## Introduction to the

# Cray Programming Environment

October 10, 2019

Heidi Poxon



a Hewlett Packard Enterprise company

### **Technical Data Rights**



All materials contained in, attached to, or referenced by this document that are marked Cray Confidential or with a similar restrictive legend may not be disclosed in any form without the advance written permission of Cray, a Hewlett Packard Enterprise company. These data are submitted with limited rights under Government Contract No. B626589 and Lease Agreement 4000167127. These data may be reproduced and used by the Government with the express limitation that they will not, without written permission of Cray, be used for purposes of manufacture nor disclosed outside the Government.

This notice shall be marked on any reproduction of these data, in whole or in part.

### **Copyright and Trademark Acknowledgements**



©2016-2019 Cray, a Hewlett Packard Enterprise company. All Rights Reserved.

Portions Copyright Advanced Micro Devices, Inc. ("AMD") Confidential and Proprietary.

The following are trademarks of Cray, and are registered in the United States and other countries: CRAY and design, SONEXION, URIKA, and YARCDATA. The following are trademarks of Cray: APPRENTICE2, CHAPEL, CLUSTER CONNECT, CLUSTERSTOR, CRAYDOC, CRAYPAT, CRAYPORT, DATAWARP, ECOPHLEX, LIBSCI, NODEKARE, and REVEAL. The following system family marks, and associated model number marks, are trademarks of Cray: CS, CX, XC, XE, XK, XMT, and XT. ARM is a registered trademark of ARM Limited (or its subsidiaries) in the EU and/or elsewhere. ThunderX, ThunderX2, and ThunderX3 are trademarks or registered trademarks of Cavium Inc. in the U.S. and other countries. The registered trademark LINUX is used pursuant to a sublicense from LMI, the exclusive licensee of Linus Torvalds, owner of the mark on a worldwide basis. Intel, the Intel logo, Intel Cilk, Intel True Scale Fabric, Intel VTune, Xeon, and Intel Xeon Phi are trademarks or registered trademarks of Intel Corporation in the U.S. and/or other countries. Lustre is a trademark of Xyratex. NVIDIA, Kepler, and CUDA are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and/or other countries.

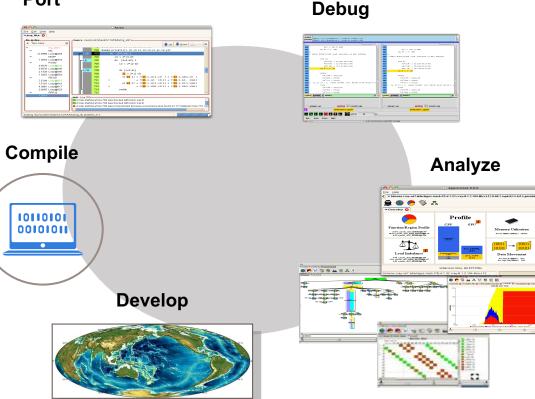
Other trademarks used in this document are the property of their respective owners.

#### FORWARD LOOKING STATEMENTS

This presentation may contain forward-looking statements that involve risks, uncertainties and assumptions. If the risks or uncertainties ever materialize or the assumptions prove incorrect, the results of Hewlett Packard Enterprise Company and its consolidated subsidiaries ("Hewlett Packard Enterprise") may differ materially from those expressed or implied by such forward-looking statements and assumptions. All statements other than statements of historical fact are statements that could be deemed forward-looking statements, including but not limited to any statements regarding the expected benefits and costs of the transaction contemplated by this presentation; the expected timing of the completion of the transaction; the ability of HPE, its subsidiaries and Cray to complete the transaction considering the various conditions to the transaction, some of which are outside the parties' control, including those conditions related to regulatory approvals; projections of revenue, margins, expenses, net earnings, net earnings per share, cash flows, or other financial items; any statements concerning the expected development, performance, market share or competitive performance relating to products or services; any statements regarding current or future macroeconomic Enterprise and its financial performance; any statements of expectation or belief; and any statements of assumptions underlying any of the foregoing. Risks, uncertainties and assumptions include the possibility that expected benefits of the transaction described in this presentation may not materialize as expected; that the transaction may not be timely completed, if at all; that, prior to the completion of the transaction, Cray's business may not perform as expected due to transaction-related uncertainty or other factors; that the parties are unable to successfully implement integration strategies; the need to address the many challenges facing Hewlett Packard Enterprise's businesses; the competitive pressures faced by Hewlett Packard Enterprise's businesses; risks associated with executing Hewlett Packard Enterprise's strategy; the impact of macroeconomic and geopolitical trends and events: the development and transition of new products and services and the enhancement of existing products and services to meet customer needs and respond to emerging technological trends; and other risks that are described in our Fiscal Year 2018 Annual Report on Form 10-K, and that are otherwise described or updated from time to time in Hewlett Packard Enterprise's other filings with the Securities and Exchange Commission, including but not limited to our subsequent Quarterly Reports on Form 10-Q. Hewlett Packard Enterprise assumes no obligation and does not intend to update these forward-looking statements.



Cray PE



- Intuitive behavior and best performance with the least amount of effort
- Mature programming environment to develop, debug, analyze, and optimize applications for production supercomputing
- Complete developer software stack with unique functionality built from close interaction with users



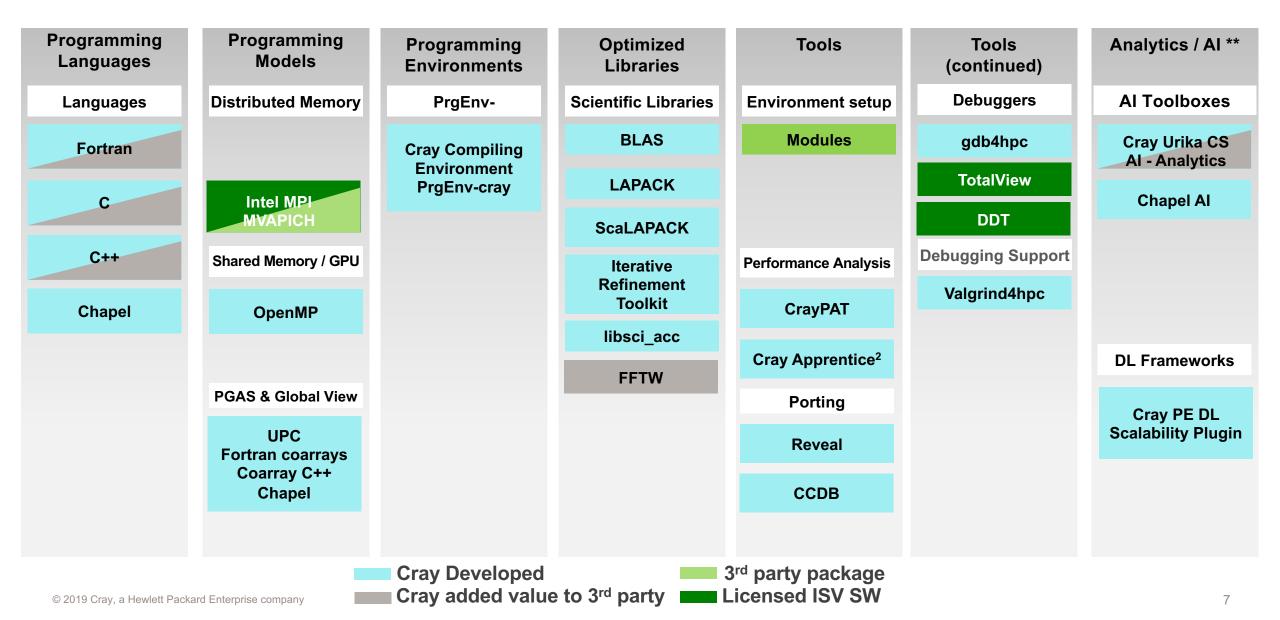
## Cray Developer Environment for Cray XC



Programming Languages	Programming Models	Programming Environments	Optimized Libraries	Tools	Tools (continued)	Analytics / AI **	
Languages	<b>Distributed Memory</b>	PrgEnv-	Scientific Libraries	Environment setup	Debuggers	AI Toolboxes	
Fortran	Cray MPI SHMEM	Cray Compiling Environment	BLAS	Modules	gdb4hpc	Cray Urika XC AI - Analytics	
С		PrgEnv-cray	LAPACK		TotalView	Chapel Al	
C			ScaLAPACK		DDT	onaporra	
C++	Shared Memory / GPU	GNU PrgEnv-gnu	Iterative	Performance Analysis	Debugging Support		
Chapel	OpenMP		Refinement Toolkit	CrayPAT	Valgrind4hpc		
		3 <sup>rd</sup> Party compilers	3 <sup>rd</sup> Party compilers (Intel, Allinea, PGI)	libsci_acc		STAT	
Python	PGAS & Global View	PrgEnv-???	FFTW	Cray Apprentice <sup>2</sup>	Abnormal	DL Frameworks	
R	UPC Fortran coarrays		I/O Libraries	Porting	Termination Processing (ATP)	Cray PE DL	
_	Coarray C++ Chapel		NetCDF	Reveal		Scalability Plugi	
	Global Arrays		HDF5	CCDB			
© 2019 Cray, a Hewlett Pack	ard Enterprise company	Cray Developed		3 <sup>rd</sup> party package Licensed ISV SW		6	

## Cray Developer Environment for Cray CS





## Cray's PE Vision for Accelerated Computing



- Most important hurdle for widespread adoption of accelerated computing in HPC is programming difficulty
  - Need a single programming model that is portable across machine types
    - **Portable** expression of heterogeneity and multi-level parallelism
    - Programming model and optimization should not be significantly difference for "accelerated" nodes and multi-core x86 processors
    - Allow users to maintain a single code base
- Cray's approach to Accelerator Programming is to provide an **ease of use** tightly coupled high level programming environment with compilers, libraries, and tools that can **hide the complexity** of the system
- Ease of use is possible with
  - Compiler making it feasible for users to write applications in Fortran, C, and C++
  - Tools to help users port and optimize for hybrid systems
  - Auto-tuned scientific libraries

