Using Hardware Performance Counters on the Cray XT

Luiz DeRose Programming Environment Director Cray Inc. Idr@cray.com

CSCS July 13-15, 2009



Luiz DeRose (Idr@cray.com) © Cray Inc.



Simplified memory hierachy on the Quad Core AMD Opteron





Hardware Performance Counters

- AMD Opteron Hardware Performance Counters
 - Four 48-bit performance counters.
 - Each counter can monitor a single event
 - Count specific processor events
 - » the processor increments the counter when it detects an occurrence of the event
 - » (e.g., cache misses)
 - Duration of events
 - » the processor counts the number of processor clocks it takes to complete an event
 - » (e.g., the number of clocks it takes to return data from memory after a cache miss)
 - Time Stamp Counters (TSC)
 - Cycles (user time)



PAPI Predefined Events

- Common set of events deemed relevant and useful for application performance tuning
 - Accesses to the memory hierarchy, cycle and instruction counts, functional units, pipeline status, etc.
 - The "papi_avail" utility shows which predefined events are available on the system – execute on compute node
- PAPI also provides access to native events
 - The "papi_native_avail" utility lists all AMD native events available on the system – execute on compute node
- Information on PAPI and AMD native events
 - pat_help counters
 - man papi_counters
 - For more information on AMD counters:
 - http://www.amd.com/us-en/assets/content_type/white_papers_and_tech_docs/26049.PDF



Hardware Counters Selection

- PAT_RT_HWPC <set number> | <event list>
 - Specifies hardware counter events to be monitored
 - A set number can be used to select a group of predefined hardware counters events (recommended)
 - CrayPat provides 19 groups on the Cray XT systems
 - Alternatively a list of hardware performance counter event names can be used
 - Maximum of 4 events
 - Both formats can be specified at the same time, with later definitions overriding previous definitions
 - Hardware counter events are not collected by default
 - Hardware counters collection is not supported with sampling on systems running Catamount on the compute nodes



Accuracy Issues

- Granularity of the measured code
 - If not sufficiently large enough, overhead of the counter interfaces may dominate
- Pay attention to what is not measured:
 - Out-of-order processors
 - Speculation
 - Lack of standard on what is counted
 - Microbenchmarks can help determine accuracy of the hardware counters
- For more information on AMD counters:
 - architecture manuals:
 - http://www.amd.com/us-en/assets/content_type/white_papers_and_tech_docs/26049.PDF





Hardware Performance Counters

Hardware performance counter events: PAPI TOT INS Instructions completed PAPI L1 DCA Level 1 data cache accesses PAPI FP OPS Floating point operations DATA CACHE MISSES Data Cache Misses CYCLES USER User Cycles (approx, from clock ticks) Estimated minimum overhead per call of a traced function, which was subtracted from the data shown in this report (for raw data, use the option: -s overhead=include): PAPI TOT INS 2021.905 instr PAPI_L1_DCA 1275.739 refs 0.000 ops PAPI FP OPS DATA CACHE MISSES 7.702 misses 0.000 cycles CYCLES USER Time 2.054 microseconds



PAT_RT_HWPC=1 (Summary with TLB)

PAPI_TLB_DM Data translation lo PAPI_L1_DCA Level 1 data cache PAPI_FP_OPS Floating point ope: DC_MISS Data Cache Miss User_Cycles Virtual Cycles	ookaside buffer misses accesses rations	PAT_RT_HWPC=1 Flat profile data
USER		Hard counts
Time% Time Imb.Time Imb.Time% Calls 0.0011 PAPI_L1_DCM 14.3651 PAPI_TLB_DM 2.1141 PAPI_L1_DCA 276.0561 PAPI_FP_OPS 382.6971 User time (approx) 1.291 Average Time per Call CrayPat Overhead : Time 0.2% HW FP Ops / User time 382.6971 HW FP Ops / User time 382.6971 HW FP Ops / WCT 369.4681 Computational intensity 0.17 MFLOPS (aggregate) 6123.161 TLB utilization 130.61 D1 cache hit,miss ratios 94.8% D1 cache utilization (M) 19.22	68.0% 1.336838 secs secs M/sec 1500.0 calls M/sec 18539562 misses M/sec 2727811 misses M/sec 356285406 refs M/sec 493918940 ops secs 2839375000 cycles 0.000891 sec M/sec 493918940 ops 4.3 M/sec ops/cycle 1.39 ops/ref M/sec refs/miss 0.255 avg uses hits 5.2% misses refs/miss 2.402 avg uses	96.5%Time 3%peak(DP)



PAT_RT_HWPC=2 (L1 and L2 Metrics)

USER			
Time%		57.5%	
Time		1.282322	secs
Imb.Time			secs
Imb.Time%			
Calls	0.001M/sec	1500.0	calls
REQUESTS_TO_L2:DATA	92.608M/sec	117033567	req
DATA_CACHE_REFILLS:			
L2 MODIFIED:L2 OWNED:			
L2_EXCLUSIVE:L2_SHARED	9.691M/sec	12247253	fills
DATA_CACHE_REFILLS_FROM_SYS	STEM:		
ALL	23.312M/sec	29461089	fills
PAPI_L1_DCA	285.477M/sec	360771229	refs
User time (approx)	1.264 secs	2780250000	cycles 98.6%Time
Average Time per Call		0.000855	sec
CrayPat Overhead : Time	0.1%		
D1 cache hit,miss ratio (R)	88.4% hits	11.6%	misses
D1 cache utilization	8.65 refs/	refill 1.081	avg uses
D2 cache hit,miss ratio	74.8% hits	25.2%	misses
D1+D2 cache hit,miss ratio	91.8% hits	8.2%	misses
D1+D2 cache utilization	12.25 refs/	miss 1.531	avg uses
System to D1 refill	23.312M/sec	29461089	lines
System to D1 bandwidth	1422.878MB/sec	1885509672	bytes
L2 to Dcache bandwidth	591.504MB/sec	783824176	bytes



