

Memory Subsystem Profiling with the Sun Studio Performance Analyzer

CScADS, July 20, 2009

Marty Itzkowitz, Analyzer Project Lead

Sun Microsystems Inc.

marty.itzkowitz@sun.com

Outline

- Memory performance of applications
 - > The Sun Studio Performance Analyzer
- Measuring memory subsystem performance
 - > Four techniques, each building on the previous ones
 - First, clock-profiling
 - Next, HW counter profiling of instructions
 - Dive deeper into dataspace profiling
 - Dive still deeper into machine profiling
 - What the machine (as opposed to the application) sees
 - > Later techniques needed if earlier ones don't fix the problems
- Possible future directions

No Comment



The Message

- Memory performance is crucial to application performance
 - > And getting more so with time
- Memory performance is hard to understand
 - > Memory subsystems are very complex
 - All components matter
 - > HW techniques to hide latency can hide causes
- Memory performance tuning is an art
 - > We're trying to make it a science
- The Performance Analyzer is a powerful tool:
 - > To capture memory performance data
 - > To explore its causes

Memory Performance of Applications

- Operations take place in registers
 - > All data must be loaded and stored; latency matters
- A load is a load is a load, but
 - > Hit in L1 cache takes 1 clock
 - > Miss in L1, hit in L2 cache takes ~10-20 clocks
 - > Miss in L1, L2, hit in L3 cache takes ~50 clocks
 - > Fetch from memory takes ~200-500 clocks (or more)
 - > Page-in from disk takes milliseconds
 - Costs are typical; each system is different
- What matters is total stalls in the pipeline
 - > If latency is covered, there's no performance cost

Why Memory Performance is Hard

- SW developers know code, algorithms, data structures
 - > What the HW does with them is magic
 - Many, if not most, SW developers can't even read assembler
- HW engineers know instruction, address streams
 - > How the SW generates them is magic
- HW performance optimizations further confuse the issue
- Difficulty lies in bridging the gap
 - > Get data to show HW perspective to SW developers
- The rest of this talk will show how we do so

Memory Performance Problems

- Some causes of memory performance problems:
 - > Initial cache miss, capacity misses
 - Layout and padding; lack of prefetch
 - > Conflict cache misses within a thread
 - Striding through arrays
 - > Coherence misses across thread
 - Sharing: unavoidable misses
 - False sharing: avoidable miss, not a real conflict
 - Threads refer to different fields in same cache line
 - Different processes use same VA for different PA's
 - > Cache and Page coloring
 - Mappings from addresses to cache lines

The Sun Studio Performance Analyzer

- Integrated set of tools for performance measurement
 - > Data collection
 - > Data examination
- Many types of data:
 - > Clock-profiling, HW counter profiling, ...
 - > Special support for OpenMP, MPI, Java
- Common command-line and GUI interface for all
- Available on SPARC and X86, Solaris and Linux
 - > It's FREE!
- You've seen it before....

Clock Profiling

- Periodic statistical sampling of callstacks
 - > `collect -p <interval> target`
 - Note: many other tools do clock-profiling, too
- Shows expensive functions, instructions
 - > Is it the heart of the computation, or is it stalled?
 - > If it's stalled,
 - Is it stalled waiting for a previous operation?
 - Is it stalled waiting for a load?
 - Is it stalled trying to do a store?
 - > Can only guess with clock profiling

Measuring Memory Costs

- Need better data to understand more
 - > See: Zaghera, *et.al.*, SC '96
- Use HW counters to trigger sampling
 - > `collect -h <cntr1>, <val1>, <cntr2>, <val2>, ...`
 - As many counters as chip supports
 - `collect` with no arguments prints list for that machine
- Collect counter name, overflow value, callstack
 - > Cache misses/references, TLB misses, instructions, ...
 - > Cycles, L1-Cache stalls, L2-Cache stalls, ...
 - Measured in cycles; convertible to time
- Shows memory costs based on the counters

Memory Performance Example

- Test code: 8 copies of vector-matrix multiply
 - > 8 functions named: **`dgemv_<opt-flag><order>`**
 - Same computation, different performance
 - > Two loop orders
 - Row, column and column,row
 - **`<order> = 1, 2`**
 - > Four optimization levels
 - Compile with **`-g, -O, -fast, and -fast -autopar`**
 - **`<opt-flag> = _g, _opt, _hi, and _p`**

Detailed Memory Performance

Sun Studio Analyzer [test.memory.1.er, ...]

