>	no_mmap_for_malloc
0x21160e	4	3	1.33	global	default	elf-init.c		libc_csu_fini	libc_csu_fini
0x21158e	103	36	2.86	global	default	elf-init.c		libc csu init	libc csu init

[[ssitaram@login05.frontier build]\$ cat omnitrace-GhostExchange.inst-output/2024-05-22\_16.24/instrumentation/instrumented.txt

StartAddress	AddressRange	#Instructions	Ratio Linkage	Visibility	Module	Function	FunctionSignature
0x20c4d0	7069	1509	4.68 unknown	n unknowr	/ccs/home/ssitaram/git/HPCTrainingExamples/MPI-examples/GhostExchange/GhostEx	boundarycondition_update	boundarycondition_update
0x20e070	10868	2403	4.52 unknown	n unknowr	/ccs/home/ssitaram/git/HPCTrainingExamples/MPI-examples/GhostExchange/GhostEx	ghostcell_update	ghostcell_update
0x209d00	10184	2198	4.63 unknown	n unknowr	/ccs/home/ssitaram/git/HPCTrainingExamples/MPI-examples/GhostExchange/GhostEx	haloupdate_test	haloupdate_test

#### Functions that were instrumented

Functions that could be instrumented

AMD @ ORNL

#### Understanding output from omnitrace-run

Omnitrace ASCII art is proof that Omnitrace is running, shows version used



omnitrace v1.11.2 (rev: 1df597e049b240fb263e7fcd7bddc78097d27f00, tag: v1.11.2, compiler: GNU v7.5.0, rocm: v5.7.x)

[omnitrace][130341][perfetto]> Outputting '/ccs/home/ssitaram/git/HPCTrainingExamples/MPI-examples/GhostExchange/GhostExchange\_ArrayAssign/Orig/build/omnitrace-GhostExchange.inst-output/2 024-05-22\_16.50/perfetto-trace-1.proto' (765.48 KB / 0.77 MB / 0.00 GB)... Done

[omnitrace][130342][perfetto]> Outputting '/ccs/home/ssitaram/git/HPCTrainingExamples/MPI-examples/GhostExchange/GhostExchange\_ArrayAssign/Orig/build/omnitrace-GhostExchange.inst-output/2 024-05-22\_16.50/perfetto-trace-2.proto' (765.47 KB / 0.77 MB / 0.00 GB)... Done

[omnitrace][130343][perfetto]> Outputting '/ccs/home/ssitaram/git/HPCTrainingExamples/MPI-examples/GhostExchange/GhostExchange\_ArrayAssign/Orig/build/omnitrace-GhostExchange.inst-output/2 024-05-22\_16.50/perfetto-trace-3.proto' (765.47 KB / 0.77 MB / 0.00 GB)... Done

[omnitrace][130340][perfetto]> Outputting '/ccs/home/ssitaram/git/HPCTrainingExamples/MPI-examples/GhostExchange/GhostExchange\_ArrayAssign/Orig/build/omnitrace-GhostExchange.inst-output/2 024-05-22\_16.50/perfetto-trace-0.proto' (765.48 KB / 0.77 MB / 0.00 GB)... Done

Paths to output trace files

#### Visualizing Omnitrace .proto files

#### Copy .proto file to local workstation or laptop, open in Perfetto: <u>https://ui.perfetto.dev/</u>



5/29/2024

### Orig CPU implementation of Ghost Exchange

#### **Orig: First look at Omnitrace profile for Rank 0**

∧ ./GhostExchange.inst 1179	99		
GhostExchange.i 11799		GhostExchange.inst	CPU functions
CPU Context Switches (S)	~	5K	
CPU Frequency [0] (S)	~	5K	
CPU Frequency [1] (S)	~	2.5 K	
CPU Frequency [2] (S)	~	2.5 K	
CPU Frequency [3] (S)	~	2.5 K	
CPU Frequency [4] (S)	~	2.5 K	
CPU Frequency [5] (S)	~	2.5 K	
CPU Frequency [6] (S)	~	2.5 K	
CPU Freauencv [7] (S)	~	2.5 K	
<snip></snip>		CPU characteristics	
CPU Frequency [126] (S)	~	2.5K	
CPU Frequency [127] (S)	~	2.5K	
CPU Kernel Time (S)	~	25	
CPU Memory Usage (S)	~	2.5K	
CPU Page Faults (S)	~	0.5 M	
CPU Peak Memory (S)	~	2.5K	
CPU User Time (S)	~	25	
CPU Virtual Memory Usage (S)	~	7.5K	
5/29/2024		AMD @ ORNL	together we advance

#### **Orig: First profile – zoom in**



#### **Orig: Keep profile short - sample one CPU core**

• Update config file to sample only 1 CPU core:

OMNITRACE\_SAMPLING\_CPUS

• Just rerun, no need to instrument again

<ul> <li>./GhostExchange.inst 73952</li> </ul>		
GhostExchange.i 73952		GhostExchange.inst
CPU Context Switches (S)	$\sim$	4K
CPU Frequency [0] (S)	$\sim$	
CPU Kernel Time (S)	$\sim$	1000m
CPU Memory Usage (S)	$\sim$	2K
CPU Page Faults (S)	$\sim$	450K
CPU Peak Memory (S)	~	2K
CPU User Time (S)	~	15
CPU Virtual Memory Usage (S)	$\sim$	7K

= 0



#### **Orig: Generate CPU-side wall clock times in profile**

• Enable tracking of function durations in config file

OMNITRACE\_PROFILE

= true

• Omnitrace generates wall\_clock files with durations of each instrumented function

[omnitrace][57701][wall\_clock]> Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.44/wall\_clock-3.txt'
[omnitrace][57699][wall\_clock]> Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.44/wall\_clock-1.txt'
[omnitrace][57698][wall\_clock]> Outputting 'omnitrace-GhostExchange.inst-output/2024-05-22\_16.44/wall\_clock-0.txt'

Look for durations of MPI calls here

			REAL-CLO	CK TIMER	(I.E. WALL-CL	_OCK TIMER)					
LABEL	COUNT	DEPTH	METRIC	UNITS	SUM	MEAN	MIN	MAX	VAR	STDDEV	% SELF
0>>> GhostExchange.inst		0	   wall_clock	   sec	12.179875	12.179875	12.179875	12.179875	0.000000	0.000000	97.7
0>>>  _MPI_Init	1	1	wall_clock	sec	0.194357	0.194357	0.194357	0.194357	0.000000	0.000000	100.0
0>>>  _MPI_Comm_rank	4	1	wall_clock	sec	0.000278	0.000069	0.000001	0.000270	0.000000	0.000134	100.0
0>>>  _MPI_Comm_size	2	1	wall_clock	sec	0.000011	0.000006	0.000001	0.000010	0.000000	0.000006	100.0
0>>>  _MPI_Allreduce	1	1	wall_clock	sec	0.001482	0.001482	0.001482	0.001482	0.000000	0.000000	100.0
0>>>  _boundarycondition_update	1	1	wall_clock	sec	0.000352	0.000352	0.000352	0.000352	0.000000	0.000000	100.0
0>>>  _ghostcell_update	101	1	wall_clock	sec	0.081816	0.000810	0.000415	0.002370	0.000000	0.000408	43.2
0>>>  _MPI_Irecv	404	2	wall_clock	sec	0.003434	0.000009	0.000003	0.000171	0.000000	0.000011	100.0
0>>>  _MPI_Isend	404	2	wall_clock	sec	0.002443	0.000006	0.000003	0.000038	0.000000	0.000007	100.0
0>>>  _MPI_Waitall	202	2	wall_clock	sec	0.040556	0.000201	0.000009	0.001980	0.000000	0.000321	100.0