PAT_RT_HWPC=3 (Bandwidth)

==				
US	SER / mlwxyz_			
	Time%		44.0%	
	Time		5.393606	
	Imb.Time		0.054000	
	Imb.Time%		1.0%	
	Calls		10	
	QUADWORDS_WRITTEN_TO_SYSTEM:			
	ALL	76.516M/sec	410363958	ops
	DATA_CACHE_REFILLS:			
	L2 MODIFIED:L2 OWNED:			
	L2 EXCLUSIVE:L2 SHARED	14.494M/sec	77731399	fills
	DATA_CACHE_REFILLS_FROM_SYSTEM:			
	ALL	18.999M/sec	101891701	fills
	DATA_CACHE_LINES_EVICTED:ALL	52.589M/sec	282042348	ops
	User time (approx)	5.363 secs	11262496875	cycles
	Average Time per Call		0.539361	sec/call
	Cycles	5.363 secs	11262496875	cycles
	User time (approx)	5.363 secs	11262496875	cycles
	Utilization rate		99.4%	
	D2 cache hit ratio		43.3%	
	System to D1 refill	18.999M/sec	101891701	lines
	System to D1 bandwidth	1159.587MB/sec	6521068888	bytes
	L2 to Dcache bandwidth	884.629MB/sec	4974809544	bytes
	L2 to System BW per core	583.773MB/sec	3282911662	bytes



PAT_RT_HWPC=5 (Floating point mix)

USER			
Time%		58.5%	
Time		1.166749	secs
Imb.Time			SECS
Imb.Time%			
Calls	0.001M/sec	1500.0	calls
RETIRED_MMX_AND_FP_INS	TRUCTIONS:		
PACKED_SSE_AND_SSE2	481.704M/sec	544927850	instr
PAPI_FML_INS	153.030M/sec	173115267	ops
PAPI_FAD_INS	283.583M/sec	320803673	ops
PAPI_FDV_INS	7.258M/sec	8210206	ops
User time (approx)	1.131 secs	2601875000	cycles 97.0%Time
Average Time per Call		0.000778	sec
CrayPat Overhead : Tim	e 0.2%		<i>i</i>
HW FP Ops / Cycles		0.19	ops/cycle
HW FP Ops / User time	436.613M/sec	493918940	ops 4.7%peak(DP)
HW FP Ops / WCT	423.329M/sec		
FP Multiply / FP Ops		35.0%	
FP Add / FP Ops		65.0%	
MFLOPS (aggregate)	6985.81M/sec		



PAT_RT_HWPC=12 (QC Vectorization)

=======================================	================	=============	-======	============
USER				
Time%		62.6%		
Time		1.251600	secs	
Imb.Time			secs	
Imb.Time%				
Calls	0.001M/sec	1500.0	calls	
RETIRED_SSE_OPERATIONS:				
SINGLE_ADD_SUB_OPS:				
SINGLE_MUL_OPS		0	ops	
RETIRED_SSE_OPERATIONS:				
DOUBLE_ADD_SUB_OPS:				
DOUBLE_MUL_OPS	199.842M/sec	248803518	ops	
RETIRED_SSE_OPERATIONS:				
SINGLE_ADD_SUB_OPS:				
SINGLE_MUL_OPS:OP_TYPE		0	ops	
RETIRED_SSE_OPERATIONS:				
DOUBLE_ADD_SUB_OPS:				
DOUBLE_MUL_OPS:OP_TYPE	396.722M/sec	493918940	ops	
User time (approx)	1.245 secs	2863500000	cycles	99.5%Time
Average Time per Call		0.000834	sec	
CrayPat Overhead : Time	0.2%			



Vectorization Example

USER / calc2				
Time% Time		28.2% 0.600875	secs	
Imb.Time		0.069872	secs	
Calls	864.9 /sec	500.0	calls	
RETIRED SSE OPERATIONS: SINGLE ADD SUB OPS:				
SINGLE_MUL_OPS		0	ops	
DOUBLE ADD SUB OPS:				
DOUBLE_MUL_OPS	369.139M/sec	213408500	ops	
RETIRED SSE OPERATIONS: SINGLE ADD SUB OPS:				
SINGLE_MUL_OPS:OP_TYPE		0	ops	
RETIRED SSE OPERATIONS:				
DOUBLE_MUL_OPS:OP_TYPE	369.139M/sec	213408500	ops	
User timē (approx) —	0.578 secs	1271875000	cycles	96.2%Time
When compiled with fastsse	:			
======================================			=======	
Time% Time		24.3%	5005	
Imb.Time		0.146551	secs	
Imb.Time%	0.00114/	26.4%		
Calls RETIRED SSE OPERATIONS:	0.001M/sec	500.0	calls	
SINGLE_ADD_SUB_OPS:				
SINGLE MUL OPS Retired Sse Operations:		0	ops	
DOUBLE_ADD_SUB_OPS:				
DOUBLE MUL OPS DETIDED SEE ODEDATIONS	208.641M/sec	103016531	ops	
SINGLE ADD SUB OPS:				
SINGLE MUL OPS: OP TYPE		0	ops	
DOUBLE ADD SUB OPS:				
DOUBLE_MUL_OPS:OP_TYPE	415.628M/sec	205216531	ops	
User time (approx)	0.494 secs	1135625000	cycles	100.0%Time



How do I interpret these derived metrics?