File View Timeline Help

Find Text:

Functions	Callers-Callees	Source	Disassembly	Timeline	Experiments
⌚ User CPU (sec.)	⌚ CPU Cycles (sec.)	⌚ D\$ and E\$ Stall Cycles (sec.)	⌚ E\$ Stall Cycles (sec.)	⌚ I\$ Stall Cycles (sec.)	Name
59.892	47.227	26.330	21.580	0.041	<Total>
13.289	7.920	4.256	3.568	0.008	dgemv_g1_
12.028	6.693	4.242	3.592	0.007	dgemv_opt1_
8.686	8.240	0.	0.003	0.003	load_arrays_
6.004	5.440	5.060	4.172	0.007	dgemv_p2_ -- MP doall from line 40 [_\$dlB40.dgemv
5.934	5.467	5.064	4.182	0.005	dgemv_p1_ -- MP doall from line 16 [_\$dlA16.dgemv
4.133	4.133	1.045	0.776	0.001	dgemv_g2_
3.262	3.027	2.802	2.267	0.004	dgemv_hil_
3.152	3.013	2.827	2.252	0.004	dgemv_hi2_
2.992	2.880	1.032	0.767	0.001	dgemv_opt2_
0.210	0.187	0.001	0.001	0.	barrier_ -- MP doall from line 45 [_\$dlB45.barrie
0.200	0.213	0.	0.	0.	barrier_ -- MP doall from line 20 [_\$dlA20.barrie
0.	0.	0.	0.	0.	<OMP-idle>
0.	0.013	0.	0.	0.	<OMP-implicit_barrier>
0.	0.	0.	0.	0.	MAIN
0.	0.	0.	0.	0.001	__collector_ext_usage_sample
0.	0.	0.	0.	0.	__collector_log_write
0.	0.	0.	0.	0.	__f90_esfw
0.	0.	0.	0.	0.	__f90_init
0.	0.	0.	0.	0.	__f90_opencat

Summary Event

Selected Object:

Name: dgemv_g1_

PC Address: 2:0x00003428

Size: 1480

Source File: dgemv_g.f90

Object File: cachetest

Load Object: <cachetest>

Mangled Name:

Aliases:

Metrics for Selected Object:

	⌚ Exclusive	
User CPU:	13.289 (22.19%)	
Wall:	13.740 (23.66%)	
Total Thread:	13.740 (20.70%)	
System CPU:	0.260 (11.61%)	
Wait CPU:	0.030 (7.69%)	
User Lock:	0. (0. %)	
Text Page Fault:	0. (0. %)	
Data Page Fault:	0.010 (6.67%)	
Other Wait:	0.150 (7.08%)	
E\$ Stall Cycles:	1.284 (8.22%)	

Separate out costs of the various caches
Two experiments, combined in Analyzer

Memory Performance Problems

- Data shows where in program problems occur
 - > High cache misses, TLB misses
 - Does not show why
- Cause is striding through memory
 - > Clue from differences between loop order versions
 - > In this example, the compiler can diagnose
 - Studio compilers generate commentary to say what they did
 - See next slide
- In general, diagnosing these problems is hard
 - > This one is easy – other cases are more difficult

Annotated Source of dgemv_hi1

Sun Studio Analyzer [test.cpi.1.er]

File View Timeline Help

Find Text:

Functions Callers-Callees Source Disassembly Timeline Experiments

User CPU (sec.)	CPU Cycles (sec.)	Instructions Executed	Source File: dgemv_hi.f90
0.	0.	0	Object File: cachetest
			Load Object: <cachetest>
			Source loop below has tag L3
			L3 interchanged with L2
			L3 cloned for unrolling-epilog. Clone is L12
			All 8 copies of L12 are fused together as part of unroll and jam
			L12 scheduled with steady-state cycle count = 9
			L12 unrolled 4 times
			L12 has 9 loads, 1 stores, 8 prefetches, 8 FPadds, 8 FPMuls, and 0
			L12 has 0 int-loads, 0 int-stores, 11 alu-ops, 0 muls, 0 int-divs and
			L3 scheduled with steady-state cycle count = 2
			L3 unrolled 4 times
			L3 has 2 loads, 1 stores, 1 prefetches, 1 FPadds, 1 FPMuls, and 0 F
			L3 has 0 int-loads, 0 int-stores, 4 alu-ops, 0 muls, 0 int-divs and
0.	0.	0	18. DO i = 1, m
			Source loop below has tag L2
			L2 interchanged with L3
			L2 cloned for unrolling-epilog. Clone is L10
			L10 is outer-unrolled 8 times as part of unroll and jam
0.	0.	0	19. DO j = 1, n
2.362	2.413	290 000 047	20. a(i) = a(i) + b(i,j) * c(j)

Summary Event

Selected Object:

Name: line 4 in "dgemv_hi.f90"

PC Address: 2:0x0000432C

Size: 0

Source File: dgemv_hi.f90

Object File: cachetest

Load Object: <cachetest>

Mangled Name:

Aliases:

Metrics for Selected Object:

	Exclusive	
User CPU:	0.	(0. %)
Wait:	0.	(0. %)
Total Thread:	0.	(0. %)
System CPU:	0.	(0. %)
Wait CPU:	0.	(0. %)
User Lock:	0.	(0. %)
Text Page Fault:	0.	(0. %)
Data Page Fault:	0.	(0. %)
Other Wait:	0.	(0. %)
CPU Cycles:	0.	(0. %)

Loop interchange – compiler knows best order of loops
 Compiler commentary from **-fast** compilation

Dive Deeper

- We understand program instructions, not data
- Want better performance data
 - > The data addresses that trigger the problems
 - > The data objects that trigger the problems
 - *i.e.*, Source references
- Hard to get data reference address:
 - > HW counters skid past triggering instruction
 - Interrupt PC != Trigger PC
 - Current registers may not reflect state at time of event
- Solution: Dataspace profiling
 - > Built on top of HW counter profiling

Dataspace Profiling Technology

- Extend HW counter to capture more data
 - > collect -h +<cntr1>,<val1>,+<cntr2>,<val2>,...
 - + sign in front of counter name
- Causes backtracking at HW profile event delivery
 - > Capture trigger PC (might fail)
 - > Capture virtual and physical data addresses (might fail)
 - Track register changes that might affect address
 - > Post-process to see if branch-target crossed
 - Typically, 95% of backtracking succeeds
- SPARC-only functionality, alas
 - > Backtracking not possible on x86/x64
 - But instruction sampling can extend it to x86/x64

Dataspace Profiling Example

- **mcf** from SPEC cpu2000 benchmark suite
 - > Single depot vehicle scheduler; network simplex
 - Single-threaded application
- Collect two experiments
 - > `-p on -h +ecstall,lo,+ecrm,on`
 - > `-p off -h +ecref,on,+dtlbn,on`
- Combine in Analyzer
- See Itzkowitz, *et.al.*, SC|03 for details

Dataspace Profiling: Function List

Performance Analyzer [mcf.er, ...]

File View Timeline Help

Find Text:

Functions	Callers-Callees	Source	Lines	Disassembly	PCs	DataLayout	DataObjects	Timeline	LeakL
Total LWP (sec.) (%)	User CPU (sec.) (%)	E\$ Stall Cycles (%)	E\$ Read Misses (%)	E\$ Refs (%)	DTLB Misses (%)	Name			
552.677	100.0	549.404	100.0	100.0	100.0	<Total>			
282.488	51.1	280.706	51.1	61.9	62.3	38.4	88.0	refresh_potential	
128.120	23.2	127.319	23.2	30.3	29.6	13.8	9.4	primal_bea_mpp	
120.294	21.8	120.134	21.9	3.8	4.0	41.7	0.8	price_out_impl	
6.174	1.1	6.154	1.1	1.9	2.0	0.3	0.2	flow_cost	
3.763	0.7	3.733	0.7	0.9	1.0	0.4	0.1	dual_feasible	
2.702	0.5	2.682	0.5	0.6	0.6	0.8	0.1	update_tree	
1.871	0.3	1.861	0.3	0.3	0.2	0.4	0.3	primal_iminus	
0.690	0.1	0.690	0.1	0.1	0.1	0.0	0.7	write_circulations	
0.	0.	0.	0.	0.1	0.1	0.2	0.1	collector_record_counters	
0.290	0.1	0.280	0.1	0.1	0.0	0.1	0.	refresh_neighbour_lists	
0.050	0.0	0.050	0.0	0.0	0.0	0.	0.	primal_feasible	
0.060	0.0	0.060	0.0	0.0	0.0	0.0	0.	primal_start_artificial	
4.933	0.9	4.933	0.9	0.0	0.	3.4	0.	sort_basket	
0.	0.	0.	0.	0.0	0.0	0.0	0.0	collector_final_counters	
0.	0.	0.	0.	0.	0.	0.	0.	__collector_log_write	
0.090	0.0	0.090	0.0	0.	0.	0.1	0.	__doscan_u	
0.	0.	0.	0.	0.	0.	0.0	0.	__fsetlocking	
0.010	0.0	0.	0.	0.	0.	0.	0.	__open	
0.010	0.0	0.010	0.0	0.	0.	0.0	0.	__doprnt	
0.	0.	0.	0.	0.	0.	0.	0.	__endopen	

Summary Event Legend Leak

Data for Selected Object:

Name: refresh_potential

PC Address: 2:0x000030E0

Size: 412

Source File: mcf.er/archives/mcfutil.c

Object File: src/mcf

Load Object: <mcf>

Mangled Name:

Aliases:

Process Times (sec.) / Counts

	Exclusive	Incl
User CPU:	280.706 (51.1%)	280.
Wall:	282.488 (51.1%)	282.
Total LWP:	282.488 (51.1%)	282.
System CPU:	1.731 (56.2%)	1.
Wait CPU:	0.050 (27.8%)	0.
User Lock:	0. (0.%)	0.
Text Page Fault:	0. (0.%)	0.
Data Page Fault:	0. (0.%)	0.
Other Wait:	0. (0.%)	0.