## Cray PE Technology Applied to Accelerators

#### • Fortran, C, and C++ compilers

- OpenMP target offload directives to drive compiler optimization
  - Compiler optimizations take advantage of accelerator and multi-core X86 hardware appropriately
- Support for new programming frameworks, such as RAJA and Kokkos
- Advanced users can mix HIP or CUDA functions with compiler-generated accelerator code



- Scientific Libraries tuned to take advantage of multi-core and accelerators appropriately
- Cray Reveal, built upon compiler knowledge of the application
  - Scoping tool to help users port and optimize applications
- Cray Performance tools for whole-program view of performance
- Cray MPI GPU to GPU for direct communications between GPUs on node and off node



🛧 Up 🔻 🛙







# Environment Setup



### **Programmability Focused Environment**



- Modules simplify build environment
  - Complexity of compile and link lines (-h –l –l –L) reduced
- Multiple product versions, compilers, and compiler versions available on system at the same time offers more flexibility and convenience
- Product agnostic drivers (cc, CC, ftn) are used to compile for supported Programming Environments
  - Customer-integrated and Cray libraries share the same driver interface
- Support available to plug 3<sup>rd</sup> party software into Cray software environment (craypkggen)

#### Which Software Versions Are Available?



		<pre>/opt/cray/pe/modulefiles</pre>		
		3 cce/8.7.6	cce/8.7.9	
cce/8.7.1	cce/8.7.	4 cce/8.7.7	cce/8.7.10	
cce/8.7.2	cce/8.7.	5 cce/8.7.8	cce/9.0.0(default)	
user@login:~>		PrgEnv /opt/cray/pe/modulefiles		
PrgEnv-cray/6.		PrgEnv-gnu/6.0.3		
PrgEnv-cray/6.		PrgEnv-gnu/6.0.4		
PrgEnv-cray/6.	0.5( <b>default</b> )	PrgEnv-gnu/6.0.5(default)	<pre>PrgEnv-intel/6.0.5(default)</pre>	
user@login:~>				
		<pre>/opt/cray/pe/modulefiles</pre>	cray-fftw/3.3.8.2	
	6.1	<pre>/opt/cray/pe/modulefiles</pre>		

### **Targeting Processors / Accelerators**



#### user@login:~> module avail craype

craype-accel-host craype-accel-nvidia20 craype-accel-nvidia35 craype-accel-nvidia52 craype-accel-nvidia60 craype-broadwell craype-haswell craype-hugepages1G craype-hugepages2G

craype-hugepages2M craype-hugepages4M craype-hugepages8M craype-hugepages16M craype-hugepages32M craype-hugepages64M craype-hugepages128M craype-hugepages256M craype-hugepages512M

/opt/cray/pe/craype/2.6.1/modulefiles -craype-intel-knc craype-ivybridge craype-mic-knl craype-network-aries craype-network-none craype-sandybridge craype-x86-cascadelake craype-x86-skylake

/opt/cray/pe/modulefiles ------\_\_\_\_\_ craype/2.6.1(default) craype-dl-plugin-py3/19.09.1(default) craype-dl-plugin-py2/19.09.1(default)

### Choosing Different Compilers on Cray XC



- To access a different compiler:
  - Load or swap to the corresponding Programming Environment (PE) module
    - PrgEnv-cray for CCE
    - PrgEnv-intel for Intel
    - PrgEnv-gnu for GNU
  - Once one of these is loaded, you can then select a compiler version
    - CCE: module avail cce
    - GNU: module avail gcc

#### • With PE 19.06 the default linking for all compilers (CCE, Intel, GCC, etc.) is now dynamic

• Static linking still a non-default option, where supported

## Using the Compilers



- Cray Systems come with compiler wrappers to simplify building parallel applications (similar the mpicc/mpif90)
  - Fortran Compiler: ftn
  - C Compiler: cc
  - C++ Compiler: CC
- Using these wrappers ensures that your code is built for the compute nodes and linked against important libraries
  - Cray MPT (MPI, Shmem, etc.)
  - Cray LibSci (BLAS, LAPACK, etc.)
  - ....
- Do not call the Cray, Intel, GNU compilers directly
- Cray Compiler wrappers try to hide the complexities of using the proper header files and libraries

#### **Release Notes Easily Accessible**



user@login:~> module help cce

```
----- Module Specific Help for 'cce/9.0.0' -----
```

The modulefile, cce, defines the system paths and environment variables needed to run the Cray Compile Environment.

Type "module avail cce" to see if other versions of this product are available on this system. Use "module switch" to change versions.

Cray Compiling Environment (CCE) 9.0.0

Release Date:

June 20, 2019

Purpose:

\_\_\_\_\_

- CCE 9.0.0 provides Fortran, C, and C++ compilers for Cray XC systems and Cray CS systems.
  - With this release, the default C and C++ compilers are based on Clang/LLVM
  - The previous C and C++ compilers are also provided and referred to as Cray Classic C and C++
  - See S-5212 Cray Compiling Environment Release Overview (9.0) for additional information

The following key enhancements are included in this CCE release:

- Dynamic linking is the default link mode
  - This change is in craype included with the June PE (19.06) release
  - Applies to all compilers (CCE, Intel, GCC, etc.) when the new craype is loaded
  - Override the default with the -static flag on the command line,
    - or by setting CRAYPE\_LINK\_TYPE=static in the environment
- C++17 support for both Classic and Clang-based C and C++ compilers

-...

# Cray MPI



## Cray MPI Overview



- ANL MPICH3.2 implementation base
- Fully compliant with the MPI-3.1 Specification
- Optimize for scale in the following areas: I/O, collectives, point-to-point, one-sided communication
- Optimization examples
  - SMP-aware collectives
  - High performance single-copy on-node communication via Cray's XPMEM
  - MPI rank reordering
  - Supports GPU to GPU direct communication
  - Ran with 2,001,150 ranks on Trinity in 2016
- Highly tunable through environment variables
  - Defaults should generally be best, but some cases benefit from fine tuning

### Example Environment Variable Control



• By default, MPI uses 2M hugepages for Aries GNI mailboxes internally

- Used even if user didn't select hugepages or if user selected different hugepage size
- User can change the 2M size by setting the MPICH\_GNI\_HUGEPAGE\_SIZE environment variable
- Other buffers that get allocated within MPI (e.g. temporary buffers for collectives) will abide by the hugepage module the user has selected

### MPICH\_GNI\_HUGEPAGE\_SIZE



Specifies the hugepage size in bytes

that will be used for the GNI internal mailbox memory. The default size is 2MB. Jobs that scale to high process counts and have a high connectivity pattern may benefit from using a larger hugepage size for this memory, as this can reduce the number of Aries PTE misses. If setting MPICH\_GNI\_HUGEPAGE\_SIZE to a larger value, you may also want to increase the MPICH GNI MBOXES PER BLOCK value.

The supported values are 2M, 4M, 8M, 16M, 32M, 64M, 128M, 256M, and 512M. For CLE 6.0.UP05 and later, values of 1G and 2G are also supported.

Default: 2M

See the MPI man page for environment variable descriptions

#### Interaction with Cray PE and other MPI Libraries



- Integrated into the Cray Programming Environment
  - Compiler drivers manage compile flags and linking automatically
  - Profiling through Cray performance tools
- MPICH ABI compatibility
  - Programs built with other MPICH ABI compatible vendors like Intel MPI can get native Aries performance without recompiling
  - <u>http://wiki.mpich.org/mpich/index.php/ABI\_Compatibility\_Initiative</u>
  - cray-mpich-abi module exposes this feature

#### Cray Performance Tools



### **Cray Performance Tools**



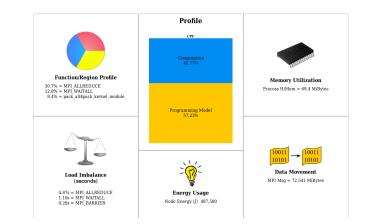
- Reduce the time investment associated with porting and tuning applications on Cray systems
- Analyze whole-program behavior across many nodes to identify critical performance bottlenecks within a program
- Improve profiling experience by using simple and/or advanced interfaces for a wealth of capability that targets analyzing the largest HPC jobs