- The following thresholds are guidelines to identify if optimization is needed:
 - Computational Intensity: < 0.5 ops/ref
 - This is the ratio of FLOPS by L&S
 - Measures how well the floating point unit is being used
 - FP Multiply / FP Ops or FP Add / FP Ops: < 25%</p>
 - Vectorization: < 1.5



Memory Hierarchy Thresholds

TLB utilization: < 90.0%</p>

- Measures how well the memory hierarchy is being utilized with regards to TLB
- This metric depends on the computation being single precision or double precision
 - A page has 4 Kbytes. So, one page fits 512 double precision words or 1024 single precision words
- TLB utilization < 1 indicates that not all entries on the page are being utilized between two TLB misses
- D1 cache utilization: < 1 (D1+D2 cache utilization: < 1)</p>
 - A cache line has 64 bytes (8 double precision words or 16 single precision words)
 - D1 cache utilization < 1 indicates that not all entries on the cache line are being utilized between two cache misses
- D1 cache hit (or miss) ratios: < 90% (> 10%)
- D2 (L2) cache hit (or miss) ratios: < 95% (> 5%)
- D1 + D2 cache hit (or miss) ratios: < 92% (> 8%)
 - D1 and D2 caches on the Opteron are complementary
 - This metric provides a view of the Total Cache hit (miss) ratio

Detecting Load Imbalance on the Cray XT

Luiz DeRose Programming Environment Director Cray Inc. Idr@cray.com

CSCS July 13-15, 2009



Luiz DeRose (Idr@cray.com) © Cray Inc.



Motivation for Load Imbalance Analysis

- Increasing system software and architecture complexity
 - Current trend in high end computing is to have systems with tens of thousands of processors
 - This is being accentuated with multi-core processors
- Applications have to be very well balanced In order to perform at scale on these MPP systems
 - Efficient application scaling includes a balanced use of requested computing resources
- Desire to minimize computing resource "waste"
 - Identify slower paths through code
 - Identify inefficient "stalls" within an application



Cray Tools Load Imbalance Support

- Very few performance tools focus on load imbalance
 - Need standard metrics
 - Need intuitive way of presentation
- CrayPat support:
 - Imbalance time and %
 - MPI sync time
 - OpenMP Performance Metrics
 - MPI rank placement suggestions
- Cray Apprentice² support:
 - Load imbalance visualization



Imbalance Time

- Metric based on execution time
- It is dependent on the type of activity:
 - User functions
 - Imbalance time = Maximum time Average time
 - Synchronization (Collective communication and barriers)
 Imbalance time = Average time Minimum time
- Identifies computational code regions and synchronization calls that could benefit most from load balance optimization
- Estimates how much overall program time could be saved if corresponding section of code had a perfect balance
 - Represents upper bound on "potential savings"
 - Assumes other processes are waiting, not doing useful work while slowest member finishes



Load balance metric - rationale



Imbalance %

Imbalance% = 100 X
$$\frac{\text{Imbalance time}}{\text{Max Time}}$$
 X $\frac{\text{N}}{\text{N}-1}$

- Represents % of resources available for parallelism that is "wasted"
- Corresponds to % of time that rest of team is not engaged in useful work on the given function
- Perfectly balanced code segment has imbalance of 0%
- Serial code segment has imbalance of 100%

Call Tree Visualization (Swim3d)

Discrete Unit of Help (DUH Button)

Load Distribution

Profile with Load Distribution by Groups

Т	able 1:	Profile by 1	Function Gro	oup and F	unction	
	Time % 	Time I 	mb. Time 1	Imb. 0 	Calls G 	roup Function PE='HIDE'
1	100.0%	0.482144			2530 т	otal
	83.7%	0.403314			303	USER
	32.48	0.156028	0.009882	6.8%	98	calc3_
	27.7%	0.133643	0.007400	6.0%	100	calc2_
	21.0%	0.101406	0.002552	2.8%	100	calc1
İ	2.0%	0.009696	0.000287	3.3%	1	inital_
	16.3%	0.078830		======== 	======= 2227	======== MPI
	12.7%	0.061266	0.078133	64.1%	351	mpi_waitall_
ĺ	2.2%	0.010607	0.011582	59.7%	936	mpi_isend
İ	1.4%	0.006945	0.004463	44.7%	936	mpi_irecv_

MPI Sync Time

- Measure load imbalance in programs instrumented to trace MPI functions to determine if MPI ranks arrive at collectives together
- Separates potential load imbalance from data transfer
- Sync times reported by default if MPI functions traced
- If desired, PAT_RT_MPI_SYNC=0 deactivated this feature