Which functions have memory performance issues

Dataspace Profiling: Data Layout

Performance Analyzer [mcf.er, ...]

File View Timeline Help

Find Text:

Functions Callers-Callees Source Lines Disassembly PCs DataLayout DataObjects Timeline LeakL

E\$ Stall Cycles (%)	E\$ Read Misses (%)	E\$ Refs (%)	DTLB Misses (%)	Name
41.9	39.4	41.5	29.8	{structure:node -}
0.0	0.0	0.0	0.0	/ +0 .{long number}
0.0	0.0	0.0	0.0	+8 .{pointer+char ident}
2.4	1.3	1.8	0.5	+16 .{pointer+structure:node pred}
9.8	8.2	3.7	0.1	\ +24 .{pointer+structure:node child}
0.1	0.0	2.2	0.0	/ +32 .{pointer+structure:node sibling}
0.0	0.0	0.0	0.0	+40 .{pointer+structure:node sibling_prev}
0.1	0.1	0.4	0.2	+48 .{long depth}
20.7	21.4	8.0	24.3	\ +56 .{long orientation}
1.8	2.0	1.5	0.1	/ +64 .{pointer+structure:arc basic_arc}
0.0	0.0	0.1	0.0	+72 .{pointer+structure:arc firstout}
0.0	0.0	0.0	0.0	+80 .{pointer+structure:arc firstin}
6.1	5.2	16.7	4.5	\ +88 .{cost_t=long potential}
0.1	0.1	0.2	0.0	/ +96 .{flow_t=long flow}
0.0	0.0	1.4	0.0	+104 .{long mark}
0.8	0.9	5.5	0.0	+112 .{long time}
0.1	0.0	2.2	0.0	{structure:basket -}
0.1	0.0	0.4	0.0	/ +0 .{pointer+arc_t=structure:arc a}
0.0	0.0	0.0	0.0	+8 .{cost_t=long cost}

Summary Event Legend Leak

Data for Selected Object:

Data Object: {structure:node -}

Scope: {Global}

Type: structure:node

Member of:

Offset:

Size: 120

Elements: 15

Process Times (sec.) / Counts

Data-derived

E\$ Stall Cycles: 124.4

" count: 112140858

E\$ Read Misses: 622918

E\$ Refs: 10312057

DTLB Misses: 76302

Show costs against Data Structure Layout, not code

Dataspace Example Conclusions

- Hot references all are to **node** and **arc** fields
- Structures not well-aligned for cache
 - > Need to pad to cache-line boundary
 - > Reorganize structures to put hot fields on same line
 - Note: reorganizing might move hot fields, but not improve perf.
- High TLB misses imply need for large heap pages
- These changes led to ~21% improvement
 - > But not following SPEC cpu2000 rules
 - That does not matter for real codes, of course

Dive Still Deeper

- We understand instructions and data, but not machine
 - > So far, problems have been in a single thread
- Use same data to explore interactions among threads
 - > Sample questions to answer:
 - Which cache lines are hot?
 - Is usage uniform across lines, or is there one very hot line?
 - Which threads refer to those lines?
 - Which addresses are being referred to by which threads?

Advanced Diagnostic Strategy

- Iterative analysis of the data:
 - > Slice and dice data into sets of “objects”
 - Cache lines, pages, TLB entries, CPUs, ...
 - Threads, Processes, Time intervals
 - > Find the hot objects of one set
 - > Filter to include data only for those hot objects
 - > Look at other types of objects to see why
 - It is non-trivial to know which ones to look at
 - > Repeat as needed

Advanced Diagnostic Techniques

- Collect Dataspace profiling data
 - > `collect -h +<cntr1>,<val1>,+<cntr2>,<val2>,...`
- Collect over all threads and processes
- Slice into “Index Object” or “Memory Object” sets
 - > Each set has formula for computing an index from records
 - > Analyzer has a Tab for each object set
 - > Each Tab shows metrics for the objects in each set
 - e.g., Threads, L2-cache lines, ...