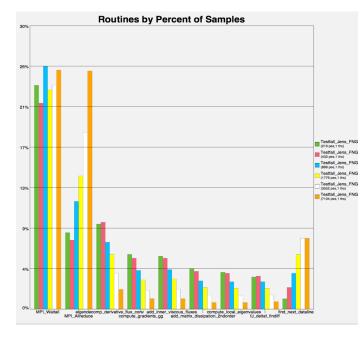
a Hewlett Packard Enterprise compa

### **Cray Tools Provide Various Levels of Detail**



Table 1:	Time and	Bytes Tro	unsferred for	Accelerator	• Regions	
Host	Host	Acc	Acc Copy I	Acc Copy I	Events	Function=[max10]
Time%	Time	Time	In l	Out I		PE=HIDE
1	I		(MiBytes)	(MiBytes)		Thread=HIDE
100.0%	173.71	163.57 I	2,544	2,544	1,071,949	Total
1 50.4%	87.64	83.26			320,970	l cudaLaunchKernel
I 13.7%	23.71	22.55	2,544		285,936	CloverleafCudaChunk::unpackBuffer
11.8%	20.46	19.37		2,544	285,938	CloverleafCudaChunk::packBuffer
10.0%	17.32	15.54			67,975	CloverleafCudaChunk::update_array
I 5.9%	10.29	9.53			23,640	<pre>CloverleafCudaChunk::advec_mom_kernel</pre>
I 3.3%	I 5.65 I	5.26		0.09	29,550	<pre>thrust::tuple&lt;&gt; thrust::cuda_cub::extrema::extrema&lt;&gt;</pre>
I 1.9%	I 3.30 I	3.08			11,820	<pre>void thrust::cuda_cub::core::AgentLauncher&lt;&gt;::launch&lt;&gt; const</pre>
=======						





	amg2006.xf.ap2	
Elle Compare Yiew Help		
About Apprentice2 🧝 arrg2006.xl.ep2 💥		
🏟 🥐 🐳 🌆 🚼 🛛		
🕶 Overview 💥 🖛 Activity 💥 🖛 Call Tree 💥 🖜 Load Bal		
PE	Load Balance: hypre_CSRMatrixMatvec	
E #10		
E #63		
#31 #07		
A07 A25		
A25 A52		
A45		
61		
20 204		
/60		
437		
E #87		
413		
473		
E#10		
E #05		
E #15		
E #17		
E #47 E #83		
E #83 E #85		
765		
E #80 E #79		
E #59		
E 466		
Fato		
490		
E #91		
#92		
E #94		
E #90		
E #95		
9	3.8+40	8.8+40 1.1++04
	Wallclock time: 310.860223s	

#### /\* Turn data recording on for two regions of interest. \*/ PAT\_record(PAT\_STATE\_ON);

PAT\_region\_begin(1, "step 1");

PAT\_region\_end(1);

PAT\_region\_begin(2, "step 2");

...
PAT\_region\_end(2);

/\* Turn data recording off again. \*/
PAT\_record(PAT\_STATE\_OFF);

# Two Modes of Use



• Lite modes: simple interface for convenience



• Advanced interface for in-depth performance investigation and tuning assistance

- Both offer:
  - Whole program analysis across many nodes
  - Indication of causes of problems
  - Ability to easily switch between the two interfaces

### What About Different Compilers?



Cray Performance Tools support the following compilers

- Cray (CCE), Intel, GCC, and Arm Allinea compilers on Cray XC systems
  - AMD compiler support coming next year
- Cray (CCE) compiler on Cray CS systems

### Using the Simple Interface



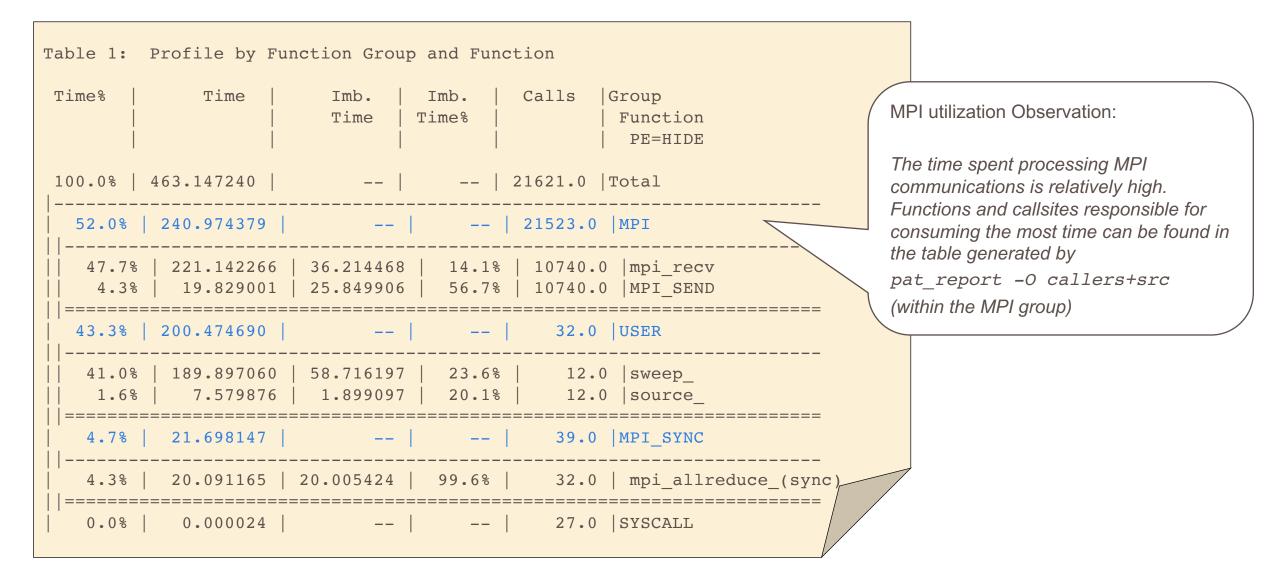
• user@login> module load perftools-lite

• Build program

- Run program
- View report sent to STDOUT (and .rpt file in experiment directory)
  - Example data directory: stencil\_order+49144-225s/

### Find Top Bottlenecks





## Find Any Program Load Imbalance



Table 1:	Profile by H	Function Gro	oup and F	unction		
Time%     	Time     	Imb.   Time	Imb.   Time%   	Calls  0   	Group Function PE=HIDE	Look for function execution imbalance as well as late arrivers
100.0%	1.957703			42,970.8	Iotal	to synchronization points
60.0%	1.174021			3,602.0	USER	
30.8%	0.375117	•	26.0%	1,200.0	function3_  function2_  function1_	
36.0%	======================================			9,613.0	MPI_SYNC	
25.8%		•			<pre> mpi_barrier_(syne)  mpi_init_(sync)</pre>	c)
4.0%	0.078736			29,754.8	  MPI	
2.3%    1.1%		0.003531			MPI_BARRIER  MPI_ISEND ====================================	

#### **Reduce Communication Distance**



#### MPI Grid Detection:

There appears to be point-to-point MPI communication in a 96 X 8 grid pattern. The 52% of the total execution time spent in MPI functions might be reduced with a rank order that maximizes communication between ranks on the same node. The effect of several rank orders is estimated below.

A file named MPICH\_RANK\_ORDER.Grid was generated along with this report and contains usage instructions and the Custom rank order from the following table.

Rank Order	On-Node Bytes/PE	On-Node Bytes/PE% of Total Bytes/PE	MPICH_RANK_REORDER_METHOD
Custom	2.385e+09	95.55%	3
SMP	1.880e+09	75.30%	1
Fold	1.373e+06	0.06%	2
RoundRobin	0.000e+00	0.00%	0

### Use Auto-Generated MPI Rank Order File



<pre># The 'USER_Time_hybrid' rank</pre>	73,395,81,427,57,459,17,419,	53,399,85,431,21,463,61,391,	19,392,75,424,59,456,83,384,	257,345,265,313,281,305,273,
order in this file targets	113,491,49,387,89,451,121,48	109,423,93,455,117,495,125,4	107,416,91,488,115,448,123,4	337,609,369,577,377,617,329,
nodes with multi-core	3	87	80	513,529
<pre># processors, based on Sent</pre>	6,436,102,468,70,404,38,412,	2,530,34,562,66,538,98,522,1	132,401,196,441,164,409,228,	545,297,633,361,625,321,585,
Msg Total Bytes collected	14,444,46,476,110,508,78,500			537,601,289,553,353,593,521,
for:	86.396.30.428.62.460.54.492.	18,514,74,586,58,626,82,546,	188,497	569,561
#			252,505,140,425,212,457,156,	256, 373, 261, 341, 264, 349, 280,
" "	4	10		317,272,381,269,309,285,333,
<pre># Program:</pre>	129,563,193,531,161,571,225,	125 215 167 220 100 247 250		277,365
/WORKSHOP/bh2o-			131,534,195,542,163,566,227,	·
demo/Rank/sweep3d/src/sweep3		247,299		304,360,312,376,293,296,368,
d			187,606	336,344
		1/3,303,139,323,143,333,233,		,
<pre># Ap2 File: sweep3d.gmpi-u.ap2</pre>	619,1//,515,145,5/9,209,54/, 217,611	291,207,275,183,283,151,267, 215,223	251,590,211,630,179,638,139,	
				370,766,306,710,378,742,330, 678,362
# Number PEs: 768	7,405,71,469,39,437,103,413,	133,406,197,438,165,470,229,		,
<pre># Max PEs/Node: 16</pre>			761,660,737,652,705,668,745,	
#	111,397,63,461,55,429,87,421	253,398		290,734,662,686,670,726,702,
# To use this file, make a		157,510,189,462,173,430,205,		694,654
copy named MPICH RANK ORDER,	85		729,732,681,756,721,716,764,	
and set the	134,402,198,434,166,410,230,		640 700	351,286,319,278,342,287,350,
<pre># environment variable</pre>		130,316,260,340,194,372,162,		279,374
MPICH RANK REORDER METHOD to	246,474		760,528,736,536,704,560,744,	
3 prior to	190,498,254,426,142,458,150,	250,300		382,326,303,327,367,366,335,
<pre># executing the program.</pre>	386,182,418,206,490,214,450,	202,364,186,324,154,356,138,	640,600	302,334
" executing the program.	222,482		728,584,680,624,720,512,696,	
#	128,533,192,541,160,565,232,	268,146		669,767,655,743,671,749,695,
0,532,64,564,32,572,96,540,8	525,224,573,240,597,184,557,	4,535,36,543,68,567,100,527,	648,576	679,703
,596,72,524,40,604,24,588	248,605	12,599,44,575,28,559,76,607		677,727,751,693,647,701,717,
104,556,16,628,80,636,56,620	168,589,200,517,152,629,136,	52,591,20,631,60,639,84,519,		687,757,685,733,725,719,735,
,48,516,112,580,88,548,120,6	549,176,637,144,621,208,581,	108,623,92,551,116,583,124,6	730,723	645,759
12	216,613	15	722,731,763,658,642,755,739,	
1,403,65,435,33,411,97,443,9	5,439,37,407,69,447,101,415,	3,440,35,432,67,400,99,408,1	675,707,650,682,715,698,666,	
,467,25,499,105,507,41,475	13,471,45,503,29,479,77,511		690,747	