MPI Sync Time Statistics

	Time	୫ 			Ti	.me	: : 	Imł	э.	Ti	me	 	Ti	Imb Ime	• 8 	С	alls 	Gr F	oup unction PE='HIDE	,	
	100.	0%	7	.19	37	'14	- 1					1			- 1	1	7604	То	tal	_	
	76	.5%	1	5.5	500	07	8	I				-					4752	טן	SER	_	
	9	 6.0% 3.2%		5.	27	 77 73	91		0.	17	184	48 82		3	 .3% 1%		12	2	sweep_	-	
		0.38 0.28	i	0.	01	.85	88		0.	00	052 301	27 23		2	.98 88		12	- 2)	flux_err_	-	
		0.20 0.18 0 18	i	0.	00	50	32		0.	00	014	44 54		2	.00 .98 28		1	/ 	initializ	ze_	
		0.1% ====	 ==	0.	.00	28	19	 ===	0.	00	17	73 ===	 ===	40	.3%	י ו ==	2280)	rcv_real	-	
i	16	. 6 %		1.1	.97	32	1					-					4603	M	IPI	-	
	9: .	3.9% 5.9% 0.2%		1. 0. 0.	12 07 00	:42 '04 22	27 81 10	 	0. 0. 0.	27 01 00	787 443 108	78 37 38	 	20 17 34	.7% .7% . 4 %	 	2280 2280 32)) 2	<pre>mpi_recv mpi_send mpi_allre</pre>	- educe_	_
	6	.3%		0.4	153	09	1					-		·			39	M	IPI_SYNC	_	
	6: 3:	1.1% 8.7% 0.1%	 	0 . 0 . 0 .	27 17 00	70 55 05	12 64 15	 	0. 0. 0.	21 27 00	560 004 020	08 49 65	 	45 63 35	.78 .28 .58	 	4 32 3	1 2 3	mpi_bcast mpi_allre mpi_barri	c_(syneduce_ ler_(s	nc) _(sync) sync)

OpenMP (Ideal) Instrumentation

CrayPat OpenMP Performance Metrics

- Per-thread timings
- Overhead incurred at enter/exit of
 - Parallel regions
 - Worksharing constructs within parallel regions
- Load balance information across threads
- Sampling performance data without API
- Separate metrics for OpenMP runtime and OpenMP API calls

OpenMP Data from pat_report

- Default view (no options needed to pat_report)
 - Focus on where program is spending its time
 - Shows imbalance across all threads
 - Assumes all requested resources should be used
 - Highlights non-uniform imbalance across threads

MPI Rank Reorder

MPI rank placement with environment variable

- Distributed placement
- SMP style placement
- Folded rank placement
- User provided rank file

Rank Reorder Example - hycom

pat_report -0 load_balance									
Table 2: Load Balance across PE's by FunctionGroup									
Time % Cum. Time Calls Group Time % PE[mmm]									
100.0% 100.0% 482.705844 7446623155 Total									
57.7% 57.7% 278.657370 7329740077 USER									
0.5% 0.5% 361.310805 33130409 pe.201									
0.4% 58.2% 311.898417 34020074 pe.45									
0.0% 100.0% 23.780267 320096 pe.184									
42.3% 100.0% 204.048383 116783478 MPI									
0.9% 0.9% 476.662251 399087 pe.184									
0.3% 61.9% 167.921814 422197 pe.37									
0.2% 100.0% 119.123503 514637 pe.201									

Rank Reorder Example - hycom

pat report -O load balance -s pe=ALL Table 2: Load Balance across PE's by FunctionGroup Time % | Cum. | Time | Calls |Group | Time % | | I PE 100.0% | 100.0% | 482.705844 | 7446623155 |Total 57.7% | 57.7% | 278.657370 | 7329740077 |USER || 0.5% | 0.5% | 361.310805 | 33130409 |pe.201 || 0.5% | 1.0% | 349.849557 | 30460022 |pe.182 || 0.5% | 1.5% | 346.919713 | 33685730 |pe.200 || 0.5% | 2.0% | 342.844256 | 34879988 |pe.188 || 0.5% | 2.5% | 342.308415 | 34913960 |pe.172 || 0.1% | 99.8% | 45.464691 | 3000260 |pe.248 || 0.1% | 99.9% | 35.970972 | 399401 |pe.213 || 0.0% | 99.9% | 27.431543 | 340673 |pe.232 || 0.0% | 100.0% | 25.142167 | 117620 |pe.240 || 0.0% | 100.0% | 23.780267 | 320096 |pe.184

Rank Reorder Example - hycom

After custom placement (10% performance improvement):

Table 2: Load Balance with MPI Sent Message Stats

Time % 	Time Sent M Cou	Msg Sent M unt Total Byt	Asg Avg Sent ces Msg Size	Group PE[mmm]
100.0% 437.43	18783 171618	829 2893282858	340 16858.83	Total
60.2% 263.2	11966			USER
0.5% 322.0 0.4% 286.1	019049 179471	 		pe.158 pe.126
0.0% 23.3 ==================================	318648 ====================================	 ===================================	 	pe.184 ====================================
	091071 62	2224 635942	2368 10220.21	pe.184
0.3% 151.2 0.3% 115.2	242560 68 396258 68	8002 1039329 8002 1039329	9136 15283.80 9136 15283.80	pe.126 pe.158

MPI Rank Placement Suggestions

- When to use?
 - Point-to-point communication consumes significant fraction of the program time and have a significant imbalance
 - pat_report -O mpi_sm_rank_order ...
 - When there seems to be a load imbalance of another type
 - Can get a suggested rank order file based on user time
 - pat_report -O mpi_rank_order ...
 - Can have a different metric for load balance
 - pat_report -O mpi_rank_order -s mro_metric=DATA_CACHE_MISSES ...
- Information in resulting report
 - Available if MPI functions traced (-g mpi)
- Custom placement files automatically generated

MPI Rank Placement Suggestions (cont'd)

- See table notes in resulting report from pat_report
- Report provides quad core and dual core suggestions
- Set MPICH_RANK_REORDER_METHOD environment variable
 - Set to numerical value or MPICH_RANK_ORDER file from pat_report

Example: -O mpi_rank_order (asura)

Notes for table 1:

To maximize the locality of point to point communication, choose and specify a Rank Order with small Max and Avg Sent Msg Total Bytes per node for the target number of cores per node.