#### **Orig: Use flat profiles for finding hotspots**

• To flatten the hierarchy in the wall clock profile, enable flat profile:

OMNITRACE\_FLAT\_PROFILE

= true

• Now each function appears once, all timings are consolidated for each function:

			REAL-CL	DCK TIMER	(I.E. WALL-(	CLOCK TIMER)					
LABEL	COUNT	DEPTH	METRIC	UNITS	SUM	MEAN	MIN	MAX	VAR	STDDEV	% SELF
0>>> GhostExchange.inst	1	0	wall_clock	sec	12.252487	12.252487	12.252487	12.252487	0.000000	0.000000	100.0
0>>> MPI_Init  0>>> MPI_Comm_rank	1   4	0	wall_clock wall_clock	sec   sec	0.198300   0.000474	0.198300 0.000118	0.198300 0.000001	0.198300 0.000466	0.000000 0.000000	0.000000   0.000232	100.0     100.0
0>>> MPI_Comm_size	2	0	wall_clock	sec sec		0.000009	0.000001	0.000017	0.000000		100.0
0>>> boundarycondition_update		0	wall_clock	sec	0.000304	0.000304	0.000304	0.000304	0.000000	0.000000	100.0
0>>> ghostcell_update  0>>> MPI_Irecv	101   404	0	wall_clock wall_clock	sec   sec	0.079793   0.004957	0.000790 0.000012	0.000439	0.002256 0.000165	0.000000 0.000000	0.000282     0.000013	100.0   100.0
0>>> MPI_Isend	404   202	0	wall_clock	sec sec	0.002565	0.000006	0.000003	0.000043	0.000000	0.000008	100.0
			CIOCK								

Not much here other than MPI calls not being the bottleneck

### Ver1

First GPU implementation of Ghost Exchange OpenMP offload + Managed Memory Programming Model

AMD @ ORNL

#### Ver1: Code changes

Pragma for unified memory added to each translation unit

#pragma omp requires unified\_shared\_memory

• Target offload pragma added to all compute loops

#pragma omp target teams distribute parallel for collapse(2)

- Data still resides in host memory, but accessed from compute kernels and MPI calls
- On Frontier nodes, this means data is moved across AMD Infinity Fabric<sup>™</sup> link between CPU and GPU
- Needs environment variable to enable OS managed page migration

export HSA\_XNACK=1

#### Ver1: Profile shows offloaded compute regions

Many rows in profile, pin select rows to bring them to the top

5/29/2024

GPU kernels

								Flo	ow Eve	ents	
HIP Activity Device 4, Queue 0	Ŧ	omp_offloadir	ng_10008e_2228c6e6_main	_1140	omp_offlo	Id	omp_cifloading_1000		omp		
GhostExchange.i 83591	Ŧ		hipStreamSynchronize		hipStreamSynch		hipStreamSynchronize	GhostExchange.in MPI_W	nst hipStr		
✔ Misc Global Tracks											
∧ ./GhostExchange.inst 83591											
GhostExchange.i 83591	Ŧ		hipStreamSynchronize		hipStreamSynch		hipStreamSynchronize	GhostExchange.ir	nst hipStr		
CPU Context Switches (S)	~	25K									
CPU Frequency [0] (S)	$\sim$	4К				_					
CPU Kernel Time (S)	$\sim$	3									
CPU Memory Usage (S)	~	2К			I						
CPU Page Faults (S)	~	50K									
CPU Peak Memory (S)	~	2К			[					1	
CPU User Time (S)	$\sim$	45									
CPU Virtual Memory Usage (S)	$\sim$	45K									
GPU Busy [0] (S)	$\sim$	100								<u> </u>	
CDI Ruey [1] (0)		0	AMD @ OR	NL							

#### Ver1: Observe GPU characteristics from rocm-smi in profile

∧ ./GhostExchange.inst 83591																
GhostExchange.i 83591		h h h h	h h	h h h	h h	h h	h h	h h h	h	e	hostExcha	nge.inst h h		h h h	h h	h
CPU Context Switches (S)	$\sim$	25K								 						
CPU Frequency [0] (S)	$\sim$	4K		11 1	111	1 111 1	1 1 1	111		 		1	1			1
CPU Kernel Time (S)	$\sim$	3								 						
CPU Memory Usage (S)	$\sim$	2К								 				·		
CPU Page Faults (S)	$\sim$	_50K							-	 						
CPU Peak Memory (S)	$\sim$	2К		_		I				 						
CPU User Time (S)	$\sim$	45								 						
CPU Virtual Memory Usage (S)	$\sim$	45K								 						
GPU Busy [0] (S)	~	100				· · · ·	,		<u> </u>	 · · · · · ·	1	· · · · ·		· · · · · ·		
GPU Busy [1] (S)	$\sim$	0														
GPU Busy [2] (S)	$\sim$	0								 						
GPU Busy [3] (S)	$\sim$	0														
GPU Busy [4] (S)	$\sim$	0								 						

Shows details of all GPUs by default, we see activity only on GPU 0

[Public]

#### **Ver1: Trim profile to GPU of interest**

• Indicate which GPU to sample in config file

#### OMNITRACE\_SAMPLING\_GPUS

<ul> <li>./GhostExchange.inst 14163</li> </ul>											
GhostExchange.i 14163								GhostExchange.	inst		
<u>j</u>		hipStreamSynchron	hipStream	Synchronize	hipStreamSynchro	hipStreamSynchro	hipSt	reamSynchro	hipStreamS	ynchro	hipStreamSynch
CPU Context Switches (S)	~	15K									
CPU Frequency [0] (S)	~	4K		Л		Π	F				
CPU Kernel Time (S)	$\sim$	2									
CPU Memory Usage (S)	$\sim$	2К							I		
CPU Page Faults (S)	~	50K									
CPU Peak Memory (S)	~	2К				1		1			
CPU User Time (S)	$\sim$	45									
CPU Virtual Memory Usage (S)	~	45K									
GPU Busy [0] (S)	~	100	L.		- <u>-</u>	-					
GPU Memory Usage [0] (S)	~	2K									
GPU Power [0] (S)	~	150									
GPU Temperature [0] (S)	~	45						<u></u>	<u></u>		
HIP Activity Device 4, Queue 0		omp_offloading	omp_offloadin	g_10008e_2228c	omp_offloading	omp_offloading	on	p_offloading	omp_offle	oading	omp_offloadin