### View Memory Traffic per NUMA Domain



MPI

Table 3: Memory Bandwidth by Numanode (limited entries shown)

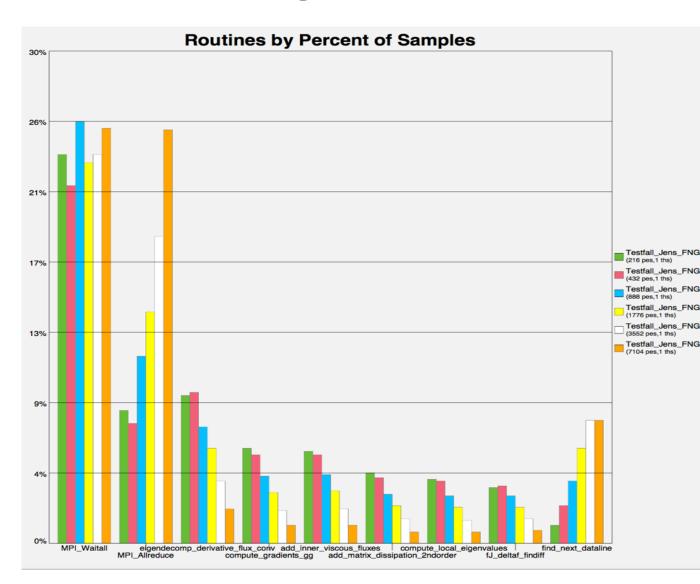
Memory   Traffic   GBytes   	Local   Memory   Traffic   GBytes   	Remote   Memory   Traffic   GBytes   	Memory   Traffic   GBytes   / Sec   	Memory Traffic / Nominal Peak	Numanode   Node Id=[max3,min3]   PE=HIDE   
172.95	171.48	1.48	8.75	11.4%	numanode.0
172.77    172.09    171.20    162.51    162.28    161.75	171.48 170.61 169.93 161.07 160.82 160.29	1.48 1.27 1.43 1.46	8.90   9.02   9.71   8.26   8.22   8.19	11.6% 11.7% 12.6% 10.8% 10.7% 10.7%	nid.68   nid.63   nid.62   nid.71   nid.72   nid.70
168.69	166.81	1.89	8.53	11.1%	numanode.1
168.69    167.74    166.66    161.68    161.60    157.32	166.81   166.03   164.88   160.07   159.99   156.01	1.71	8.67   8.61   8.67   8.17   8.23   8.72	11.3% 11.2% 11.3% 10.6% 10.7% 11.4%	nid.62   nid.63   nid.61   nid.71   nid.70   nid.72

#### MPI + OpenMP

Table 3: Memory Bandwidth by Numanode (limited entries shown)

Memory   Traffic   GBytes   	Local   Memory   Fraffic   ' GBytes   	Remote   Memory   Traffic   GBytes   	GBytes	Memory   Traffic   /   Nominal   Peak	Numanode Node Id=[max3,min3] PE=HIDE Thread=HIDE
184.47	173.59	10.89	15.93	20.7%	numanode.0
183.50    182.61    178.55    178.10    178.08    178.01	173.59 172.40 167.75 168.14 168.07 167.20	9.91   10.21   10.80   9.96   10.01   10.82	15.86   15.77   15.44   15.40   15.40   15.38	20.7%   20.5%   20.1%   20.1%   20.1%   20.0%	nid.63   nid.61   nid.71   nid.62   nid.68   nid.70
60.36	14.73	45.62	6.65	8.7%	numanode.1
60.36 59.88 59.48 58.78 58.67 58.53	14.73 14.33 14.19 13.70 13.87 13.86	45.62   45.55   45.29   45.08   44.81   44.67	6.65   6.60   6.56   6.48   6.47   6.46	8.7%   8.6%   8.5%   8.4%   8.4%   8.4%	nid.63   nid.62   nid.68   nid.70   nid.69   nid.71

#### **Check Scaling**





 pat\_view takes multiple experiment directories as input

 Helpful when assessing performance differences between runs

 Good for function or overall program scaling analysis

### View Application Profile with GPU Information



0 0				🛛 🛛 Аррі	rentice2					
<u>F</u> ile <u>H</u> elp										
🕶 About Appre	entice2 🔞	🕶 himeno_mpi.	ap2 🔞							
۱	🥭 🧇	> ~								
🔻 Overview 🚺	🔯 🔻 Call Ti	ree 🙆 🔻 Tex	1 🚳							
(For per -s per	rcentages r rcent=r[ela	elative to no tivel)	ext leve	lup, spe	cify:					
		unction Group	and Fu	nction					PU waits J executes	
Time%     	Time     	Imb.   : Time   T: 	[mb.   ime%   	Calls     	Group Function PE=HIDE		>			
100.0%   25	5.675919			55473.0	Total					
96.5%   2	24.765662		I	38169.9	USER					
10.2%      6.9%	2.612815 1.761435	0.000594   0.011827	0.0% 0.8%	1003.0   1003.0	jacobiACC_  jacobiACC_  jacobiACC_  jacobiACC_	SYNC_WAIT@li.3 COPY@li.271	31 82			
3.3%	0.850054		I	15066.0						
				3009.0	mpi_waitall_					
0.2% 0.1% 0.0%	0.046805   0.013341   0.000057		   	1005.0 225.1	•					
								Data tra	nsfer to and	
	Obse	rvations and	suggest	ions ===				from	the GPU	
Number of ac	celerators	used: 8 of	8							
		Observations							-	
	III							>		
			Wa	allclock tim	e: 26.309761s					
himeno_mpi.ap	p2 (1,600 eve	ents in 0.444s)							11.	

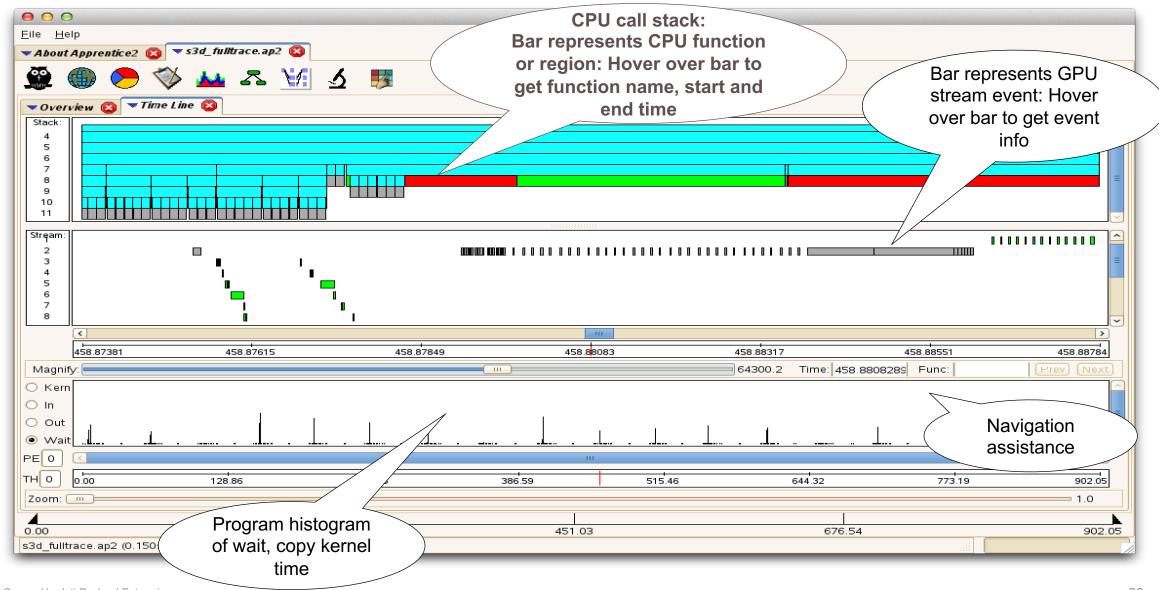
#### Focus on Accelerator Statistics



Table 1:	Time and E	Bytes Transfe	erred for Acc	celerator Re	egions		Host time (wallclock) Time spent on GPU (wallclock)
Host	Host	Acc	Acc Copy	Acc Copy	Calls	Calltree	Data copies to and from GPU
Time%	Time	Time	In	Out		PE=HIDE	Time to copy data
			(MBytes)	(MBytes)			
100.0%	2.750	2.015	2812.760	13.568	103	Total	
	2.750	2.015	2812.760	13.568	103		
     -	 	 	 	 	 		ACC_DATA_REGION@li.104
							.ACC_COPY@li.104
3   22.1	.8 0.60	0.088	3   12.304	1   12.30	04   3	36  streaming_	
4    20.	6% 0.5	66 0.04	46   12.30	04   12.3	304	27  streaming	exchange
5							g_exchangeACC_DATA_REGION@li.526
		ACC_DATA_REG	 GION@li.526(@	 exclusive)		1	
4    1.	6%   0.0	0.043	12   -			9  streaming_	ACC_DATA_REGION@li.907
5    1.	.1%   0.0	0.03	31   -			4   streaming	gACC_REGION@li.909