To specify a Rank Order with a numerical value, set the environment variable MPICH RANK REORDER METHOD to the given value.

To specify a Rank Order with a letter value 'x', set the environment variable MPICH RANK REORDER METHOD to 3, and copy or link the file MPICH RANK ORDER.x to MPICH RANK ORDER.

Table 1: Sent Message Stats and Suggested MPI Rank Order

Sent Msg Total Bytes per MPI rank

	Max		Avg		Min	Max	Min
Total	Bytes	Total	Bytes	Total	Bytes	Rank	Rank

378638104 271474542 169280552 56 109

Quad core: Sent Msg Total Bytes per node

, 63
, 63
10,111
52
,117

Example: File MPICH_RANK_ORDER.u (asura)

```
Suggested custom rank placement:
#
#
#
      pat report -0 mpi sm rank order \
#
        /home/crayadm/ldr/ASURA/asura10it.x+apa+4442-824tdt.ap2
#
#
    Targets multi-core processors, based on Sent Msg Total Bytes.
#
#
    Program:
                 /work/crayadm/ldr/ASURA/run/asura10it.x
#
    Number PEs: 128
#
    Cores/Node:
                  4
#
#
    Heuristic: u
86, 27, 108, 63, 13, 67, 23, 39, 70, 3, 113, 17, 21, 46, 40, 89
28,36,34,10,7,127,41,105,94,25,12,38,6,75,57,60
56,109,106,68,42,66,43,79,72,45,85,80,33,111,49,107
14,103,114,9,126,52,78,2,55,88,87,118,119,64,15,16
90,102,122,31,37,123,29,59,71,53,98,82,92,124,35,91
5,125,115,11,97,95,30,54,19,4,69,0,62,110,51,112
26, 32, 121, 77, 65, 100, 76, 24, 58, 74, 1, 18, 101, 116, 84, 50
44,96,93,20,83,61,104,47,99,81,120,73,8,117,22,48
```


MPI + OpenMP? (some ideas)

- When does it pay to add OpenMP to my MPI code?
 - Only add OpenMP when code is network bound
 - Adding OpenMP to memory bound codes will most likely hurt performance rather than help it
 - Look at collective time, excluding sync time: this goes up as network becomes a problem
 - Look at point-to-point wait times: if these go up, network may be a problem

Parallel Performance Analysis and Visualization on the Cray XT

Luiz DeRose Programming Environment Director Cray Inc. Idr@cray.com

CSCS July 13-15, 2009

Luiz DeRose (Idr@cray.com) © Cray Inc.

Cray Apprentice²

- Call graph profile
- Communication statistics
- Time-line view
 - Communication
 - I/O
- Activity view
- Pair-wise communication statistics
- Text reports
- Source code mapping

- Cray Apprentice²
- is target to help and correct:
 - Load imbalance
 - Excessive communication
 - Network contention
 - Excessive serialization
 - I/O Problems

Statistics Overview

Switch Overview display

Function Profile

Load Balance View (Aggregated)

Call Tree View

Call Tree View – Function List

Load Balance View (from Call Tree)

Ele v sweep3d+tr-u+mpi96p.ap2 v swim+tr16p.ap2 v Sweep3d+tr-u+mpi96p.ap2 v swim+tr16p.ap2 v sw		Min, Avg, and N Values	Max Help
PE Calls PE #33	Load Balance: MPI_Bcast	Time (in secs) Image:	-1, +1 Std Dev marks
	0	1.2e-05	2.2e-05 2.6e-05 3.14