Concise trace, easier to analyze

= 0

5/29/2024

AMD @ ORNL

#### Ver1: Wall clock profile shows OMP offload kernels and HIP APIs

REAL-CLOCK	TIMER (I	.E. WALL-	CLOCK TIMER)								   
LABEL	COUNT	DEPTH	METRIC	UNITS	SUM	MEAN	MIN	MAX	VAR	STDDEV	% SELF
	34	0	   wall_clock	   sec	0.000226	0.00	Quintim	a incra	hase	-12	2002
0>>> hipRuntimeGetVersion	1	0	wall_clock	sec	0.000132	0.00	Vununn		aseu		5005
0>>> hipDeviceGet	2	0	wall_clock	sec	0.000005	0.00003	0.000002	0.000004	0.000000	0.000001	100.0
10>>> hinGetDeviceCount	1	0	wall clock	sec	0.000006	0.000006	0.000006	0.000006	0.000000	0.000000	100.0
0>>> GhostExchange.inst	1	0	wall_clock	sec	42.833165	42.833165	42.833165	42.833165	0.000000	0.000000	100.0
0>>> MPI_Init	1	0	wall_clock	sec	0.195783	0.195783	0.195783	0.195783	0.000000	0.000000	100.0
0>>> MPI_Comm_rank	4	0	wall_clock	sec	0.000037	0.000009	0.000001	0.000029	0.000000	0.000013	100.0
0>>> MPI_Comm_size	2	0	wall_clock	sec	0.000012	0.000006	0.000002	0.000010	0.000000	0.000006	100.0
0>>> MPI_Allreduce	1	0	wall_clock	sec	0.000173	0.000173	0.000173	0.000173	0.000000	0.000000	100.0
0>>> hipDeviceComputeCapability	1	0	wall_clock	sec	0.000002	0.000002	0.000002	0.000002	0.000000	0.000000	100.0
0>>> hipDeviceGetName	1	0	wall_clock	sec	0.000001	0.000001	0.000001	0.000001	0.000000	0.000000	100.0
0>>> hipDeviceGetAttribute	1	0	wall_clock	sec	0.000003	0.00003	0.000003	0.000003	0.000000	0.000000	100.0
0>>> hipGetDeviceProperties	1 1	0	wall_clock	sec	0.000005	0.000005	0.000005	0.000005	0.000000	0.000000	100.0
0>>> hipGetDevice	543	0	wall_clock	sec	0.001256	0.000002	0.000001	0.000016	0.000000	0.000002	100.0
0>>> hipHostMalloc	1	0	wall_clock	sec	0.000101	0.000101	0.000101	0.000101	0.000000	0.000000	100.0
0>>> hipEventCreate	2	0	wall_clock	sec	0.000013	0.000007	0.000001	0.000012	0.000000	0.000008	100.0
0>>> hipStreamCreate	j 1	j 0	wall_clock	sec	0.562414	0.562414	0.562414	0.562414	0.000000	0.000000	100.0
0>>> hipModuleLoadData	1 1	0	wall_clock	sec	0.002285	0.002285				-	
0>>> hipPointerGetAttributes	1	0	wall_clock	sec	0.000004	0.000004	I Mo	stlv wa	itina fa	hr GPI	
0>>> hipModuleGetGlobal	38	0	wall_clock	sec	0.000120	0.000003					9
0>>> hipModuleGetFunction	7	0	wall_clock	sec	0.000026	0.000004	kor	nole to	comp	loto	
0>>> hipFuncGetAttribute	14	0	wall_clock	sec	0.000017	0,00001	L VCI		Comp	IELE	
0>>> hipFuncSetCacheConfig	7	0	wall_clock	sec	0.000009	0.000001	0.000001	0.000002	0.000000	0.000000	100.0
0>>> hipModuleLaunchKernel	506	i 0	wall clock	sec	0.003388	/0.000007	0.00003	0.000033	0.000000	0.000003	100.0
0>>> hipStreamSynchronize	506	0	wall_clock	sec	41.607516	0.082228	0.000232	0.697088	0.026855	0.163875	100.0
0>>>omp_offloading_10008e_2228c6e6_main_1115_cce\$noloop\$form	1	0	wall_clock	sec	0.000007	0.000007	0.000007	i 0.000007	0.000000	i 0.000000	i 100.0 i
<pre>0&gt;&gt;&gt;omp_offloading_10008e_2228c6e6_main_1125_cce\$noloop\$form</pre>	1	0	wall_clock	sec	0.000000	0.00000	0.		ما ما م		المعينيات جاج
0>>>omp_offloading_10008e_2228c6e6Z24boundarycondition_updatePPdiiiiiii_l179_cce\$noloop\$form	101	0	wall_clock	sec	0.000017	0.00000	0. <u> </u>	<u> </u>	<u>U-SIQE</u>	<u>raun</u>	<u>ch durati</u>
0>>> MPI_Irecv	404	0	wall_clock	sec	0.003960	0.000010	0.				
0>>> MPI_Isend	404	0	wall_clock	sec	0.002559	0.000006	0. O	Open	VIP off	load_k	ernels
0>>> MPI_Waitall	202	0	wall_clock	sec	0.309031	0.001530	0.				
<pre> 0&gt;&gt;&gt;omp_offloading_10008e_2228c6e6Z24boundarycondition_updatePPdiiiiiii_1197_cce\$noloop\$form</pre>	101	0	wall_clock	sec	0.000022	5.000000	0. st	<u>nown</u> ir	ı wall-	clock	orofile
0>>>omp_offloading_10008e_2228c6e6Z16ghostcell_updatePPdiiiiiiiiii1252_cce\$noloop\$form	101	0	wall_clock	sec	0.000026	0.00000	0.		- vvan		
0>>>omp_offloading_10008e_2228c6e6Z16ghostcell_updatePPdiiiiiiiii1273_cce\$noloop\$form	101	0	wall_clock	sec	0.000020	0.00000	0.00000	0.000000	0.000000	0.000000	100.0
0>>>omp_offloading_10008e_2228c6e6_main_1140_cce\$noloop\$form	100	0	wall_clock	sec	0.000017	0.00000	0.000000	0.000000	0.000000	0.000000	100.0

Managed memory affects kernel performance – but profile does not show in what way, yet

- To implement OpenMP offload capability,
  - the AMD compiler uses the HSA layer
  - the Cray compiler uses the HIP layer
- Set up config file to see HSA activity in profile:

OMNITRACE\_ROCTRACER\_HSA\_ACTIVITY = true
OMNITRACE\_ROCTRACER\_HSA\_API = true



More details about <u>HIP</u> and <u>HSA</u> runtime libraries

[Public]