### Analyze CPU and GPU Overlap

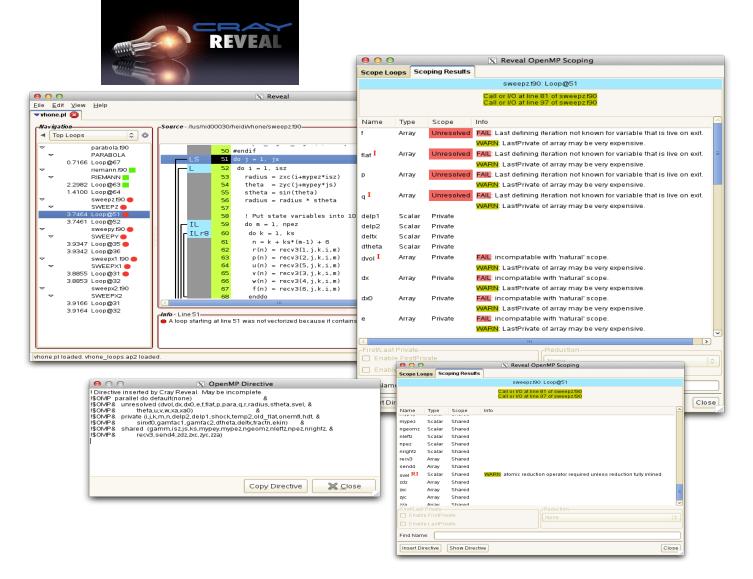




## Adding OpenMP with Cray Reveal



#### Reveal





- Reduce effort associated with adding OpenMP to MPI programs
- Get insight into optimizations performed by the Cray compiler
- Add OpenMP as a first step to parallelize loops that will target GPUs

## The Problem – How Do I Parallelize This Loop?



- How do I know this is a good loop to parallelize?
- What prevents me from parallelizing this loop?
- Can I get help building a directive?

```
subroutine sweepz
do j = 1, js
 do i = 1, isz
   radius = zxc(i+mypez*isz)
   theta = zyc(j+mypey*js)
   do m = 1, npez
    do k = 1, ks
     n = k + ks*(m-1) + 6
     r(n) = recv3(1,j,k,i,m)
     p(n) = recv3(2,j,k,i,m)
     u(n) = recv3(5,j,k,i,m)
     v(n) = recv3(3,j,k,i,m)
     w(n) = recv3(4,j,k,i,m)
     f(n) = recv3(6,j,k,i,m)
    enddo
   enddo
   call ppmlr
   do k = 1, kmax
     \mathbf{n} = \mathbf{k} + \mathbf{6}
     xa(n) = zza(k)
     dx(n) = zdz(k)
     xaO(n) = zza(k)
     dxO(n) = zdz(k)
     e(n) = p(n)/(r(n)*gamm)+0.5 \&
        *(u(n) **2+v(n) **2+w(n) **2)
```

enddo

call ppmlr

enddo enddo

```
subroutine ppmlr
```

```
call boundary
call flatten
call paraset(nmin-4, nmax+5, para, dx, xa)
```

```
call parabola(nmin-4,nmax+4,para,p,dp,p6,pl,flat)
call parabola(nmin-4,nmax+4, para,r,dr,r6,rl,flat)
call parabola(nmin-4,nmax+4,para,u,du,u6,u1,flat)
```

```
call remap < contains more calls
```

```
call volume(nmin,nmax,ngeom,radius,xa,dx,dvol)
```

```
call remap ← contains more calls
```

return end

#### Loop Work Estimates



Gather loop statistics using the Cray performance tools and the Cray Compiling Environment (CCE) to determine which loops have the most work

- Helps identify high-level serial loops to parallelize
  - Based on runtime analysis, approximates how much work exists within a loop

Inclusive and Exclusive Time in Loops

Time

Loop Hit

Loop Incl

© 2019 Cray, a Hewlett Packard Enterprise company

Loop

Incl

Time%

99.4%

98.7%

26.5%

24.6%

24.28

22.5%

22.5% 22.5%

18.9%

17.1% 15.7%

5.0%

#### Create Loop Profile to Find Tuning Candidates

Loop

Trips

Loop

Trips

Use with compiler listing to understand compiler generated optimizations

Function=/.LOOP[.]

PE=HIDE

		Avg	Min	Max	
333.895923	1	500.0	500	500	les3dLOOP.3.li.216
331.721721	500	2.0	2	2	les3dLOOP.4.li.272
89.032566	1,000	96.0	96	96	fluxkLOOP.1.li.28
82.681435	96,000	97.0	97	97	fluxkLOOP.2.li.29
81.356609	1,000	96.0	96	96	fluxjLOOP.1.li.28
75.770180	96,000	97.0	97	97	fluxjLOOP.2.li.29
75.458432	1,000	96.0	96	96	fluxiLOOP.1.li.21
75.453469	96,000	96.0	96	96	fluxiLOOP.2.li.22
63.574836	9,312,000	96.0	96	96	visckLOOP.1.li.344
57.529187	9,312,000	96.0	96	96	viscjLOOP.1.li.340
52.794857	9,216,000	97.0	97	97	visciLOOP.1.li.782
16.924522	1,000	99.0	99	99	extrapiLOOP.1.li.128

Loop

Trips



## View Source and Optimization Information



000	🔀 Reveal	
<u>F</u> ile <u>E</u> dit <u>V</u> iew <u>H</u> elp		Reveal - Loopmark Legend (on g _ 🗆 🖂
vhone.pl 🐰		▶ A Pattern Matched ▼ C Collapsed
Mavigation         ▲       Loop Performance         ▶       4.0423       SWEEPX2@32 ★         ▶       3.8576       SWEEPZ@51 ★         ▶       3.8573       SWEEPZ@52 ★         ▶       3.2068       RIEMANN@63 ★         ▶       1.2299       RIEMANN@64         ▼       0.8068       PARABOLA@67         0.0146       Instance #1       0.0163         0.0163       Instance #5       0.0163         0.0163       Instance #7       ▼         ✓	Source - /home/users/heidi/reveal/parabola.f90	A loop nest has been collapsed into one loop  D Deleted  G Cloned  G Accelerated  I Inlined  L Loop  M Multithreaded  A loop or block of code is multi-threaded  R Region  S Scoping Analysis  V Vectorized  A Atomic Memory Operation  B Blocked  C Conditional and/or Computed  I Fused  P Partial  F Cunrolled  S Shortloop  A loop was converted to a single vector iteration  W Unwound  C Market Construction  A loop construction  C Market C
whone.pl loaded. vhone_loops.ap2 loa		

#### **Review Scoping Results**



<b>e e</b>	Loops with scoping
<u>F</u> ile <u>E</u> dit ⊻iew <u>H</u> elp	information are
Navigation Loop Performance 🛛 🗸 🔅	flagged. Red needs
▶ 4.0778 SWEEPY@35  ★	
▶ 4.0773 SWEEPY@36	user assistance
▶ 4.0529 SWEEPX1@31 \varTheta 📩	<b>52</b> do i = 1, isz
▶ 4.0526 SWEEPX1@32 ●🔨	53 radius = zxc(i+mypez*isz)
▶ 4.0425 SWEEPX2@31	$\frac{33}{54}  \text{theta} = zyc(j+mypey*js)$
▶ 4.0423 SWEEPX2@32	55 stheta = sin(theta)
▶ 3.8576 SWEEPZ@51 🛑☆	56 radius = radius * stheta
▶ 3.8573 SWEEPZ@52	57
▶ 2.2068 RIEMANN@63 = +	58 ! Put state variables into 1D arrays, padding with 6 ghost zones
▶ 1.2299 RIEMANN@64	<b>FS</b> 59 dom = 1, npez
▶ 0.8068 PARABOLA@67	r = Fr8 = 60 do k = 1, ks
▶ 0.5429 PARABOLA@44	61   n = k + ks*(m-1) + 6
▶ 0.5331 PARABOLA@75	$\frac{62}{62} r(n) = recv3(1, j, k, i, m)$
▶ 0.4244 REMAP @83 <mark></mark>	$\frac{63}{63} p(n) = recv3(2, j, k, i, m)$
▶ 0.3341 PARABOLA@30	64   u(n) = recv3(5, j, k, i, m)
0.2966 PARABOLA@84	65 v(n) = recv3(3, j, k, j, m)
▶ 0.2915 PARABOLA@53	$\frac{66}{66} = w(n) = recv3(4, i, k, i, m)$
0.2287 RIEMANN@44	
▶ 0.2028 PARABOLA@36	-Info - Line 51
▶ 0.2009 PARABOLA@117	A loop starting at line 51 was scoped with errors. See Scoping Tool for more information.
▶ 0.1858 PARABOLA@24	"ppmlr" (called from "sweepz") was not inlined because I/O was detected in "volume".
▶ 0.1847 SWEEPY@86 ★	ppmlr" (called from "sweepz") was not inlined because the enclosing loop body did not completely flatten.
▶ 0.1771 STATES@64	A loop starting at line 105 is flat (contains no external calls).
▶ 0.1723 EVOLVE@70 ■★	A loop starting at line 105 was not vectorized because it does not map well onto the target architecture.
▶ 0.1638 REMAP@111 ■★	A loop starting at line 105 was not vectorized because it does not map wen onto the target architecture.
▶ 0.1619 PARABOLA@129	A loop starting at line 103 was unrolled a times. A loop starting at line 51 was not vectorized because it contains a call to subroutine "ppmlr" on line 81.
▶ 0.1070 PARABOLA@139	A loop starting at line 51 was not vectorized because it contains a call to subroutine "ppmlr" on line 81.
▶ 0.0938 SWEEPZ@120	
▶ 0.0936 SWEEPZ@121	A loop starting at line 59 is flat (contains no external calls).
▶ 0.0930 SWEEPZ@122	A loop starting at line 59 was not vectorized because a better candidate was found at line 60.
0.0925 SWEEPX1@59	A loop starting at line 60 is flat (contains no external calls).
▶ 0.0901 SWEEPZ@22	A loop starting at line 60 was not vectorized because it does not map well onto the target architecture.
▶ 0.0898 SWEEPZ@23 <mark>▲</mark> ★	A loop starting at line 60 was unrolled 8 times.
	A loop starting at line 71 is flat (contains no external calls).
▶ 0.0892 STATES@50	
<ul> <li>▶ 0.0892 STATES@50</li> <li>▶ 0.0880 SWFFP7@105</li></ul>	A loop starting at line 71 was vectorized.