Source Mapping from Call Tree

EleHelp• sweep3d+mpi96p+tr.xml.gz• Overview• Traffic Report• Activity• Call Graph• sweep.f• Overview• If (k2.1t.0.or.• kbc.eq.0) then• f• Overview• If (k2.1t.0.or.• kbc.eq.0) then• f• Overview• If (do_dsa) then• f• f• Overview• If (do_dsa) then• f• f• If (do_dsa)• hen• f• f• If (do_i = 1, it• f• ff• If (do_i = 1, it•	🗙 Appre	entice2 2.3	
<pre>vseep3d+mpi96p+tr xml.gz</pre>	<u>Fi</u> le		<u>H</u> elp
Image: Second state of the second s	▼sweep3c	d+mpi96p+tr.xml.gz	
Overview Traffic Report ▲ Activity C call Graph ■ sweep.f 165 166 167 168 169 169 170 171 170 171 170 171 170 172 173 173 173 174 175 175 176 176 176 177 177 178 177 178 178 178 178 178 179 179 179 179 179 178 179 179 178 179 178 179 178 179 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 178 179 179 178 179 178 179 178 177 178 177 178 177 178 178 178 178		ی 🛃 🐝 🔢 🏎 📰 ک	
<pre>165 166 c angle pipelining loop (batches of mmi angles) 167 167 168 D0 mo = 1, mmo 170 171 c K-inflows (k=k0 boundary) 172 c 167 if (k2.lt.0.or.kbc.eq.0) then 174 do j = 1, jt 175 do i = 1, it 177 phikb(i,j,mi) = 0.0d+0 178 end do 180 end do 180 else 182 if (do_dsa) then 183 leak = 0.0 184 k = k0 - k2 185 do mi = 1, mmi 186 m = mi + mio 187 do j = 1, jt 188 do i = 1, it 189 do i = 1, jt 188 do i = 1, jt 189 do i = 1, jt 189 do i = 1, jt 190 leak = leak 191 & + wtsi(m)*phikb(i,j,mi)*di(i)*dj(j) 192 face(i,j,k+K3,3) = face(i,j,k+K3,3) 194 end do 195 end do 196 end do 197 leak = leak 191 & + wtsi(m)*phikb(i,j,mi) 193 & end do 194 end do 195 end do 195 end do 196 end do 197 leak = leak 197 leak = leak 197 leak = leak 197 leak = leak 197 leak = leak 197 leak = leak 197 leak = leak 197 leak = leak 197 leak = leak 198 leak ge(5) = leakage(5) + leak 199 leak = leak 197 leak = leak 197 leak = leak 197 leak = leak 198 leak ge(5) = leakage(5) + leak 199 leak = leak 199 leak ge(5) = leakage(5) + leak 190 leak = le</pre>	🔻 Overvie	ew 💌 Traffic Report 🔍 Activity 🔍 Call Graph 🔍 sweep.f	
180 else 182 if (do_dsa) then 183 leak = 0.0 184 k = k0 - k2 185 do mi = 1, mmi 186 m = mi + mio 187 do j = 1, jt 188 do i = 1, it 189 phikb(i,j,mi) = phikbc(i,j,m) 190 leak = leak 191 & 192 face(i,j,k+k3,3) = face(i,j,k+k3,3) 193 & 194 end do 195 end do 196 end do 197 leakage(5) = leakage(5) + leak 100 2.13	166 167 168 169 170 171 172 173 174 175 176 177 178 179	<pre>c angle pipelining loop (batches of mmi angles) c D0 mo = 1, mmo mio = (mo-1)*mmi c K-inflows (k=k0 boundary) c if (k2.lt.0 .or. kbc.eq.0) then do mi = 1, mmi do j = 1, jt do i = 1, it phikb(i,j,mi) = 0.0d+0 end do</pre>	
193 & + wtsi(m)*phikb(i,j,mi) 194 end do 195 end do 196 end do 197 leakage(5) = leakage(5) + leak 100 2.13 4.27 6.40	180 181 182 183 184 185 186 187 188 189 190 191	<pre>end do else if (do_dsa) then leak = 0.0 k = k0 - k2 do mi = 1, mmi m = mi + mio do j = 1, jt do i = 1, it phikb(i,j,mi) = phikbc(i,j,m) leak = leak</pre>	
0.00 2.13 4.27 6.40 8.53	192 193 194 195 196 197 197	<pre></pre>	. ▼
	0.00	2.13 4.27 6.40	8.53

Function Profile

vim+iompi+156	6td.ap2	▼ I+nv	/1+swp+	io+mpi+48p	.ap2				
	8	١	16		≫ _	🏂 🚧 🧱			
Iverview 🔻 F	unction								
Time	Percent	Hits	Callsites	imbalance %	Potential Savings	Function	Line	File	
124.175511	63.29	576	1	5.63	0.15	sweep_	116	/lus/hid00036/ldr/Apps/sweep3d/sweep.f	
40.211774	20.50	118080	1	23.40	0.25	mpi_recv_	0	==NA==	
16.319527	8.32	48	1	48.26	0.30	exit	35	/hotbackedup/users/rsrel/rs64.REL_1_4_33.060914.Thu/pe/computelibs/glibc/stdlib/exit.c	
6.173236	3.15	1536	3	50.00	0.12	mpi_allreduce_	0	==NA==	
2.760376	1.41	118080	1	17.58	0.01	mpi_send_	0	==NA==	
2.250029	1.15	576	1	2.62	0.00	source_	18	/lus/hid00036/ldr/Apps/sweep3d/source.f	
1.984620	1.01	144	1	2.59	0.00	mpi_barrier_	0	==NA==	
0.867359	0.44	192	2	2.47	0.00	mpi_bcast_	0	==NA==	
0.416231	0.21	576	1	2.60	0.00	flux_err_	17	/lus/hid00036/ldr/Apps/sweep3d/flux_err.f	
0.382130	0.19	118080	2	10.98	0.00	snd_real_	135	/lus/hid00036/ldr/Apps/sweep3d/mpi_stuff.f	
0.237772	0.12	309	1	95.76	0.07	fwrite	- 36	/hotbackedup/users/rsrel/rs64.REL_1_4_33.060914.Thu/pe/computelibs/glibc/libio/iofwrite.c	
0.185067	0.09	118080	2	17.30	0.00	rcv_real_	164	/lus/hid00036/ldr/Apps/sweep3d/mpi_stuff.f	
0.067832	0.03	48	1	4.56	0.00	initialize_	42	Aus/hid00036/ldr/Apps/sweep3d/initialize.f	
0.059407	0.03	48	1	4.99	0.00	initxs_	77	Aus/hid00036/ldr/Apps/sweep3d/initialize.f	
0.041371	0.02	48	1	23.84	0.00	inner_	72	Aus/hid00036/ldr/Apps/sweep3d/inner.f	
0.023948	0.01	8	1		0.00	fputc	35	/hotbackedup/users/rsrel/rs64.REL_1_4_33.060914.Thu/pe/computelibs/glibc/libio/fputc.c	
0.016902	0.01	68	1		0.00	getc	36	/hotbackedup/users/rsrel/rs64.REL_1_4_33.060914.Thu/pe/computelibs/glibc/libio/getc.c	
0.008104	0.00	4992	2	28.14	0.00	octant_	17	Aus/hid00036/ldr/Apps/sweep3d/octant.f	
0.002457	0.00	576	1	18.79	0.00	global_real_max_	321	/lus/hid00036/ldr/Apps/sweep3d/mpi_stuff.f	
0.002083	0.00	48	1	69.55	0.00	MAIN_	72	/lus/hid00036/ldr/Apps/sweep3d/driver.f	
0.001588	0.00	576	1	39.10	0.00	global_int_sum_	373	/lus/hid00036/ldr/Apps/sweep3d/mpi_stuff.f	
0.001393	0.00	48	1	10.23	0.00	inner_auto_	69	Aus/hid00036/ldr/Apps/sweep3d/inner_auto.f	
0.000982	0.00	48	1	97.74	0.00	task_init_	24	/lus/nid00036/ldr/Apps/sweep3d/mpi_stuff.f	
0.000739	0.00	384	2	27.87	0.00	global_real_sum_	347	/lus/nid00036/ldr/Apps/sweep3d/mpi_stuff.f	
0.000662	0.00	2	1		0.00	fopen	106	/hotbackedup/users/rsrel/rs64.REL_1_4_33.060914.Thu/pe/computelibs/glibc/libio/lofopen.c	
0.000499	0.00	48	1	7.54	0.00	initsnc_	175	/lus/nid00036/ldr/Apps/sweep3d/initialize.f	
									Þ
0				1.1	5			2.30 3.45	