### Ver2

Manually instrument code with roctx ranges to study regions of code

#### **Ver2: Profile shows roctx ranges**

roctx region

GhostExchange.i 115987								GhostExchange.inst		
			Stencil	Βοι	undaryUpdate	e		GhostCellUpdate		
			hipStreamSynchronize	hipStrea	amSyn hip	oStre L	oadLeftRight	MPILeftRightExchange	Unp	MPIUpDownE
						hi	pStreamSyn	MPI_Waitall	hip	MPI_Waitall
CPU Context Switches (S)	$\sim$	15K								
CPU Frequency [0] (S)	$\sim$	4K								
CPU Kernel Time (S)	$\sim$	2								
CPU Memory Usage (S)	$\sim$	2K		l				1		
CPU Page Faults (S)	$\sim$	50K								
CPU Peak Memory (S)	$\sim$	2K		1				1		
CPU User Time (S)	$\sim$	45								
CPU Virtual Memory Usage (S)	$\sim$	45K								
GPU Busy [0] (S)	$\sim$	100		· · ·						
GPU Memory Usage [0] (S)	$\sim$	_2K								
GPU Power [0] (S)	$\sim$	150								
GPU Temperature [0] (S)	$\sim$	40								
HIP Activity Device 4, Queue 0		om	o_offloading_10008e_2237056d	or	mp_of		omp_o			

5/29/2024

#### Ver2: Wall clock files show timings of roctx regions

REAL-CLOCK	TIMER (I.	E. WALL-C	CLOCK TIMER)								 
LABEL	COUNT	DEPTH	METRIC	UNITS	SUM	MEAN	MIN	MAX	VAR	STDDEV	% SELF
	34	0	wall_clock	sec	0.000207	0.000006	0.000003	0.000070	0.000000	0.000011	100.0
<pre>3&gt;&gt;&gt; hipRuntimeGetVersion</pre>	1	0	wall_clock	sec	0.000152	0.000152	0.000152	0.000152	0.000000	0.000000	100.0
<pre>d&gt;&gt;&gt; hipDeviceGet</pre>	2	0	wall_clock	sec	0.000005	0.000003	0.000002	0.000003	0.000000	0.000001	100.0
<pre>&gt;&gt;&gt; hipGetDeviceCount</pre>	1	0	wall_clock	sec	0.000016	0.000016	0.000016	0.000016	0.000000	0.000000	100.0
>>>> GhostExchange.inst	1	0	wall_clock	sec	42.948441	42.948441	42.948441	42.948441	0.000000	0.000000	100.0
>>> MPI_Init	1	0	wall_clock	sec	0.222832	0.222832	0.222832	0.222832	0.000000	0.000000	100.0
>>> MPI_Comm_rank	4	0	wall_clock	sec	0.000039	0.000010	0.000001	0.000030	0.000000	0.000014	100.0
>>> MPI_Comm_size	2	0	wall_clock	sec	0.000012	0.000006	0.000002	0.000011	0.000000	0.000007	100.0
>>> MPI_Allreduce	1	0	wall_clock	sec	0.000697	0.000697	0.000697	0.000697	0.000000	0.000000	100.0
>>> hipDeviceComputeCapability	1	0	wall_clock	sec	0.000003	0.000003	0.000003	0.000003	0.000000	0.000000	100.0
>>> hipDeviceGetName	1	0	wall_clock	sec	0.000001	0.000001	0.000001	0.000001	0.000000	0.000000	100.0
>>> hipDeviceGetAttribute	1	0	wall_clock	sec	0.000003	0.000003	0.000003	0.000003	0.000000	0.000000	100.0
>>> hipGetDeviceProperties	1	0	wall_clock	sec	0.000002	0.000002	0.000002	0.000002	0.000000	0.000000	100.0
>>> hipGetDevice	543	0	wall_clock	sec	0.001649	0.000003	0.000001	0.000018	0.000000	0.000002	100.0
>>> hipHostMalloc	1	0	wall_clock	sec	0.000110	0.000110	0.000110	0.000110	0.000000	0.000000	100.0
>>> hipEventCreate	2	0	wall_clock	sec	0.000007	0.000003	0.000001	0.000005	0.000000	0.000003	100.0
>>> hipStreamCreate	1	0	wall_clock	sec	0.557462	0.557462	0.557462	0.557462	0.000000	0.000000	100.0
>>> hipModuleLoadData	1	0	wall_clock	sec	0.002125	0.002125	0.002125	0.002125	0.000000	0.000000	100.0
>>> hipPointerGetAttributes	1	0	wall_clock	sec	0.000005	0.000005	0.000005	0.000005	0.000000	0.000000	100.0
>>> hipModuleGetGlobal	38	0	wall_clock	sec	0.000120	0.000003	0.000001	0.000010	0.000000	0.000002	100.0
>>> hipModuleGetFunction	7	0	wall_clock	sec	0.000023	0.000003	0.000002	0.000005	0.000000	0.000001	100.0
>>> hipFuncGetAttribute	14	0	wall_clock	sec	0.000019	0.000001	0.000001	0.000003	0.000000	0.000001	100.0
>>> hipFuncSetCacheConfig	7	0	wall_clock	sec	0.000008	0.000001	0.000001	0.000001	0.000000	0.000000	100.0
>>> hipModuleLaunchKernel	506	0	wall_clo								
>>> hipStreamSynchronize	506	.0	<mark>  wall_c</mark> lo	Call	count s	shows	we allo	ncate ł	Nuffers	: man\	/ time
>>>omp_offloading_10008e_2237056d_main_1116_cce\$noloop\$form	1	0	wall_cl							<u>, many</u>	
>>>omp_offloading_10008e_2237056d_main_1126_cce\$noloop\$form	1	0	wall_clock	sec	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	100.0
>>> MPTRequest	101	0	wall_clock	sec	0.000101	0.000001	0.000001	0.000012	0.000000	0.000001	100.0
>>> BufAlloc	101	0	wall_clock	sec	0.000110	0.000001	0.000001	0.000009	0.000000	0.000001	100.0
>>> LoadLeftRight	101	0	wall_clock	sec	0.032994	0.000327	0.000122	0.002041	0.000000	0.000264	100.0
>>>omp_offloading_10008e_2237056dZ24boundarycondition_updatePPdiiiiiiii191_cce\$noloop\$form	101	0	wall_clock	sec	0.000020	0.000000	0.000000	0.000000	0.000000	0.000000	100.0
>>> MPILeftRightExchange	101	0	wall_clock	sec	0.039038	0.000387	0.000068	0.012803	0.000002	0.001273	100.0
>>> MPI_Irecv	404	0	wall_clock	sec	0.003454	0.000009	0.000005	0.000035	0.000000	0.000005	100.0
>>> MPI_Isend	404	0	wall_clock	sec	0.002998	0.000007	0.000005	0.000033	0.000000	0.000004	100.0
>>> MPI_Waitall	202	0	wall_clock	sec	0.047174	Tim		nt maa			to lea
>>> UnpackLeftRight	101	0	wall_clock	sec	0.028112		ie spe	ni mos	suy in s	compl	ле ке
>>>omp_offloading_10008e_2237056dZ24boundarycondition_updatePPdiiiiiiii1209_cce\$noloop\$form	101	0	wall_clock	sec	0.000030	0.000000	0.000000	0.000001	0.000000	0.000000	100.0
>> MPTIIpDownExchange	101	Q	wall_clock	SPC	0.020360	0.000202	0.000067	0.001836	0.000000	0.000257	100.0
>>> Stencil	100	0	wall_clock	sec	41.254781	0.412548	0.405178	0.684404	0.000755	0.027483	100.0
>>>omp_offloading_10008e_2237056dZ16ghostcell_updatePPdiiiiiiiiii1269_cce\$noloop\$form	101	0	wall_clock	sec	0.000031	0.000000	0.000000	0.000001	0.000000	0.000000	100.0
>>> BoundaryUpdate	100	0	wall_clock	sec	0.096786	0.000968	0.000297	0.003505	0.000000	0.000704	100.0
>>>omp_offloading_10008e_2237056dZ16ghostcell_updatePPdiiiiiiiiii1294_cce\$noloop\$form	101	0	wall_clock	sec	0.000024	0.000000	0.000000	0.000001	0.000000	0.000000	100.0
>>>omp_offloading_10008e_2237056d_main_1144_cce\$noloop\$form	100	0	wall_clock	sec	0.000021	0.000000	0.000000	0.000001	0.000000	0.000000	100.0
	100	0	unall alaak		1 0 1000/F	0 001000	0 000502	0 017501	0 000002	0 001710	100 0 1