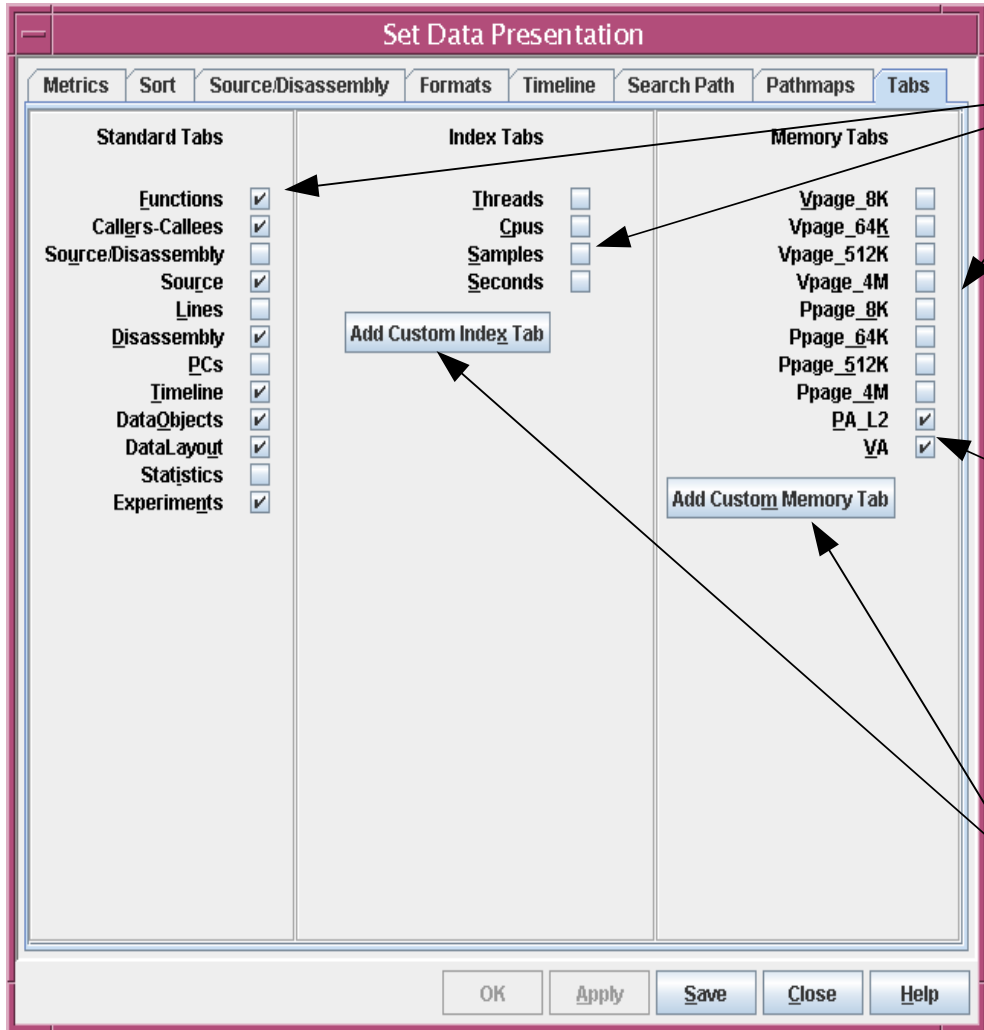
Memory and Index Objects

- Index Objects: formula does not use VADDR or PADDR
 - > Formula fields present in all records
 - > Can be used for all data
 - > Some are predefined
- Memory Objects: formula uses VADDR or PADDR
 - > Address fields present only dataspace records
 - > Definitions depend on the specific physical machine
 - Cache structure, page size, TLB organization
 - Not yet captured automatically, but could be
- Define in **.er.rc** file, based on specific machine

Sample Object Definitions

- Each definition uses one (or more) fields
 - > Thread
 - `indxobj_define Threads THRID`
 - > Virtual address
 - `mobj_define VA VADDR`
 - > L2 cache line
 - `mobj_define PA_L2 (PADDR&0x7ffc0)>>6`
- Can be a lot more complicated
 - > e.g., Niagara-2 level-2 data cache line set
 - `mobj_define UST2 L2DCacheSet \`
`(((((PADDR>>15) ^ PADDR) >>9) &0x1f0) | \`
`(((((PADDR>>7) ^ PADDR) >>9) &0xc) | \`
`(((PADDR>>9) &3))`

Displaying Objects in Analyzer Tabs



Predefined Tabs

2 Tabs from `.er.rc` file

Buttons to add custom Tabs

Example: `mttest`

- Analyzer test code
 - > Organized as series of tasks
 - Each task queues 4 blocks, spawns 4 threads
 - Threads synchronize differently for each task
 - Each thread calls one of the **compute*** functions for its block
- We will explore why **computeB** is different
 - > Takes almost 3X as much time as the others
- Collect experiment:
 - > **collect -p on -h +ecstall,on**

Demo

Function List

Sun Studio Analyzer [test.1.er]

File View Timeline Help

Find Text:

Functions	Threads	Source	DataObjects	VA	PA_L2	
User CPU (sec.)	User CPU (sec.)	E\$ Stall Cycles (sec.)	E\$ Stall Cycles (sec.)			Name
0.	34.234	0.	8.521			cache_trash
0.	16.411	0.	3.984			cache_trash_even
0.	17.822	0.	4.536			cache_trash_odd
0.	8.466	0.	0.			calladd
12.879	12.879	0.001	0.001			compute
12.179	12.179	0.	0.			computeA
34.234	34.234	8.521	8.521			computeB
12.088	12.088	0.	0.			computeC
12.699	12.699	0.001	0.001			computeD
12.179	12.179	0.	0.			computeE
3.903	8.466	0.	0.			computeF
12.048	12.048	0.	0.			computeG
12.219	12.219	0.	0.			computeH
12.689	12.689	0.	0.			computeI
0.	12.048	0.	0.			cond_global
0.	0.	0.	0.			cond_sleep_queue
0.	0.	0.	0.			cond_timedwait
0.	12.219	0.	0.			cond_timeout_global
0.	0.	0.	0.			cond_wait

Summary

Selected Object:

Name: addone

PC Address: 2:0x00004AC0

Size: 44

Source File: mttest.c

Object File: mttest

Load Object: <mttest>

Mangled Name:

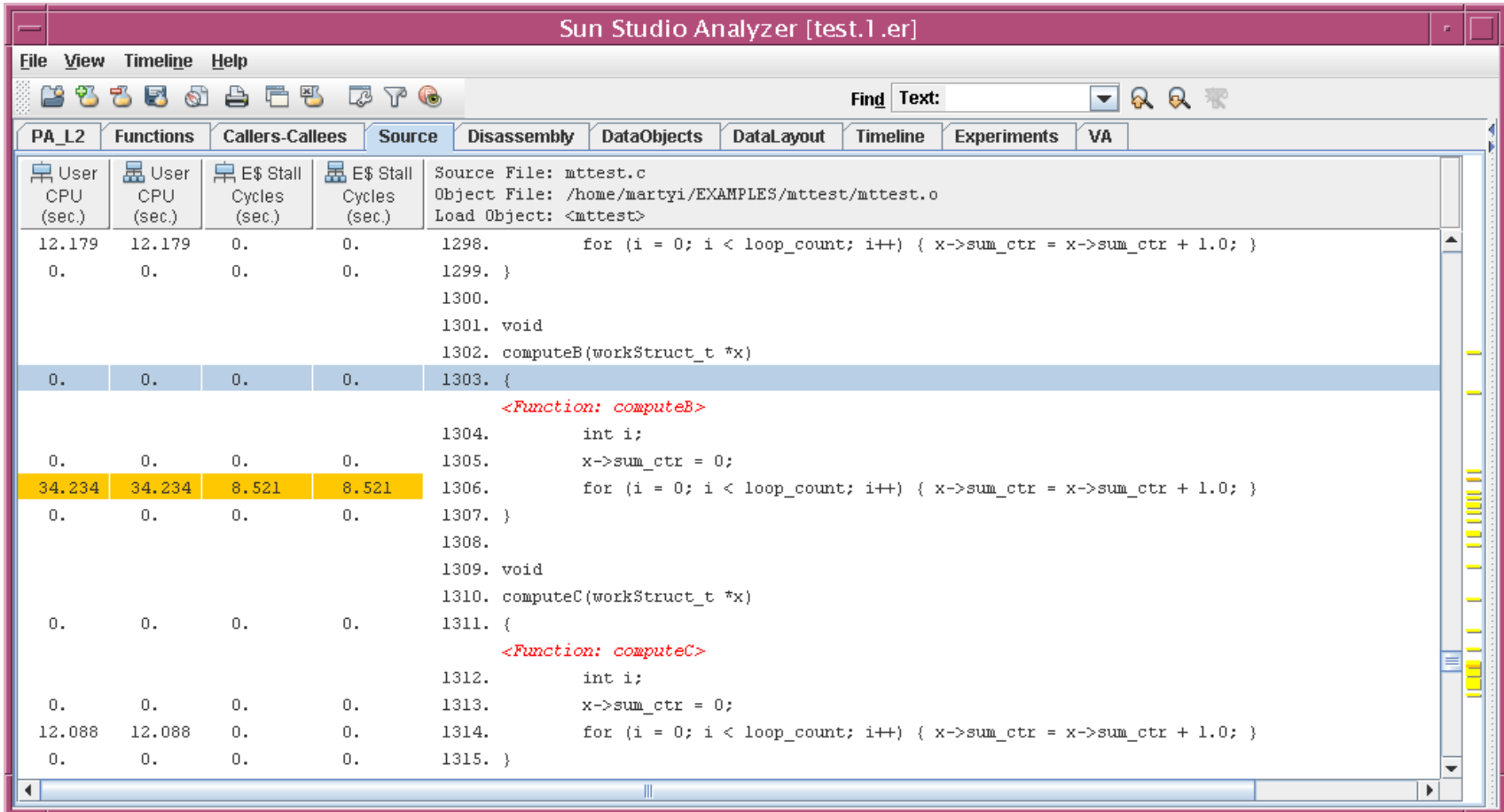
Aliases:

Metrics for Selected Object:

	Exclusive	Incl
User CPU:	4.563 (2.85%)	4.5
Wall:	0. (0. %)	0.
Total Thread:	9.016 (2.08%)	9.0
System CPU:	0.020 (2.17%)	0.0
Wait CPU:	4.353 (4.14%)	4.3
User Lock:	0. (0. %)	0.
Text Page Fault:	0. (0. %)	0.
Data Page Fault:	0. (0. %)	0.
Other Wait:	0.080 (0.08%)	0.0

Alphabetical (name) sort – note **computeB** vs. others

Source for compute*



The screenshot shows the Sun Studio Analyzer interface for a file named 'test.1.er'. The 'Source' tab is active, displaying the source code for 'computeB' and 'computeC'. The code is organized into columns for performance metrics: PA_L2, Functions, Callers-Callees, Source, Disassembly, DataObjects, DataLayout, Timeline, Experiments, and VA. The 'Source' column shows the following code:

```

Source File: mttest.c
Object File: /home/martyi/EXAMPLES/mttest/mttest.o
Load Object: <mttest>

1298.     for (i = 0; i < loop_count; i++) { x->sum_ctr = x->sum_ctr + 1.0; }
1299. }
1300.
1301. void
1302. computeB(workStruct_t *x)
1303. {
    <Function: computeB>
1304.     int i;
1305.     x->sum_ctr = 0;
1306.     for (i = 0; i < loop_count; i++) { x->sum_ctr = x->sum_ctr + 1.0; }
1307. }
1308.
1309. void
1310. computeC(workStruct_t *x)
1311. {
    <Function: computeC>
1312.     int i;
1313.     x->sum_ctr = 0;
1314.     for (i = 0; i < loop_count; i++) { x->sum_ctr = x->sum_ctr + 1.0; }
1315. }

```

The performance metrics for the highlighted lines are as follows:

PA_L2	Functions	Callers-Callees	Source	Disassembly	DataObjects	DataLayout	Timeline	Experiments	VA
12.179	12.179	0.	0.	1298.					
0.	0.	0.	0.	1299.					
0.	0.	0.	0.	1300.					
0.	0.	0.	0.	1301.					
0.	0.	0.	0.	1302.					
0.	0.	0.	0.	1303.					
0.	0.	0.	0.	1304.					
0.	0.	0.	0.	1305.					
34.234	34.234	8.521	8.521	1306.					
0.	0.	0.	0.	1307.					
0.	0.	0.	0.	1308.					
0.	0.	0.	0.	1309.					
0.	0.	0.	0.	1310.					
0.	0.	0.	0.	1311.					
0.	0.	0.	0.	1312.					
0.	0.	0.	0.	1313.					
12.088	12.088	0.	0.	1314.					
0.	0.	0.	0.	1315.					

Lines 1298,1306,1314 are *identical*
 But they perform *differently*

Function List Filtered on computeB

The screenshot shows the Sun Studio Analyzer interface with the 'Functions' tab selected. The function list is filtered to show only callers of 'computeB'. The 'computeB' function is highlighted in blue. The summary panel on the right shows details for the selected object 'computeB' and a table of metrics.

	Exclusive	Inclusive
User CPU:	34.234 (100.00%)	34.23
Wall:	0. (0. %)	0.
Total Thread:	77.494 (100.00%)	77.49
System CPU:	0.360 (100.00%)	0.36
Wait CPU:	41.679 (100.00%)	41.67
User Lock:	0. (0. %)	0.
Text Page Fault:	0. (0. %)	0.
Data Page Fault:	0.060 (100.00%)	0.06
Other Wait:	1.161 (100.00%)	1.16

Function list only shows callers of **computeB**
 (In this example **computeB** is a leaf function)

Threads, VA, and PA_L2 Tabs

The image shows three overlapping windows of Sun Studio Analyzer, each displaying a different view of performance data. The top window shows the 'Threads' tab with a table of threads. The middle window shows the 'VA' (Virtual Address) tab with a table of virtual addresses. The bottom window shows the 'PA_L2' (L2 cache lines) tab with a table of cache lines. A 'Summary' panel on the right of the bottom window shows metrics for the selected object.