Scope Loops	Scoping R	esults	
		Coll or VO o	sweepz.f90: Loop@51
		4: /hon	it line 81 of sweepz.190 🔶 📩 ne/users/heidi/reveal/volume.190:34
			ne/users/heidi/reveal/evolve.f90:21 ne/users/heidi/reveal/ppmlr.f90:73
		1: /hon	ne/users/heidi/reveal/sweepz.190:81 It line 81 of sweepz.190
			nolucore/hoidi/royool/volumo f00:25
Name	Туре	Scope	Info
wl@remap_ I	Scalar	Unresolved	FAIL: Possible recurrence involving this object.
			FAIL: Possible resolvable recurrence involving this object.
xa	Array	Unresolved	FAIL: Possible recurrence involving this object.
			FAIL: Possible resolvable recurrence involving this object.
			WARN: LastPrivate of array may be very expensive.
xa0	Array	Unresolved	FAIL: Possible recurrence involving this object.
			FAIL: Possible resolvable recurrence involving this object.
			WARN: LastPrivate of array may be very expensive.
i	Scalar	Private	N
j	Scalar	Private	
k	Scalar	Private	
m	Scalar	Private	
n	Scalar	Private	
stheta	Scalar	Private	
theta	Scalar	Private	
gamm	Scalar	Shared	
isz	Scalar	Shared	Parallelization
js	Scalar	Shared	
ks	Scalar	Shared	/ inhibitor messages
mypey	Scalar	Shared	
First/Last Privat			are provided to
Enable First			assist user with
Enable Lasti	Private		
_			analysis 🦯
Find Name:			
Insert Directive	Show	Directive	Close

#### Review Scoping Results (continued)

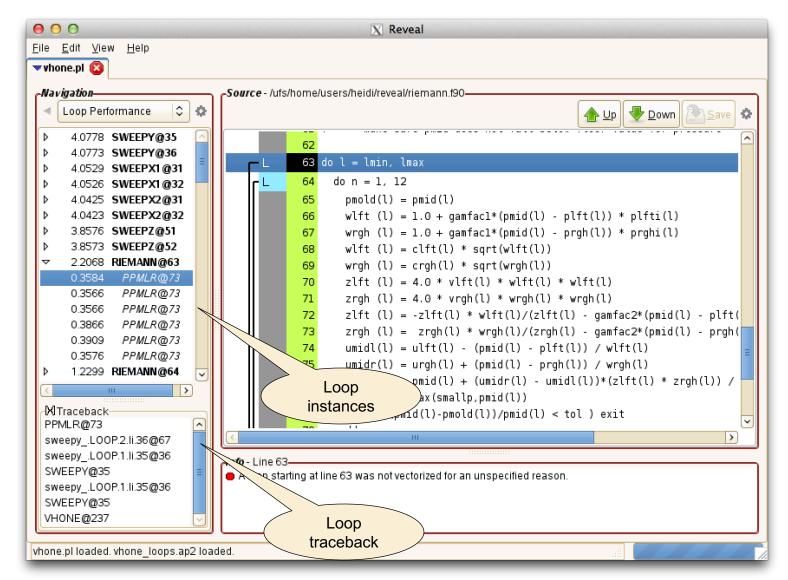


• • •			X Reveal OpenMP Scoping		
Scope L	oops Sco	ping Resul	s Footnote		
			m_mat_an.c: Loop @39		
Name	Туре	Scope	Info		
a0i	Scalar	Private			
a0r	Scalar	Private		I	
a1i	Scalar	Private		I	
a1r	Scalar	Private		I	Scope Loops Scoping Results
a2i	Scalar	Private		I	
a2r	Scalar	Private		I	
bOi	Scalar	Private		I	
bOr	Scalar	Private		I	
b1i	Scalar	Private		I	
b1r	Scalar	Private		I	
b2i	Scalar	Private		I	Assume no overlap b
b2r	Scalar	Private		I	Assume no over cap b
j	Scalar	Private			
a	Scalar	Shared	WARN: Assuming no overlap with other objects.		
			INFO: additional detail.	I	
b	Scalar	Shared	WARN: Assuming no overlap with other objects.	I	
			INFO: additional detail.	I	-
с	Scalar	Shared	WARN: Assuming no overlap with other objects.	I	
			INFO: additional detail.	I	
				I	
First/Las			Reduction		
	le FirstPriv		None	<b>▼</b>	
🗌 Enab	le LastPriv	ate			
Find Nar	ne:				
	,				
Insert D	irective	Show Dire	tive	Close	

🛑 😑 🔵 Scope Loops	Scoping Results	Reveal OpenMP Scoping  Footnote	
		Scoping Footnote	
Assur	ne no overlap b	etween lattice[*].mom[*] and tempmom[*][*]	 ▼
			Close

## View Loops through Call Chain





#### **Generate OpenMP Directives**

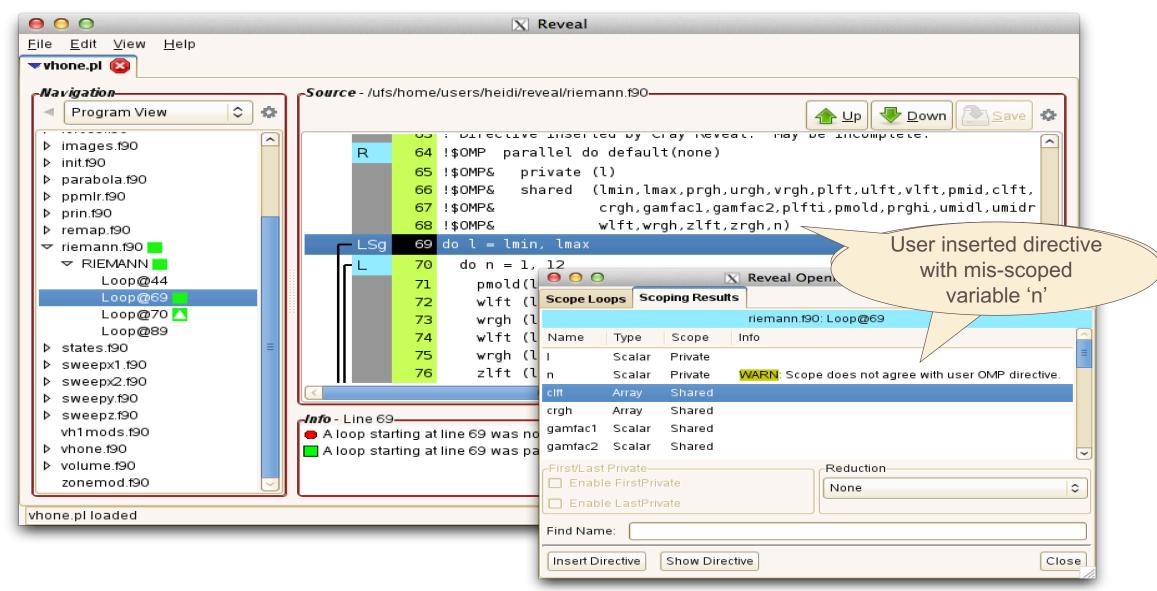


```
! Directive inserted by Cray Reveal. May be incomplete.
!$OMP parallel do default(none)
                                                                       &
!$OMP& unresolved (dvol,dx,dx0,e,f,flat,p,para,q,r,radius,svel,u,v,w,
                                                                       &
!$OMP&
               xa,xa0)
!$OMP& private (i,j,k,m,n,$$ n,delp2,delp1,shock,temp2,old flat,
                                                                       &
               onemfl,hdt,sinxf0,gamfac1,gamfac2,dtheta,deltx,fractn, &
!$OMP&
!$OMP&
               ekin)
!$OMP& shared (gamm,isy,js,ks,mypey,ndim,ngeomy,nlefty,npey,nrighty, &
!$OMP&
               recv1,send2,zdy,zxc,zya)
do k = 1. ks
do i = 1, isy
 radius = zxc(i+mypey*isy)
 ! Put state variables into 1D arrays, padding with 6 ghost zones
 do m = 1, npey
  do j = 1, js
  n = j + js^{*}(m-1) + 6
  r(n) = recv1(1,k,j,i,m)
  p(n) = recv1(2,k,j,i,m)
  u(n) = recv1(4,k,j,i,m)
  v(n) = recv1(5,k,j,i,m)
  w(n) = recv1(3,k,j,i,m)
  f(n) = recv1(6,k,j,i,m)
  enddo
 enddo
 do j = 1, jmax
  n=j+6
```

Reveal generates OpenMP directive with illegal clause marking variables that need addressing