Distribution by PE, by Call, & by Time

File Ele vsweep3d+mpi96p+tr.xml.gz File vsweep3d+mpi96p+tr.xml.gz vsweep3d+mpi96p+tr.xml.gz	Help
	<u>H</u> elp
Verview Traffic Report Activity Call Graph Sweep.f Mosaic	
▼ Overview ▼ Traffic Report ▼ Activity ▼ Call Graph ▼ sweep.f ▼ Mosaic ▼ Environment ▼ Function ▼ Text Report	
PE #0 (total=NA) BY CALL: func time calls func time calls	
0. MPI_Recv 0.375000 1440 0. MPI_Recv 79.578125 247680 1. MPI_Allreduce 0.250000 23 1 MPI_Recv 5959275 247590	
2. MPI_Send 0.015625 1440 2. mpi_recv_ 0.250000 247680 3. mpi_send 0.140625 247680	
PE #1 (total=NA) 4. MPI Allreduce 12.437500 3072 fmc 5. mpi_allreduce_ 0.000000 3072	
Image: Control of the control of t	
1. MPI Barrier 0.296875 3 8. MPI Bcast 26.375000 384 2. MPI Bcast 0.281250 4 9. MPI Barrier 27.812500 288	
3. MPI_Allreduce 0.250000 32 10. mpi_barrier_ 0.000000 288 4. MPI_Send 0.093750 2160 11. MPI_Comm_rank 0.000000 96	
5. mpi_recv_ 0.015625 2160 12. MPI_Comm_SLZE 0.000000 96 13. mpi_comm_rank_ 0.000000 96	
PE #2 (total=NA) 14. mp1_com_size_ 0.000000 36 PE #2 (total=NA) 15. MP1_Finalize 0.000000 96 fma 15. MP1_Finalize 0.000000 96	
17Exit 0.000000 96	
1. MPI Beast 0.281250 4 2. MPI Barrier 0.281250 3 BY TIME:	
3. MPI_Allreduce 0.234375 32 func time calls 4. MPI_send 0.031250 2160	
5. mpi_send_ 0.015625 2160 0. MPI_Recv 79.578125 247680 1. MPI Barrier 27.812500 288	
PE #3 (total=NA) 2. MPI Bcast 26.375000 384 3. MPI Allreduce 12.437500 3072	
func time calls 4. MPI_Send 5.859375 247680	
0. MPI Recv 0. 703125 2160 6. mpi send 0. 140625 247680 1. MPI Barrier 0.296875 3 7. MPI Finalize 0.000000 96	
2. MPI_BCast 0.281250 4 8. MPI_Init 0.000000 96 3. MPI_Allreduce 0.203125 32 9. MPI_Comm_rank 0.000000 96	
4. MPI_Send 0.031250 2160 10. MPI_Comm_size 0.000000 96 5. mpi_recv_ 0.015625 2160 10. MPI_Comm_size 0.000000 96	
PE #4 (total=NA) PE time calls	
0 0.656250 5877 1 1.609375 8757	
0.00 2.13 4.27 2 1.50000 8757 3 1.531250 8757	
0.00 2.13 4.27 6.40	8.53