Threads Tab Data:

User	E\$ Stall CPU (sec.)	E\$ Stall Cycles (sec.)	Name
34.234	8.521	<Total>	
10.277	2.790	Thread 12	
8.216	2.028	Thread 10	
8.196	1.956	Thread 11	
7.545	1.747	Thread 13	

VA Tab Data:

E\$ Stall Cycles (sec.)	Name
8.521	<Total>
8.415	PA_L2 Memory Object 6880
0.081	PA_L2 Memory Object 4777
0.024	<Unknown>

PA_L2 Tab Data:

E\$ Stall Cycles (sec.)	Name
8.521	<Total>
8.415	PA_L2 Memory Object 6880
0.081	PA_L2 Memory Object 4777
0.024	<Unknown>

Summary Panel Data:

Selected Object: PA_L2 Memory Object 6880

Metrics for Selected Object:

- E\$ Stall Cycles: 8.415 (98.76%)
- " count: 6311615605

Four threads, four virtual addresses, one cache line

Add VA Filter for One Address

Sun Studio Analyzer [test.1.er]

File View Timeline Help

Find Text: [] [] [] []

Functions Source Threads DataObjects VA PA_L2

Display Mode: Text Graphical

E\$ Stall Cycles (sec.)	Name
8.521	<Total>
2.759	VA Memory Object 4296177704
2.006	VA Memory Object 4296177688
1.924	VA Memory Object 4296177680
1.727	VA Memory Object 4296177696
0.081	VA Memory Object 4296043124
0.024	<Unknown>

Filter Data

Simple Advanced

Filter clause: (VA IN (4296177704)) [AND] [OR] [Set]

Specify filter:

```
(LEAF IN (44)) && (VA IN (4296177704))
```

Button to add && clause to filter

OK Apply Default Close Help

Will show only events in **computeB** referring to that one address

Look at VA and Threads again

The top screenshot shows the 'Threads' view. The table below displays the thread analysis results:

User CPU (sec.)	E\$ Stall Cycles (sec.)	Name
0.	2.759	<Total>
0.	2.759	Thread 12

The right-hand pane shows the 'Metrics for Selected Object' for Thread 12:

Exclusive	Value	Percentage
User CPU:	0.	(0. %)
Wall:	0.	(0. %)
Total Thread:	0.	(0. %)

The bottom screenshot shows the 'VA' (Virtual Address) view. The table below displays the VA analysis results:

E\$ Stall Cycles (sec.)	Name
2.759	<Total>
2.759	VA Memory Object 4296177704

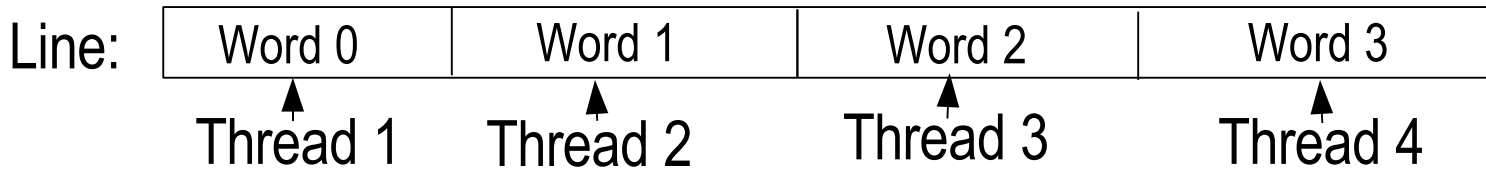
The right-hand pane shows the 'Metrics for Selected Object' for VA Memory Object 4296177704:

Data-derived	Value	Percentage
E\$ Stall Cycles:	2.759	(100.00%)
" count:	2069200513	

One thread per address; true of all four addresses

Diagnosis

- Most cache misses are on a single cache line
 - > Four threads get the misses
 - > Four addresses are referenced
 - > Each thread references only one virtual address



- > Write from one thread invalidates line for all others
- Classic manifestation of false sharing
 - > A notoriously difficult problem to spot
 - > In true sharing, multiple threads refer to each address

Potential Future Development

- Enhance data collection
 - > Support x86/x64 with instruction-based sampling
 - *Set up working group at this meeting?*
 - > Integrate configuration capture with data collection
- Improve the GUI and navigation
 - > Improve filtering grammar and syntax
 - > Other usability improvements
- Develop tuning strategy
 - > Systematic procedures for exploring problems

For more information

External Sun Studio Website

<http://developers.sun.com/sunstudio/>

External Sun Studio Performance Tools Website

http://developers.sun.com/sunstudio/overview/topics/analyzer_index.html

SC'96 paper on HW Counter Profiling

<http://portal.acm.org/citation.cfm?id=369028.369059&coll=portal&dl=ACM&CFID=33541981&CFTOKEN=50518735>

SC|03 paper on Dataspace Profiling

<http://www.sc-conference.org/sc2003/paperpdfs/pap182.pdf>

Solaris Application Programming by Darryl Gove

http://www.sun.com/books/catalog/solaris_app_programming.xml

Acknowledgments

- Nicolai Kosche, PAE
 - > Driving force for dataspace profiling enhancements
 - > Developed advanced techniques
 - > Invented term “DProfile” to refer to those techniques
- The Sun Studio Performance Analyzer team
 - > Made it all work

Thank You

Marty Itzkowitz, Analyzer Project Lead

Sun Microsystems Inc.

marty.itzkowitz@sun.com