The Ghost in the Machine Observing the Effects of Kernel Operation on Parallel Application Performance

Aroon Nataraj, Alan Morris, Allen Malony, Matthew Sottile, Pete Beckman[‡]

{anataraj, amorris, malony, matt}@cs.uoregon.edu Department of Computer and Information Science University of Oregon

beckman@mcs.anl.gov[‡] Mathematics and Computer Science Division Argonne National Laboratory





Talk Outline

- **The Problem**
 - What is Operating System / Runtime (OS/R) interference
 - Is it a problem?
 - Measurement question
- **The Solution**
 - KTAU and TAU performance systems
 - Fine-grained OS/App performance correlation
 - Noise-effect estimation technique
- **D** Evaluation
 - Demonstrating measurement and analysis of real OS noise
 - Investigating accuracy of noise-estimation analysis at scale
- **Conclusion**











Observing the Effects of Kernel Operation on Parallel Application Performance

What is the cause of imbalance?





Observing the Effects of Kernel Operation on Parallel Application Performance

What is the cause of imbalance?





Observing the Effects of Kernel Operation on Parallel Application Performance

Time lost to global noise

Is OS Noise a problem? How?

- Previous work has shown significant OS/R interference problems
 - Large variability in point-point communication latency [Mraz SC'94]
 - Measurement mismatched performance model [Petrini et al. SC'03]
 - Poor scaling performance of collectives [Jones et al. SC'03]
- □ Nature of noise matters
 - Theoretical modeling [Agarwal et al. HiPC'05]
 - Heavy-tailed and Bernoulli noise most detrimental
 - Noise Emulation [Beckman et al. CCJ'07]
 - Effect of noise on collectives
 - > Maximum noise duration determines effects
 - Large, rare noise-events are problematic







SC 2007



Noise Effects are Complex



Noise Effects are Complex

- Local noise dependent on the OS/R and its configuration, but global effects ...
 Depend on the underlying platform
 - Interconnect latency, timer resolution, TLB ...
- Depend on the OS/R configuration and noise sources
 Scheduling policy, interrupt frequency, daemons ...
- Depend on parallel application behavior
 - Synchronous communications, computational grain, load balance ...

- Can we measure the *global delay* an application experiences due to specific noise sources?
- □ Real application + real, existing noise => noise effect?



B ... bise sources

Our Approach and Contribution

□ General noise-effect estimation by direct measurement of application and OS

Contribution

- Isolate OS noise in application performance data
- Quantify global effects of noise on application
- Attribute the effects to specific noise sources
- Analyze application sensitivity



How do we measure?

- Application
 - TAU Performance System
 - Profiling and/or tracing of application events
 - Execution time, h/w performance counters ...
- Operating System
 - **O** KTAU
 - Profiling and tracing of system-level events
 - System calls, scheduling, interrupts, ...
- □ Integration
 - KTAU extended to allow fast access to system performance data
 - TAU captures OS data as counters stored with application events



Application w/ TAU

OS w/ KTAU

How do we measure?

- Application
 - TAU Performance System
 - Profiling and/or tracing of application events
 - Execution time, h/w performance counters ...
- Operating System
 - **O** KTAU
 - Profiling and tracing of system-level events
 - System calls, scheduling, interrupts, ...
- □ Integration
 - KTAU extended to allow fast access to system performance data
 - TAU captures OS data as counters stored with application events







SC 2007

The star

Observing the Effects of Kernel Operation In Parallel Application Performance





Application w/ TAU

User-level Double-Buffered Container

F. Spa

E. Sky



SC 2007

The stor

Observing the Effects of Kernel Operation In Parallel Application Performance





Application w/ TAU

User-level Double-Buffered Container

F. Spa

E. Stel



SC 2007

The stor

Observing the Effects of Kernel Operation In Parallel Application Performance





Application w/ TAU

User-level Double-Buffered Container

Ster.

E. Stel



SC 2007

The stor

Observing the Effects of Kernel Operation In Parallel Application Performance





Application w/ TAU

User-level Double-Buffered Container

F. Spa

E. Stel



SC 2007

The stor

Observing the Effects of Kernel Operation In Parallel Application Performance





Application w/ TAU

User-level Double-Buffered Container

R. SKA

E. Sky



Per-Process Virtualized OS Counters

SC 2007

Observing the Effects of Kernel Operation In Parallel Application Performance





No Daemon or System Call needed!