[Public]

### Ver3

#### Allocate MPI buffers on device

5/29/2024

AMD @ ORNL

#### Ver3: Enable GPU Aware MPI

- On Frontier nodes, there are 4 Network Interface Cards (NICs) directly attached to odd numbered GCDs
- · Using device buffers in MPI calls is expected to improve communication performance
- Need environment variable to enable GPU Aware MPI:

export MPICH\_GPU\_SUPPORT\_ENABLED=1

• In this example, communication is not our bottleneck, so we don't expect a big difference in MPI overhead



- Our attempt to achieve this is not working with ROCm 5.7.0 today on Frontier
- Known issue that is being actively investigated

### Ver4

#### Allocate all host buffers just once

REAL-CLOCK TIMER (I.E. WALL-CLOCK TIMER)											
LABEL	DEPTH	METRIC	UNITS	SUM	MEAN	MIN	MAX	VAR	STDDEV	% SELF	
	34	0	wall clock	sec	0.000204	0.000006	0.000002	0.000070	0.000000	0.000012	100.0
0>>> hipRuntimeGetVersion	1	0	wall clock	sec	0.000132	0.000132	0.000132	0.000132	0.000000	0.000000	100.0
0>>> hipDeviceGet	2	0	wall clock	sec	0.000005	0.000003	0.000002	0.000003	0.000000	0.000001	100.0
0>>> hipGetDeviceCount	1	0	wall_clock	sec	0.000006	0.000006	0.000006	0.000006	0.000000	0.000000	100.0
0>>> GhostExchange.inst	1	j 0	wall_clock	sec	43.091367	43.091367	43.091367	43.091367	0.000000	0.000000	100.0
0>>> MPI_Init	1	j 0	wall_c <u>lock</u>	sec	0.288942	0.288942	0.288942	0.288942	0.000000	0.000000	100.0
0>>> MPI_Comm_rank	4	0	wall_c				<b>A</b> 11		0.000000	0.000017	100.0
0>>> MPI_Comm_size	2	1 0	wall_c	rofile	e show	's only	1 alloc	ation	0.000000	0.000006	100.0
Allreduce	1	- 0	wall_c						0.000000	0.000000	100.0
0>>> BufAlloc	1	0	wall_clock	sec	0.000013	0.000013	0.000013	0.000013	0.000000	0.000000	100.0
0>>> hipDeviceComputeCapability	1	1 0	wall_clock	sec	0.000002	0.000002	0.000002	0.000002	0.000000	0.00000	100.0
0>>> hipDeviceGetName	1	0	wall_clock	sec	0.000001	0.000001	0.000001	0.000001	0.000000	0.000000	100.0
0>>> hipDeviceGetAttribute	1	0	wall_clock	sec	0.000003	0.000003	0.000003	0.000003	0.000000	0.000000	100.0
0>>> hipGetDeviceProperties	1	0	wall_clock	sec	0.000005	0.000005	0.000005	0.000005	0.000000	0.00000	100.0
0>>> hipGetDevice	543	0	wall_clock	sec	0.001684	0.000003	0.000001	0.000023	0.000000	0.000002	100.0
0>>> hipHostMalloc	1	0	wall_clock	sec	0.000098	0.000098	0.000098	0.000098	0.000000	0.000000	100.0
0>>> hipEventCreate	2	0	wall_clock	sec	0.000012	0.000006	0.000001	0.000011	0.000000	0.000007	100.0
0>>> hipStreamCreate	1	0	wall_clock	sec	0.481629	0.481629	0.481629	0.481629	0.000000	0.000000	100.0
0>>> hipModuleLoadData	1	0	wall_clock	sec	0.002023	0.002023	0.002023	0.002023	0.000000	0.000000	100.0
0>>> hipPointerGetAttributes	1	0	wall_clock	sec	0.000004	0.000004	0.000004	0.000004	0.000000	0.000000	100.0
0>>> hipModuleGetGlobal	38	0	wall_clock	sec	0.000118	0.000003	0.000001	0.000007	0.000000	0.000002	100.0
0>>> hipModuleGetFunction	1	0	wall_clock	sec	0.000023	0.000003	0.000002	0.000005	0.000000	0.000001	100.0
0>>> hipFuncGetAttribute	14	0	wall_clock	sec	0.000020	0.000001	0.000001	0.000004	0.000000	0.000001	100.0
0>>> nipFuncSetUacheContig	50/	0	wall_clock	sec	0.000008	0.000001	0.000001	0.000001	0.000000	0.000000	100.0
U>>> nipModuleLaunchkernel	506	0	wall_clock	sec	0.003/72	0.000007	0.000004	0.000033	0.000000	0.000003	100.0
0>>> hipStreamsynchronize	500	0	wall_clock	sec	41.886630	0.082/80	0.000013	0.699256	0.02/512	0.165867	
0>>>omp_01110ading_83_223/0944_main_1133_000\$h10fm	1		wall_clock	sec	0.000007	0.000007	0.000007	0.000007	0.000000	0.000000	
0>>>omp_offloading_83_223/0944_main_1143_0000000000000000000000000000000000	101		wall_clock	sec	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
	101		wall_clock	Sec	0.000130	0.000001		0.000010	0.000000	0.000002	
0/// Lodule:[Kiy](    0/// Lodule:[Kiy](    0/// configure 92 222700// 72/boundarycondition undateDDdijijiji 1212 confugureform	101		wall_clock	Sec	0.027090	0.000274	0.000127	0.001400	0.000000	0.000202	
0>>>OHp_intitudating_ss_zz3/0744zz4boundatycondition_updaterruiiiiii_izis_cceanoioopatoim	101		wall_clock		0.000020	0.000000	0.000000	0.000000	0.000000	0.000000	
	404		wall_clock		0.012120	0.000120	0.000000	0.001030	0.000000	0.000100	1 100.0 1
JONN MDT Isond	404	1 0	wall_clock		0.000200	0.000000	0.000004	0.0000000	0.000000	0.000004	1 100.0 1
	202	1 0	wall_clock	I Sec	0.002730	0.000007	0.000005	0.000000	0.000000	0.000004	1 100.0
	101	1 0	wall_clock	I Sec	0.0207247	0.000100	0.000000	0.000770	0.000000	0.000107	1 100.0
10>>> omposed flaghing 83 22370944 724boundarycondition undatePPdijijiji 1231 cce\$poloon\$form	101	1 0	wall_clock	Sec	0.00030	0.000270	0.0001/2	0.000000	0.000000	0.000004	1 100.0
0>>> MPILDDownExchange	101	1 0	wall clock	sec	0.020473	0.000203	0.000058	0.000538	0.000000	0.000139	1 100.0
10>>> Stencil	100	1 0	wall clock	sec	41,184269	0.411843	0.405242	0.699454	0.000845	0.029071	100.0
0>>> omp offloading 83 22370944 Z16ghostcell updatePPdiiiiiiiii 1284 cce%noloon\$form	101	1 0	wall clock	sec	0.000030	0.000000	0.000000	0.000001	0.000000	0.000000	100.0
0>>> BoundaryUbdate	100	1 0	wall clock	sec	0.102319	0.001023	0.000251	0.002748	0.000001	0.000757	100.0
0>>> omp offloading 83 22370944 Z16ghostcell updatePPdiiiiiiiii 1309 cce\$noloon\$form	101	0	wall clock	sec	0.000023	0.000000	0.000000	0.000000	0.000000	0.000000	100.0
0>>> omp offloading 83 22370944 main 1161 cce\$noloop\$form	100	0	wall clock	sec	0.000022	0.000000	0.000000	0.000001	0.000000	0.000000	100.0
0>>> GhostCellUpdate	100	0	wall_clock	sec	0.088821	0.000888	0.000478	0.002956	0.000000	0.000395	100.0