Environment & Execution Details

Eile									Help
v s	wim+ior	mpi+1566td.ap2 🔻 T+hw1	+swp+io+mpi+4	3p.ap2					
G				~\ =	N N NITTON	-	121		
Ű,	₩	∕─ 🥶 💆 !		V 🖣		<u></u>			
	Ouaruia		oment			 			
,									
	nvira	onment							
	Env Va	ers System Info Resource	Limits Heap In	fo					
	PF	Total Used Total Free	Largest Free	Fragments					
		(мв) (мв)	(MB)						
	U	97.248817 4065.597900	4065.596680	1614	start	 			
	- a -	97.283150 4065.563477	4065.559326	1654	end-start				
	- ' I	95.571546 4067.275146	4067.274414	1534	start and atort				
	2	95.531171 4067.255015	4067.231355	1534	enu-start etart				
		95 591171 4067 255615	5 4067 251953	1592	end-start				
	3	95.571548 4067.275146	4067,274414	1534	start				
		95.591171 4067.255615	5 4067.251953	1592	end-start				
	4	95.571548 4067.275146	6 4067.274414	1534	start				
		95.591171 4067.255615	4067.251953	1592	end-start				
	5	95.571548 4067.275146	6 4067.274414	1534	start				
		95.591171 4067.255615	4067.251953	1592	end-start				
	6	95.571548 4067.275146	6 4067.274414	1534	start				
		95.591171 4067.255615	4067.251953	1592	end-start				
	7	95.571548 4067.275146	6 4067.274414	1534	start				
		95.591171 4067.255615	4067.251953	1592	end-start				
	8	95.571548 4067.275148	6 4067.274414	1534	start				
		95.591171 4067.255615	5 4067.251953	1592	end-start				
	9	95.571548 4067.275146	6 4067.274414	1534	start				
		95.591171 4067.255615	5 4067.251953	1592	end-start				
	10	95.571548 4067.275146	6 4067.274414	1534	start				
		95.591171 4067.255615	4067.251953	1592	end-start				
		95.571548 4067.275146	4067.274414	1534	start				•
0.	00		1	.15		 2.30		3.45	4.61

Time Line View (Sweep3D)

Time Line View (Fine Grain Zoom)

Activity View

Pair-wise Communication

I/O Overview

I/O Rates

<u>Fi</u> le	n3d48F	2+tr an2								<u>H</u> elp	
			M	16 <	» 😼	<u>.</u>			<u>r</u> 3		
▼Overview ▼IO Overview ▼IO Traffic Report ▼IO Rates											
FD	Write Calls	Write Tot (MB)	Write Min (MB/s)	Write Avg (MB/s)	Write Max (MB/s)	Read Calls	Read Tot (MB)	Read Min (MB/s)	Read Avg (MB/s)	Read Max (MB/s)	
1	37	0.002	0.001	0.039	0.869	0	0.000	0.000	0.000	0.000	
2	48	0.001	0.001	0.003	0.136	0	0.000	0.000	0.000	0.000	
			1.34		2.67			4 01		5.34	
0.00			1.01		2.07				_	0.01	

Hardware Counters Overview

Hardware Counters Time Line

Controlling Performance File Size

- Performance files can be quite large. There are several run-time environment variables to keep data files down to reasonable sizes
- The particular run-time environment variables to use vary depending on the type of experiment being conducted
- Sampling:
 - PAT_RT_RECORD_PE
 - Collect trace for a subset of the PEs
 - PAT_RT_RECORD_THREAD
 - Collect trace for a subset of the threads
 - PAT_RT_INTERVAL
 - Specifies the interval, at which the instrumented program is sampled
 - PAT_RT_CALLSTACK
 - Limit the depth to trace the call stack
 - PAT_RT_HWPC
 - Avoid collecting hardware counters (unset)
 - PAT_RT_SIZE
 - The number of contiguous bytes in the text segment available for sampling
 - PAT_RT_WRITE_BUFFER_SIZE
 - Specifies the size, of a buffer that collects measurement data for a single thread

Controlling Trace File Size

- Tracing:
 - PAT_RT_CALLSTACK
 - Limit the depth to trace the call stack
 - PAT_RT_HWPC
 - Avoid collecting hardware counters (unset)
 - PAT_RT_RECORD_PE
 - Collect trace for a subset of the PEs
 - PAT_RT_RECORD_THREAD
 - Collect trace for a subset of the threads
 - PAT_RT_TRACE_FUNCTION_ARGS
 - Limit the number of function arguments to be traced
 - PAT_RT_TRACE_FUNCTION_LIMITS
 - Avoid tracing indicated functions
 - PAT_RT_TRACE_FUNCTION_MAX
 - Limit the maximum number of traces generated for all functions for a single process
 - PAT_RT_TRACE_THRESHOLD_PCT
 - Specifies a % of time threshold to enforce when executing in full trace mode
 - PAT_RT_TRACE_THRESHOLD_TIME
 - Specifies a time threshold to enforce when executing in full trace mode
- Use the limit built-in command for ksh(1) or csh(1) to control how much disk space the trace file can consume

Additional API Functions

- int PAT_state (int state)
 - State can have one of the following:
 - PAT_STATE_ON
 - PAT_STATE_OFF
 - PAT_STATE_QUERY
- int PAT_record (int state)
 - Controls the state for all threads on the executing PE. As a rule, use PAT_record() unless there is a need for different behaviors for sampling and tracing
 - int PAT_sampling_state (int state)
 - int PAT_tracing_state (int state)
- int PAT_trace_function (const void *addr, int state)
 - Activates or deactivates the tracing of the instrumented function
- int PAT_flush_buffer (void)

Parallel Performance Analysis and Visualization on the Cray XT

Questions / Comments Thank You!

CSCS July 13-15, 2009

Luiz DeRose (Idr@cray.com) © Cray Inc.