E. Sky

How do we analyze?

- **Trace** information
 - Application event trace with OS noise counters
 - Presence of OS noise has affected the timing of events
- **Timeline** Approximation
 - Reason about timing of events that may have occurred in the absence of noise • Remove time-duration of noise and adjust event timestamps

 - Must be careful to maintain constraints in event ordering
 - > E.g. A recv cannot end before the corresponding send
- Delay due to the global noise-effect (Accumulated Noise)
 - Difference in end times between measured and approximated timeline
- **Two main cases**
 - Noise Propagation; Noise Absorption
 - See paper for other cases





Local Noise = 100s

Local Noise = 0s

Rank 2















Rank 1 Local Noise = 100s

Rank 2 *Local* Noise = 0s



SC 2007





Rank 1 Sends; Rank 2 Receives



Accumulated Noise = Last <u>Measured</u> Timestamp - Last <u>Approximated</u> Timestamp



Observing the Effects of Kernel Operation on Parallel Application Performance

SC 2007



Rank 1 Sends; Rank 2 Receives



Propagation - Rank 2 picks up 100s of delay even though its local-noise is 0!



SC 2007







SC 2007







SC 2007







Absorption - Rank 2's Acc. Noise is only 50s, but local noise was 100s. It loses 50s!



SC 2007

Prior Work on Trace-based Timeline Approximation

- □ Wolf, Malony, Shende and Morris | HPCC'06 • Context: Measurement perturbation compensation.
- □ Sottile, Chandu and Bader | IPDPS'06
 - Trace-based simulation
 - Inject artificial noise into existing application traces and adjust timestamps
 - Reverse of our current work we remove existing, real noise effects from trace.



Noise-Effect Estimation Demonstrated

- Demonstrate estimation of delay caused in <u>application/benchmark</u> by <u>real noise</u>
- □ Platform: 32 (2x2) Opterons (P=128); GigE; Linux w/ KTAU; SDSC
- Application : Sweep 3D (kernel representative of ASC applications)
 - Repeating phases of Send/Recv followed by Allreduce
 - Problem 650³, 15 iterations, MPI based
- □ Scaling
 - O Strong
 - O P=32, P=128
- Instrumentation / Measurement
 - TAU Tracing of MPI events
 - Associating KTAU OS Metrics with the events
- □ Analysis
 - Run Trace analysis (described earlier) to calculate delay due to noise

SC 2007



Noise Sources and Metrics

OS Noise Sources

- global timer interrupt keeps time & timers, intervals (10, 4, 1 msec) • local timer interrupt - update process times (scheduling), every cpu/core • preemptive schedule - duration preempted

Image: Metrics

- Accumulated Noise
 - Estimate of delay due to global noise effects, in secs
 - > By how much time would the application have run faster w/o noise?
- Noise Amplification Ratio
 - Accumulated Noise / Local Noise
 - How much was noise amplified? How much was absorbed?
- Lesser means better for both metrics
- □ Metrics calculated seperately for each *noise source* and *combined noise*







|) | 100 | 1. | 20 | 140 |
|-------|------|--------------------|-------|-----|
| | | | | |
| | | | | _ |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | _ |
| ***** | **** | ~X~~~~ | ***** | |
| | | | | _ |
| | | | | |
| | | | | _ |
| | | | | |
| | | | | |
| | | P=32 Default Run — | | |
| | | P-32 Default Bun | | |





P=32 Default Run —+--P=128 Default Run —×--

Local noise *proportional* to runtime. Longer the run, more timer interrupts. Shorter the run, fewer are expected.



P=32 Default Run —+--P=128 Default Run —×--

Local noise *proportional* to runtime. Longer the run, more timer interrupts. Shorter the run, fewer are expected.
Overall Accumulated Noise - Effect of Scaling



Overall Accumulated Noise - Effect of Scaling



Overall Noise Amplification Ratio - Effect of Scaling



Observing the Effects of Kernel Operation on Parallel Application Performance









Confirms that scaling is increasing global noise-effect. How do the noise sources contribute?