5/29/2024

[Public]

### Ver5

Convert indexing from 2D to 1D – a step towards allocating buffers directly on device

5/29/2024

7d22:16:37 + 063 301 923 + X =		00:00:00 000 000 0	00		00:00:01 000 000 000	702ms 319us 897ns	00:00:02 000 000 000		00:00:03 000 000 00 407ms 453us 779ns-	10	-
∧ ./GhostExchange.inst 6197											
GhostExchange.i 6197		MPI_I	hipStreamCrea	ite hipStream	Synchronize	Stencil hipStreamSynchronize	Stencil hipStreamSync	Stencil hipStreamSync	Stencil hipStreamSync	Stencil Stencil hipStreamSync hipStreamSyn	nc
CPU Context Switches (S)	~	15K						I			_
CPU Frequency [0] (S)	$\sim$	4K					Л		1	I	
CPU Kernel Time (S)	~	2									
CPU Memory Usage (S)	$\sim$	2К									F
CPU Page Faults (S)	$\sim$	50K									
CPU Peak Memory (S)	$\sim$	2К									-
CPU User Time (S)	$\sim$	45	,								
CPU Virtual Memory Usage (S)	$\sim$	45K									
GPU Busy [0] (S)	$\sim$	100		-	1		ſ		· · · · · · · · · · · · · · · · · · ·		
GPU Memory Usage [0] (S)	~	2К	-								
GPU Power [0] (S)	$\sim$	150									
GPU Temperature [0] (S)	~	50			l.		~ ~ ~ ~ ~	****			-
HIP Activity Device 4, Queue 0				omp_offl	oading_83	omp_offloading_83_2237	omp_offload	omp_offload	omp_offload	omp_offloadomp_offloa	id

Quickly see durations by selecting kernel and pressing "m"

Extremely helpful to see activity in this duration on CPU, GPU, other HIP streams, etc.

### Ver6

Use explicit data management directives to allocate buffers on device and keep them on device for entire run

AMD @ ORNL

#### Ver6: Adding explicit OpenMP map directives

• Allocate buffers once on device at the beginning:

#pragma omp target enter data map(alloc: xbuf\_left\_send[0:bufcount], xbuf\_rght\_send[0:bufcount])
#pragma omp target enter data map(alloc: xbuf\_rght\_recv[0:bufcount], xbuf\_left\_recv[0:bufcount])
#pragma omp target enter data map(alloc: x[0:totcells], xnew[0:totcells])

• Release buffers at the end:

#pragma omp target exit data map(release: x, xnew)
#pragma omp target exit data map(release: xbuf\_left\_send, xbuf\_rght\_send)
#pragma omp target exit data map(release: xbuf\_rght\_recv, xbuf\_left\_recv)

- Keeping data on HBM improves performance of memory bound kernels on MI250X GPUs
- Managed memory support no longer needed:

#pragma omp requires unified\_shared\_memory
unset HSA\_XNACK

9d16:26:48 + 866 704 114				00:00:01 480 000 000		00:00:01 500 000 000	00:0 520	0:01 000 00	0  2ms 821us 145ns	00:00:01 540 000 00	0	00:0 560	00:01 000 000 5ms 474us 6	09ns	
X															
<ul> <li>./GhostExchange.inst 12468</li> </ul>															
GhostExchange.i 12468									GhostExchange	e.inst					
	hir	BufAlloc pSt	h h M	MPI_Irecv	MPIUpDov	wnExchange MPI_Waitall		St h	GhostCellUpd MPIUpDownExc	ate hange	S Stenc h hipSt	Stenc	Stenc hipSt	Stenc Stenc	Ster hip:
						_		Т	MPI_Waita	I					
CPU Context Switches (S)	≁ <u>4</u> K	(													
CPU Frequency [0] (S)	≁ 4K	(													
CPU Kernel Time (S)	~ 750	0m													
CPU Memory Usage (S)	~ 400	0													
CPU Page Faults (S)	~ 501	K													
CPU Peak Memory (S)	~ 400	0		356.3	56										
CPU User Time (S)	~ 2														
CPU Virtual Memory Usage (S)	~ 451	K								1					
GPU Busy [0] (S)	~ 100	0													
GPU Memory Usage [0] (S)	✓ <sup>15</sup>	K													
GPU Power [0] (S)	~ 200	0													
GPU Temperature [0] (S)	≁ 40														
HIP Activity Device 4, Queue 0			-								om	omp	om	om	0