#### Validate User Inserted Directives





#### Look For Vectorization Opportunities



• • • X vhone.pl <u>File E</u>dit <u>V</u>iew <u>H</u>elp Navigation-Source - /home/users/heidi/reveal/riemann.f90 Up Down Save 🔅 Compiler Messages • -62 ٠ 🛨 All Not Vectorized  $\hat{\mathbf{C}}$ 63 do l = lmin. lmax - FS iiiie 125 . 64 ▽ images.f90 Choose "Compiler 65 pmold(l) = pmid(l)line 149 66 wlft (l) = 1.0 + qamfacl\*(pmid(l) - plft(l)) \* plfti(l)⊽ init.f90 Messages" view to wrgh (l) = 1.0 + gamfacl\*(pmid(l) - prgh(l)) \* prghi(l)67 line 113 (0.000 sec) 68 wlft (l) = clft(l) \* sqrt(wlft(l))access message line 114 (0.000 sec) 69 wrqh (l) = crqh(l) \* sqrt(wrqh(l)) line 153 (0.000 sec) filtering, then select 70 zlft(l) = 4.0 \* vlft(l) \* wlft(l) \* wlft(l)line 154 (0.000 sec) 71 zrgh(l) = 4.0 \* vrgh(l) \* wrgh(l) \* wrgh(l)line 139 desired type of 72 zlft(l) = -zlft(l) \* wlft(l)/(zlft(l) - qamfac2\*(pmid(l) - plft(l))✓ prin.f90 line 125 (0.006 sec) 73 zrgh (l) = zrgh(l) \* wrgh(l)/(zrgh(l) - gamfac2\*(pmid(l) - prgh(l message line 42 (0.000 sec) 74 umidl(l) = ulft(l) - (pmid(l) - plft(l)) / wlft(l) line 43 (0.000 sec) 75 umidr(l) = urgh(l) + (pmid(l) - prgh(l)) / wrgh(l)line 127 (0.000 sec) pmid (l) = pmid(l) + (umidr(l) - umidl(l))\*(zlft(l) \* zrqh(l)) /76 line 128 (0.000 sec) 77 pmid (l) = max(smallp,pmid(l)) line 129 (0.000 sec) if (abs(pmid(l)-pmold(l))/pmid(l) < tol ) exit 78 line 104 (0.000 sec) 79 enddo Þ line 63 (0.387 sec) line 64 (0.224 sec) -Info - Line 64-A loop starting at line 64 is flat (contains no external calls). line 31 (4.053 sec) A loop starting at line 64 was not vectorized because a recurrence was found on "pmid" at line 77. line 32 (4.053 sec) line 59 (0.093 sec) Þ vhone.pl loaded. vhone loops.ap2 loaded.

# Cray Debugging Tools



## Scalable Debugging on Cray Systems



- Cray's focus is to build tools around traditional debuggers with innovative techniques for productivity and scalability
- Scalable Solutions based on MRNet from University of Wisconsin Wisconsin



- STAT Stack Trace Analysis Tool
  - Scalable generation of a single merged stack backtrace for the application
  - GUI based tool (stat-gui/stat-view) along with cli tools (stat-cl)
  - Gain insight into application behavior at a function level



- ATP Abnormal Termination Processing
  - Scalable core file generation and analysis when application crashes
  - Generates a merged stack backtrace akin to stat
  - Selection algorithm to dump unique core files

• adb4hpc



- Conventional CLI based interactive parallel debugger
- Look and feel of gdb syntax is inspired by gdb!
- Debug your application at scale
- CCDB Comparative debugging



• A data-centric paradigm

HE UNIVERSITY

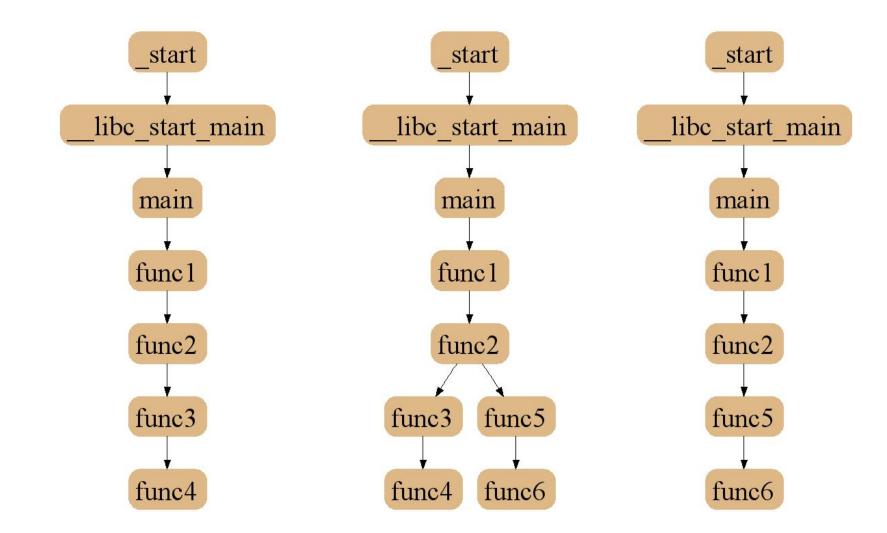
- Compare two applications side-by-side
  - Focus on the data not state and internal operations
- GUI tool that interacts with gdb4hpc



- Valgring4hpc
  - Parallel valgrind based debugging tool (memcheck)
  - Aids in detection of memory leaks and errors in parallel applications
  - Aggregates like errors across PEs/threads

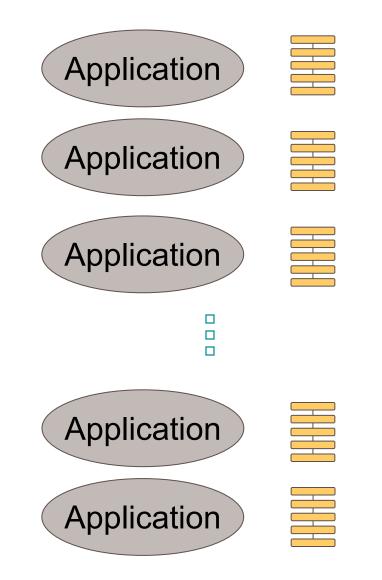
#### Stack Trace Merge Example

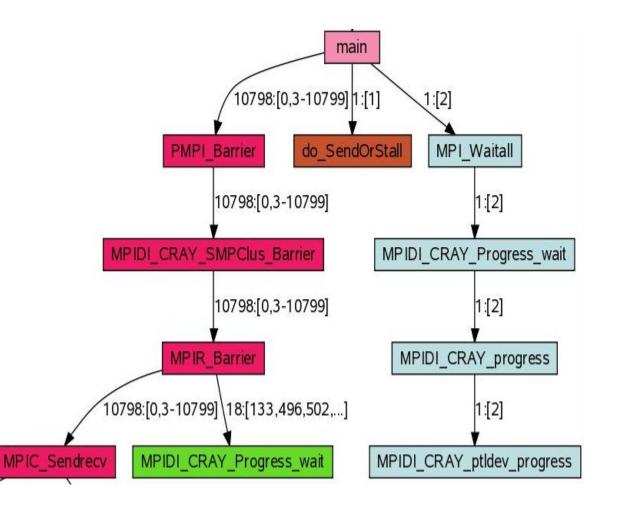




## 2D-Trace/Space Analysis









- ATP signal handler runs within an application to catch fatal errors
  - It handles the following signals:
    - SIGQUIT, SIGILL, SIGTRAP, SIGABRT, SIGFPE, SIGBUS, SIGSEGV, SIGSYS, SIGXCPU, SIGXFSZ
    - Setting the environment variables MPICH\_ABORT\_ON\_ERROR and SHMEM\_ABORT\_ON\_ERROR will cause a signal to be thrown and captured for MPI and SHMEM fatal errors
- ATP daemon running on the compute node captures signals, starts termination processing
  - All application processes are notified
  - Generates a stacktrace
  - Creates a single merged stack trace file
- The stack trace file is viewed with the stat-view tool

## ATP Can Hold Dying Application



- ATP is able to hold a dying application in stasis in order to allow the user to attach to it with a debugger
  - To do so, set the ATP\_HOLD\_TIME environment variable to the number of minutes desired
- Once attached, the debugging session can last as long as the batch system allows
  - Which in turn depends on the compute node resources you requested when you began your session
  - Use ATP\_HOLD\_TIME to define the time you need to attach to the application, not the total time needed for the debugging session
- If ATP\_HOLD\_TIME is set, core dumping is disabled

## gdb4hpc

- Why is a tool like this necessary?
  - Code developers are already familiar with gdb
  - Commercial debuggers are GUI based
  - Using gdb on each individual PE/thread is painful
    - How to attach gdb?
    - How to communicate with gdb?
  - Focus on delivering a tool with scalability in mind
- Traditional parallel debugger
  - Compilers: CCE, GNU, Intel, clang/flang
  - Languages: C, C++, Fortran, UPC
  - Programming models: MPI, SHMEM, OpenMP, pthreads
  - Provides a command set similar to gdb
- Built on top of gdb
  - Have a gdb instance for each PE/thread
  - Glued together with a communication tree
  - gdb was modified to support Fortran and UPC







## Using gdb4hpc



- user@login> module load gdb4hpc
- user@login> man gdb4hpc

```
user@login:~> gdb4hpc
gdb4hpc 3.0 - Cray Line Mode Parallel Debugger
With Cray Comparative Debugging Technology.
Copyright 2007-2018 Cray Inc. All Rights Reserved.
Copyright 1996-2016 University of Queensland. All Rights Reserved.
Type "help" for a list of commands.
Type "help" for a list of commands.
Type "help <cmd>" for detailed help about a command.
dbg all>
```

• Recommend compiling with -g -00 for best experience

### gdb4hpc Example



#### \$ gdb4hpc

•••

```
dbg all> break jacobi
App1{0..127}: Breakpoint 1: file himeno.f, line 209.
dbg all> c
App1{0..127}: Breakpoint 1, jacobi at himeno.f:209
dbg all> 1
App1{0..127}: 209
                              subroutine jacobi(nn,gosa)
App1{0..127}: 210
                                                         App1{0..127}: 211
                              IMPLICIT NONE
App1{0..127}: 212
                        С
App1{0..127}: 213
                              include 'mpif.h'
                              include 'param.h'
App1{0..127}: 214
App1{0..127}: 215
                        С
App1{0..127}: 216
                              integer :: nn,i,j,k,loop,ierr
dbg all> backtrace
App1{0..127}: #0 0x0000000000000000 in jacobi at himeno.f:209
App1{0..127}: #1 0x000000004012de in himenobmtxp at himeno.f:91
dbg all> print npe
App1{0..127}: 128
dbg all> p jmax
App1\{0...3, 12...19, 28...35, 44...51, 60...67, 76...83, 92...99, 108...115, 124...127\}: 129
App1{4..11,20..27,36..43,52..59,68..75,84..91,100..107,116..123}: 130
```