Observing the Effects of Kernel Operation on Parallel Application Performance

Noise Sources -- Accumulated Noise -- P=128



| | Comb sche local-t global-t | ined Acc. Noise – dule Acc. Noise – imer Acc. Noise – imer Acc. Noise – | -+ -× -₩ - |
|--|--|--|---------------------|
| | | | |
| _{╋┿╋╋} ┿┿┿┿┿┿┿ | _{┺┺} ┿┿┿┿┿┿┿┿┿ | _{┽┼┿} ╪╪╪╪╪╪╪╪╪╪╪╪ | |
| **** | ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ | | |
| | | | |
| XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX | ````````````````````````````````````` | ĸĸĸŴĸŴ ^{Ŕĸ} ĸĸĸĸ ŴŴŴŴŴ | |
| | | | |
| 0 | 100 | 120 | 140 |

Noise Sources -- Accumulated Noise -- P=128



| | Combi sche local-ti global-ti | ined Acc. Noise – dule Acc. Noise – imer Acc. Noise – imer Acc. Noise – | -+ -★ |
|-----------------------|--|--|----------|
| | | | |
| _{╋┿╋╋} ┿┿┿┿┿ | [┡] ╈╪╪╪╪╪╪╪╪╪╪╪╪ | ┝┼┼ [┲] ╇╍ _╋ ╪╋╋╪╋┿┿┿┿┿┿┿ | |
| **** | ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~ | | |
| nant so | ource of nois | se. | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| 0 | 100 | 120 | 140 |

Noise Sources -- Accumulated Noise -- <u>P=128</u>



| | Combine schedu local-tim global-tim | ed Acc le Acc er Acc er Acc er Acc | . Noise . Noise . Noise . Noise | + | _ |
|--|--|--|--|---|----|
| e can be ren | noved. | | | | _ |
| | | - | | | |
| | | | | | _ |
| <u>╃</u> ╾┾╍ <mark>┼╌┼╍┼╍┼╍┼╱╧┑</mark> ╌┟╍╎╴ | ₽ [╋] ╅┿┿┽┿╋ <mark>┥┥┥┥</mark> ┿ | ┿ _{╋╋} ┽╄ _{╋╋╋╋} | ╒╇╈┼┿┽┽┼ | | |
| | | | | | _ |
| ****** | ***** | ×***** | XXXXX | | |
| nant source | of noise | • | | | _ |
| | | | | | |
| | | | | | _ |
| | | | | | |
| ₭₦[₩]₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩₩ | Karanana Karananan Karananan | Ky HOK | **** | | _ |
| | | | | | |
| | | | | | _ |
| | | | | | |
| U 10 | 00 | 12 | 20 | | 40 |

Noise Sources -- Accumulated Noise -- <u>P=128</u>



Lets make a simple change that affects scheduling. **Pin** the ranks to the processors. What happens?

SC 2007

| | Combine schedu local-time global-time | ed Acc le Acc er Acc er Acc | . Noise . Noise . Noise . Noise . Noise | + | - |
|---|--|--------------------------------------|---|-------|----|
| e can be rer | noved. | | | | _ |
| | | | | | |
| | | | | | - |
| ┾ ┼ ┿┾ _{╋┿} ┿┿┿┿┿┿┿┿┿┿┿ | ╒╇╅┿┽┽┿╋╅┿┿┿┿ | ╘╪╪┼╘╕╶╽╶ ╽ | ╇╈┽┿┽┽ | | |
| | | | | | - |
| ***** | ***** | ****** | XXXXXX | | |
| nant source | of noise. | | | | - |
| | | | | | |
| | | | | | - |
| | | | | | |
| ************************************** | ******* | KANAKAKAKA | ***** | | - |
| | | | | | |
| | | | | | - |
| | | | | | |
| 0 1 | 00 | 12 | 20 | 1 | 40 |
| | | | | | |
| offoata aclos | duling | | | | |

Overall Accumulated Noise - Pinned - Effect of Strong Scaling



Observing the Effects of Kernel Operation on Parallel Application Performance



| 0 | 100 | 120 | 140 |
|-------|-------|-------|-----|
| ***** | ***** | ***** | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

Overall Accumulated Noise - Pinned - Effect of Strong Scaling





| 0 | 100 | 120 | |
|-------|-------|---|--|
| ***** | ***** | >>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>> | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | P=1 | P=128 Pinned Run —× | |

Overall Accumulated Noise - Pinned - Effect of Strong Scaling





| 0 | 100 | 120 | 140 |
|-------|-------|-------|-----|
| ***** | ***** | ***** | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |





Lets look at the *different noise sources* again...