5/29/2024

9d16:26:48 + 866 704 114		00:00:00 000 000 000	_ 	00:00:00 500 000 000	——————————————————————————————————————	00:00:01 000 000 000	00:00:01 500 000 000	00:00:02 000 000 000
<ul> <li>./GhostExchange.inst 1246</li> </ul>	8							
GhostExchange.i 12468	<b>म</b> GhostE	Exchange.i 12468			BufAlloc hipStreamCreate	GhostExchange.inst	MP	
CPU Context Switches (S)	~	4К						
CPU Frequency [0] (S)	~	4K				J		
CPU Kernel Time (S)	~	.750m				-		
CPU Memory Usage (S)	~	400						
CPU Page Faults (S)	~	.50K						
CPU Peak Memory (S)	~	400		,				
CPU User Time (S)	~	2						
CPU Virtual Memory Usage (S)	$\sim$	45K						
GPU Busy [0] (S)	$\sim$	100						๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛๛
GPU Memory Usage [0] (S)	~	15K	-		/			
GPU Power [0] (S)	n be	consider	ed startup c	ost	/			Data movement over
GPU Temperature [0] (S)	~	40						
HIP Activity Device 4, Queue 0								

5/29/2024

#### Ver6: Wall clock shows shorter durations of kernels

LABEL	COUNT	DEPTH	METRIC	UNITS	SUM	MEAN	MIN	MAX	VAR	STDDEV	% SELF
	34	0	wall_clock	sec	0.000225	0.000007	0.000002	0.000071	0.000000	0.000012	100.0
0>>> hipRuntimeGetVersion	1	0	wall_clock	sec	0.000128	0.000128	0.000128	0.000128	0.000000	0.000000	100.0
0>>> hipDeviceGet	2	0	wall_clock	sec	0.000005	0.000003	0.000002	0.000003	0.000000	0.000001	100.0
0>>> hipGetDeviceCount	1	0	wall_clock	sec	0.000006	0.000006	0.000006	0.000006	0.000000	0.000000	100.0
0>>> GhostExchange.inst	1	0	wall_clock	sec	2.319124	2.319124	2.319124	2.319124	0.000000	0.000000	100.0
0>>> MPI_Init	1	0	wall_clock	sec	0.217335	0.217335	0.217335	0.217335	0.000000	0.000000	100.0
0>>> MPI_Comm_rank	4	0	wall_clock	sec	0.000046	0.000012	0.000001	0.000031	0.000000	0.000013	100.0
0>>> MPI_Comm_size	2	0	wall_clock	sec	0.000014	0.000007	0.000002	0.000012	0.000000	0.00008	100.0
0>>> MPI_Allreduce	1	0	wall_clock	sec	0.000917	0.000917	0.000917	0.000917	0.000000	0.000000	100.0
0>>> BufAlloc	1	0	wall_clock	sec	1.246631	1.246631	1.246631	1.246631	0.000000	0.000000	100.0
0>>> hipDeviceComputeCapability	1	0	wall_clock	sec	0.000002	0.000002	0.000002	0.000002	0.000000	0.000000	100.0
0>>> hipDeviceGetName	1	0	wall_clock	sec	0.000003	0.000003	0.000003	0.000003	0.000000	0.000000	100.0
0>>> hipDeviceGetAttribute	1	0	wall_clock	sec	0.000003	0.000003	0.000003	0.000003	0.000000	0.000000	100.0
0>>> hipGetDeviceProperties	1	0	wall_clock	sec	0.000005	0.000005	0.000005	0.000005	0.000000	0.000000	100.0
0>>> hipGetDevice	550	0	wall_clock	sec	0.001796	0.000003	0.000001	0.000022	0.000000	0.000002	100.0
0>>> hipHostMalloc	1	0	wall_clock	sec	0.000106	0.000106	0.000106	0.000106	0.000000	0.000000	100.0
0>>> hipEventCreate	2	0	wall_clock	sec	0.000006	0.000003	0.000001	0.000005	0.000000	0.000003	100.0
0>>> hipStreamCreate	1	0	wall_clock	sec	1.241892	1.241892	1.241892	1.241892	0.000000	0.000000	100.0
0>>> hipMalloc	7	0	wall_clock	sec	0.000176	0.000025	0.000002	0.000044	0.000000	0.000020	100.0
0>>> hipModuleLoadData	1	0	wall_clock	sec	0.002024	0.002024	0.002024	0.002024	0.000000	0.000000	100.0
0>>> hipPointerGetAttributes	1	0	wall_clock	sec	0.000004	0.000004	0.000004	0.000004	0.000000	0.000000	100.0
0>>> hipModuleGetGlobal	38	0	wall_clock	sec	0.000133	0.000003	0.000001	0.000013	0.000000	0.000002	100.0
0>>> hipModuleGetFunction	7	0	wall_clock	sec	0.000020	0.000003	0.000002	0.000004	0.000000	0.000001	100.0
0>>> hipFuncGetAttribute	14	0	wall_clock	sec	0.000025	0.000002	0.000001	0.000005	0.000000	0.000001	100.0
0>>> hipFuncSetCacheConfig	7	0	wall_clock	sec	0.000011	0.000002	0.000001	0.000003	0.000000	0.000001	100.0
0>>> hipModuleLaunchKernel	506	0	wall_clock	sec	0.003685	0.000007	0.000004	0.000041	0.000000	0.000003	100.0
0>>> hipStreamSynchronize	506	0	wall_clock	sec	0.587353	0.001161	0.000010	0.005770	0.000005	0.002202	100.0
0>>>omp_offloading_83_22370965_main_1142_cce\$noloop\$form	1	0	wall_clock	sec	0.000010	0.000010	0.000010	0.000010	0.000000	0.000000	100.0
0>>>omp_offloading_83_22370965_main_1152_cce\$noloop\$form	1	0	wall_clock	sec	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	100.0
0>>> MPIRequest	101	0	wall_clock	sec	0.000110	0.000001					
0>>> LoadLeftRight	101	0	wall_clock	sec	0.006125	0.000061	N 4				- 61-
0>>>omp_offloading_83_223/0965224boundarycondition_updatePdililili_1227_cce\$noloop\$form	101	0	wall_clock	sec	0.000022	0.000000	. IVIEI	mory r	ouna	kernel	s, iasie
0>>> MPILETTRIGHTEXCHANGE	101	0	Wall_Clock	sec	0.0218//	0.000217		<b>,</b>			e
0>>> MPI_ITECV	404	0	Wall_clock	sec	0.021302	0.000053	l whe	en acc	essinc	i data	trom H
0>>> MPI_ISEND	404	0	Wall_clock	sec	0.003006	0.000007			eeenig	, aata	
0>>> MPI_Waltall	202	0	wall_clock	sec	0.0/1129	0.000352					100.0
0>>> UnpackLettRight	101	0	wall_clock	sec	0.008849	0.000088	0.000057	0.0002//	0.000000	0.000028	100.0
0>>>omp_otrioaaing_83_223/0965224boundarycondition_updatePdillill1245_cce%noloop%form	101	0	wall_clock	sec	0.000030	0.000000	0.000000	0.000000	0.000000	0.000000	100.0
10>>> MP10DD0wILtxChange	101	0	wall clock	Sec	0.0/941/	0.000/86	0.000063	0.0452/9	0.000024	0.004885	100.0
	100	0	wall_clock	sec	0.562033	0.005620	0.002907	0.00581/	0.000000	0.000396	100.0
w>>>ump_uiitoating_83_223/0965216gnostCell_updatePdillillill1_1300_CCe\$n0100p\$form	101	0	wall_clock	sec	0.000032	0.000000	0.000000		0.000000	0.000000	
10>>> Boundaryupdate	100	0	wall_clock	sec	0.028627	0.000286	0.000111	0.000864	0.000000	0.000169	100.0
0>>>ump_uiiiuauing_83_2223/0405210gnostceii_upaaterallillillil_1325_cce\$n0100p\$form	101	0	wall_clock	sec	0.000023	0.000000	0.000000	0.000001	0.000000	0.000000	100.0
0>>>omb_oiiingaruð22753/0400_mgtu_tt/0_cc69uotoob9ioim	100	0	wall_clock	sec			0.000000		0.000000	0.000000	100.0
v>>> GHOSTGETTOBGATE	100	0	wall_clock	sec	0.0/1/08	0.000/1/	0.000290	0.020120	0.000004	0.001010	T00.0