## **Comparative Debugger**





THE UNIVERSITY OF QUEENSLAND

- What is comparative debugging?
  - Data centric approach instead of the traditional control-centric paradigm

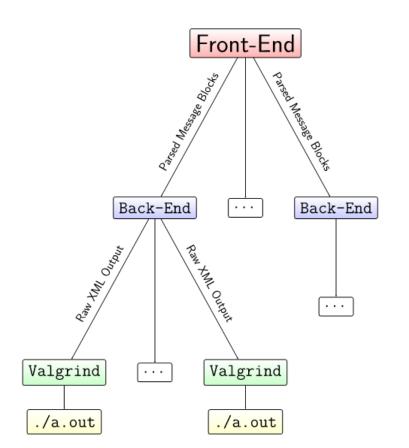
Comp

- Two applications, same data simultaneous execution of both
- Key idea: The data should match
- Quickly isolate deviating variables
- Comparative debugging tool
  - NOT a traditional debugger!
  - Assists with comparative debugging
    - Focus on data not state and internal operations
    - Creates automatic comparisons
    - Based on symbol name and type
    - Allows user to create own comparisons
    - Error and warning epsilon tolerance
    - Scalable
- How does this help me?
  - Algorithm re-writes
  - Language ports
  - Different libraries/compilers
  - New architectures

	•••				X Cray	Compara	ative	Debug	ger (CCDB)				
	<u>F</u> ile \	iew <u>T</u> ools	<u>H</u> elp			1							
				Focus:	all	- all							
	?		Applicat	<mark>ion-0 Status</mark>			_†				A	Application-1 Status	
		{015}	Stop	oed 💿 driver.f:10	5			\$ ¢	opp1{015}			Stopped 💽 sweep.f:231	
						(m)			( <b>-</b> (	_			
		reakpt + driver.				<u> </u>	=	<u> </u>	Breakpt	+ swee	_		<u>± ×</u>
	94:	nm_jbc = 1 endif				ľ		220:				= 1 + (kk-1) *mk	
	95: 96:	enair if (kbc.ne.0	) then			- 1		221: 222:				= min (k0+mk-1,kt) = k1 - k0 + 1	- 1
	97:	it_kbc = i				- 1		223:		els			- 1
	98:	jt_kbc = j				- 1		224:				= kt - (kk-1) *nk	- 1
	99:	nm_kbc = m	m			- 1		225:			k1	= max (k0-mk+1,1)	- 1
	100:	else				- 1		226:			nk	= k0 - k1 + 1	- 1
	101:	it_kbc = 1				- 1		227:		end	lif		-
	102:	jt_kbc = 1						228:					
	103:	nm_kbc = 1	L									<pre>* instead of *mk* if all phi{i,;</pre>	
	104: 105:	endif				l l	-11	230: 1	were dim			vith mmi as the second dimensior jt*mk*mmi	1
	105.	if (myid .eq	. 1) then			- 1		232:				it*mk*mmi	- 1
	107:			Method 5 -',		- 1		233:					- 1
	108:	&		d Wavefront with	h Line-H	ecurs:		<b>234:</b> c	I-inflow	s for b	100	ck (i=i0 boundary)	- 1
)			X CC	DB Comparison									- 1
	A	pplication-0				Appl	licati	on-1			_	v_rcv .ne. 0) then	
App0{015}	٢	sweep.f:160		App1{015}		<b>()</b> sv	weep	.f:160		Up Dow		<pre>ll rcv_real(ew_rcv, phiib, nib, (i2.lt.0 .or. ibc.eq.0) then</pre>	ew_t
Name		Туре		Results	Туре	App-0		op-1 C	Dp Eps			(12.10.01. Ibc.eq.0) then	
Express		Type		nesuns	Template			comp	p Ebs				
km		INTEGER*4	Click to see	e results				-	= e				4
nio		INTEGER*4	Click to see	e results				1	= e				
ık		INTEGER*4	Click to see	e results				-	= e				
i		REAL'8	Click to see					-	= e				
j		REAL'8	Click to se	e results				Ē	= e				
veta		REAL*8 (6)	Click to see	e results		None	N	one =	= e				
/mu		REAL*8 (6)	Click to see	e results		None	N	one =	= е				
											H		
are Add Com	parison Resu	lt Filter: 🔷 all 🖪	🕨 fail 🔶 v	varn 🔷 pass 🔷 to	do					Close	,		
						_	_	_	_	_	_//		

#### Valgrind4hpc





- Scalable Valgrind tool for analyzing parallel applications
- Focus: parallelizing memcheck and helgrind tools
- Memory leak/error checking, thread safety checking, bounds checking
- Performs analysis at every application rank and removes duplicates matching output blocks to provide concise summaries of program behavior
- Available via module
  - user@login> module load valgrind4hpc

## Cray PE DL Scalability Plugin



## The Cray PE DL Scalability Plugin Overview



- Goal was to design a solution for scaling TensorFlow (specifically synchronous SGD) to significantly larger node counts than existing methods allowed
  - Should require minimal changes to user training scripts and provide a more friendly user experience
- Achieve the best possible TensorFlow performance on Cray Systems
- Maintain accuracy for a given number of steps and hyper-parameter setup allowing for significantly reduced time-to-accuracy through scaling
- Can run on a range of x86-64 CPUs and currently supports Nvidia GPUs
- Ideally have a **portable solution** that would work with other deep learning frameworks

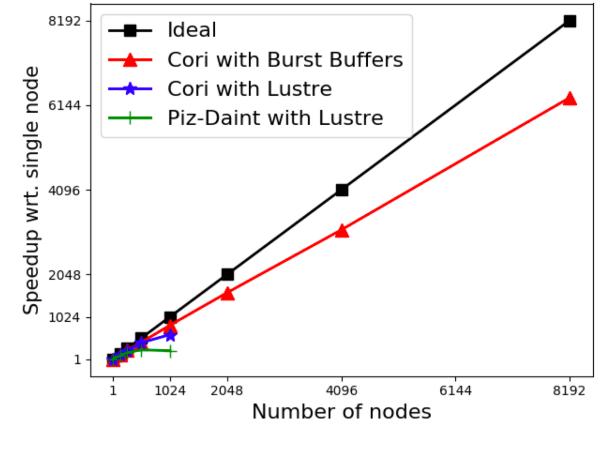
## Using the Plugin



- Available on Cray XC and CS systems
  - user@login> module load craype-dl-plugin-py3
- Plugin has both a C and Python 3 API and supports multiple DL datatypes
- Can be used on single and multi-GPU nodes of various topologies
- Can be used with popular DL frameworks or integrated into a project via its API
- Compatible with TensorFlow and PyTorch frameworks

#### **CosmoFlow Scaling Performance**

- Achieved 77% scaling efficiency at 8192 nodes on Cori
  - Fully synchronous SGD
  - Speedup of 6324X
- Measure walltime per epoch (throughput)
  - Captures end-to-end capability including:
    - Single-node computation
      - Training and validation
    - Communication
    - I/O



Note: global batch size = # nodes (local batch size of 1)





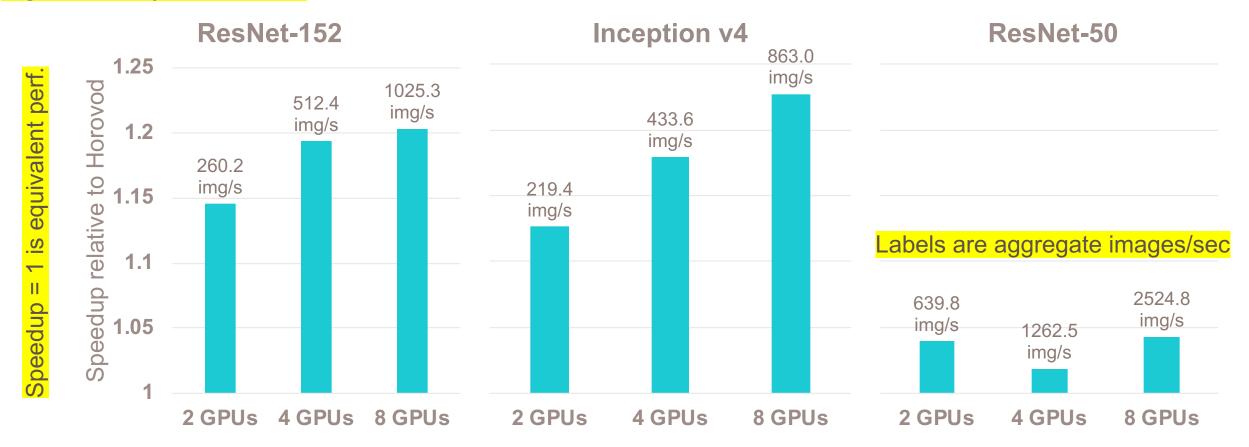
cscs

Centro Svizzero di Calcolo Scientifico Swiss National Supercomputing Cen nte

## Cray DL Plugin Perf. vs. Horovod on CS-Storm



Higher on the y-axis is better



\* TensorFlow 1.13.1 with tf\_cnn\_benchmarks, minibatch size = 32, synthetic data. Horovod 0.16.

#### QUESTIONS?



a Hewlett Packard Enterprise company