Observing the Effects of Kernel Operation on Parallel Application Performance

The Noise Sources -- Pinned -- P=128



Observing the Effects of Kernel Operation on Parallel Application Performance

| | Combined Acc schedule Acc local-timer Acc global-timer Acc | c. Noise → c. Noise → c. Noise → c. Noise → | |
|-----|---|--|-----|
| | | | |
| | | | |
| | | | |
| | | | |
| d | | | |
| | | | |
| | | global-timer | |
| 0 1 | 00 1 | 20 | 140 |



Small Noise Effects

- □ Magnitude of the Accumulated Noise 1.3 secs • Represents approx. 1% of runtime (132 secs)
- □ Small cluster (32 nodes) + Slow interconnect (GigE) O => Small global-noise effect
- □ For large OS noise related slowdowns
 - Larger scales
 - Fast interconnect

Global-Noise Estimation still detected and revealed interesting noise features



Accuracy of Noise-Estimation Analysis at Scale

- □ How accurate is the "Accumulated Noise" value provided by analysis?
- Methodology
 - Take (relatively) noise-less platform (BG/L); Inject noise (Selfish Suite/ANL) • Run parallel application without noise & then with injected noise

 - Perform trace analysis to provide the *accumulated noise* estimate
 - How far is calculated *accumulated noise* value from actual delay?
- □ Simple BSP Benchmark (Bulk Synchronous Processing)
 - Repeated phases of computation & collective communication
 - Inputs: Scaling type, No. of Phases, Type of Collective, No. of Nodes
 - Used: Strong scaling, 10000 phases, Barrier, 32 to 2048 Nodes



% Dilation in Runtime



% Dilation in Runtime



Observing the Effects of Kernel Operation on Parallel Application Performance











□ Compare access costs (in cycles) to **O** PAPI counters • *no-op"* /proc call

| # Metrics | OS Metric Access | /proc Access | PAPI Access |
|------------------|-------------------------|--------------|-------------|
| 1 | 192 | 1150 | 248 |
| 2 | 272 | _ | 304 |
| 3 | 288 | - | |
| 4 | 400 | _ | _ |

Image: Metric Access order of magnitude faster than /proc access

Image: Metric Access comparable to PAPI h/w counter access



□ Compare access costs (in cycles) to **O** PAPI counters • *no-op"* /proc call

| # Metrics | OS Metric Access | /proc Access | PAPI Access |
|------------------|-------------------------|--------------|-------------|
| 1 | 192 | 1150 | 248 |
| 2 | 272 | _ | 304 |
| 3 | 288 | _ | _ |
| 4 | 400 | | _ |

Image: Metric Access order of magnitude faster than /proc access

Image: Metric Access comparable to PAPI h/w counter access



□ Compare access costs (in cycles) to **O** PAPI counters • *no-op"* /proc call

| # Metrics | OS Metric Access | /proc Access | PAPI Access |
|------------------|-------------------------|--------------|-------------|
| 1 | 192 | 1150 | 248 |
| 2 | 272 | _ | 304 |
| 3 | 288 | _ | _ |
| 4 | 400 | _ | _ |

Image: Metric Access order of magnitude faster than /proc access

Image: Metric Access comparable to PAPI h/w counter access

What is the Measurement Perturbation?