5/29/2024



### **Omnitrace Tips and Status**

#### Tips: Reduce generated output for profiling at scale

• Turn off all options in config file except OMNITRACE\_PROFILE to reduce generated output

OMNITRACE_TRACE	= false
OMNITRACE_PROFILE	= true
OMNITRACE_FLAT_PROFILE	= true
OMNITRACE_USE_ROCTRACER	= false
OMNITRACE_USE_ROCM_SMI	= false
OMNITRACE_USE_MPIP	= true
OMNITRACE_USE_PID	= true
OMNITRACE_USE_ROCPROFILER	= false
OMNITRACE_USE_ROCTX	= false

41

[Public]

#### Tips: If Omnitrace does nothing, check app or environment

- If Omnitrace starts, but does not generate any output files, something prevented the app from running
- To check, unload Omnitrace module, build and run app. Fix errors, then profile with Omnitrace
- If app fails only when being profiled with Omnitrace, try profiling interactively using srun instead of sbatch
  - Conflict due to mismatch in loaded libraries at runtime
- If you use omnitrace-run and it complains saying "Use omnitrace-run", then try running your job
  interactively using srun instead of using sbatch
  - Conflict due to mismatch in loaded libraries at runtime

[Public]

#### Tips: To visualize very large proto files, load into memory first

Linux®

- curl -L0 https://get.perfetto.dev/trace\_processor
- chmod +x ./trace\_processor
- ./trace\_processor -httpd <path to trace file>
- Open up Chrome browser and go to <u>https://ui.perfetto.dev</u>
- When prompted, click on "Yes, use loaded trace"

#### $\mathsf{Windows}^{\mathbb{R}}$

- Open up <u>https://get.perfetto.dev/trace\_processor</u> in a browser to download the Python<sup>™</sup> script
- py trace\_processor --httpd <trace file>
  - You may need to download and install Python on your windows system
- Open up Chrome browser and go to <u>https://ui.perfetto.dev</u>
- When prompted, click on "Yes, use loaded trace"

#### **Research version of Omnitrace is brittle**

- Viewing traces of multiple ranks together was possible using simple concatenation of proto files:
   cat perfetto-trace-0.proto perfetto-trace-1.proto > merged.proto
  - Marging broken new due to change in expected date format in Derfette
  - Merging broken now due to change in expected data format in Perfetto
- Building Omnitrace with Dyninst from source requires GCC, may interfere with CCE
  - On Frontier, omnitrace/1.11.2 is set up to work with CCE and show OpenMP offload and HIP activity
- If you load ROCm, reload Omnitrace as the right build has to be made available
- Python version in your environment matters
- Production version of Omnitrace will be more robust

#### Homework

• See HIP equivalent of Ver1 here:

https://github.com/amd/HPCTrainingExamples/tree/main/MPI-examples/GhostExchange/GhostExchange\_ArrayAssign\_HIP/Ver1

- Use Omnitrace to obtain traces for a 4-rank run
- Progressively port the changes in Ver2 Ver6 in the HIP version using HIP APIs for memory copies, etc.
- Look for memory copy activity and HIP API calls in Omnitrace traces
- Submit a <u>PR</u> with your code or <u>add an issue</u> with any concerns

#### References

- Omnitrace documentation website: <u>https://rocm.github.io/omnitrace/index.html</u>
- Previous talk describing various Omnitrace options: <u>15: GPU Profiling Performance Timelines</u>
- Ghost Exchange OpenMP offload Example suite on github
- ROCm docs: <u>https://rocm.docs.amd.com/en/latest/</u>
- ROCm Blog post: Introduction to profiling tools for AMD hardware

#### Disclaimer

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions, and typographical errors. The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. Any computer system has risks of security vulnerabilities that cannot be completely prevented or mitigated. AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

THIS INFORMATION IS PROVIDED 'AS IS." AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS, OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION. AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY, OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY RELIANCE, DIRECT, INDIRECT, SPECIAL, OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

Third-party content is licensed to you directly by the third party that owns the content and is not licensed to you by AMD. ALL LINKED THIRD-PARTY CONTENT IS PROVIDED "AS IS" WITHOUT A WARRANTY OF ANY KIND. USE OF SUCH THIRD-PARTY CONTENT IS DONE AT YOUR SOLE DISCRETION AND UNDER NO CIRCUMSTANCES WILL AMD BE LIABLE TO YOU FOR ANY THIRD-PARTY CONTENT. YOU ASSUME ALL RISK AND ARE SOLELY RESPONSIBLE FOR ANY DAMAGES THAT MAY ARISE FROM YOUR USE OF THIRD-PARTY CONTENT.

© 2024 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, ROCm, Infinity Fabric, and combinations thereof are trademarks of Advanced Micro Devices, Inc. in the United States and/or other jurisdictions. Other names are for informational purposes only and may be trademarks of their respective owners.

Windows is a registered trademark of Microsoft Corporation in the US and/or other countries.

Git and the Git logo are either registered trademarks or trademarks of Software Freedom Conservancy, Inc., corporate home of the Git Project, in the United States and/or other countries

PCIe® is a registered trademark of PCI-SIG Corporation.

The OpenMP name and the OpenMP logo are registered trademarks of the OpenMP Architecture Review Board

5/29/2024

# AMDJ