- Image: Measure Overall Perturbation of NPB LU under multiple configurations • Configuration: base
 - > No instrumentation in application or OS
 - Configuration: *ktau-tau-metrics*
 - > TAU MPI tracing
 - Tracking 4 OS metrics for each MPI event

NPB LU Class C on 16 Nodes

| Configuration | base | ktau-tau-metrics |
|--------------------|--------|------------------|
| Minimum Exec. Time | 475.04 | 479.66 |
| % Min. Slowdown | | 0.97 |

• Performing OS Metric access + TAU MPI Tracing < 1% Perturbation



Conclusion

- General measurement technique that estimates the delay caused by direct OS noise effects on a parallel message passing application • Integrated OS / Application performance measurement

 - Trace timeline approximation
- Demonstrate its use on a Linux cluster to measure effects of *real, existing noise* Evaluate accuracy of analysis at scale (2048 nodes)
- • 0.5% to 4.5% estimation error over varying noise levels and computational grain
- Questions in the HPC community
 - OS/R suites for large-scale platforms Light-weight or Full-featured?
 - Can applications be changed to be less noise-sensitive?
 - DOE FAST-OS project created to investigate OS issues, including noise
 - Integrated (OS/Application) measurement and noise analysis techniques can aid in answering some of the questions



Combining Noise Sources -- Accum



Observing the Effects of Kernel Operation on Parallel Application Performance

| | Maine | D_1 | 20 |
|---------|---------|--------------------------------|----|
| iulatea | Noise . | $-\underline{P}=\underline{I}$ | 20 |

| | Combined Acc schedule Acc local-timer Acc global-timer Acc | Noise — Noise — Noise ———— Noise ———— Noise ——— | |
|--|---|---|----|
| | | | |
| ╅╺╪┥┑ ┿┿┿┿┿┿┿┿┿┿┿┿┿ | ┲╬ _{╈┿┿┿┿} ╌┲╬╅┿┿┿┿┿┿┿ | <mark>₽[╋]₩₽[╋]╋╋╋╋</mark> | |
| ****** | XxxxxXxxxxXXXXXXXXXXXXXXXXXXXXXXXXXXXX | XXXXXXX | |
| | | | |
| XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX | | | |
| | | | |
| | | | |
| 0 1 | 00 12 | 20 1 <i>i</i> | 40 |

Appendix A

Combining Noise Sources -- Accum



Observing the Effects of Kernel Operation on Parallel Application Performance

|] | atod | Noise | 5. | D = | 128 |
|-----|------|-------|------|------------|-----|
| ini | uleu | TUDBE | C.A. | | |

| | Combined Acc schedule Acc local-timer Acc global-timer Acc | Noise | |
|--|---|--|----|
| | | | |
| ୶ୄ୲୶ ୄୄୄୄ୶ _୳ ୄ୶୶୶୶୶୶୶୶୶୶୶୶୶୶ | ╒╇╕╾┼┼┼┲╇╕╴┼┼┼┼┲┿╕╪╋╌╌╌ | <mark>₽[₽]₩₽[₽]₩</mark> ₩₩ | |
| ***** | Xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx | XXXXX | |
| st source of | noise. | | |
| | | | |
| | | | |
| <u>*************************************</u> | ******* | K ^W WWWWW | |
| | | | |
| | | | |
| 0 10 | 00 1 | 20 14 | 10 |

Appendix A





Combining Noise Sources -- Accum



Why?

Observing the Effects of Kernel Operation on Parallel Application Performance

SC 2007

|] | atod | Noise | 5. | D = | 128 |
|-----|------|-------|------|------------|-----|
| ini | uleu | TUDBE | C.A. | | |

| | Combined Acc schedule Acc local-timer Acc global-timer Acc | Noise | |
|--|---|--|----|
| | | | |
| ୶ୄ୲୶ ୄୄୄୄ୶ _୳ ୄ୶୶୶୶୶୶୶୶୶୶୶୶୶୶ | ╒╇╕╾┼┼┼┲╇╕╴┼┼┼┼┲┿╕╪╋╌┝╕ | <mark>₽[₽]₩₽[₽]₩</mark> ₩₩ | |
| ***** | Xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx | XXXXX | |
| st source of | noise. | | |
| | | | |
| | | | |
| <u>*************************************</u> | ******* | K ^W WWWWW | |
| | | | |
| | | | |
| 0 10 | 00 1 | 20 14 | 10 |

Appendix A

Global Noise Estimation - Noise Combination

Process 1 Sends; Process 2 Receives



Approximated Timeline

SC 2007


Process 1 Sends; Process 2 Receives







Process 1 Sends; Process 2 Receives



Approximated Timeline

SC 2007



Process 1 Sends; Process 2 Receives





SC 2007



Process 1 Sends; Process 2 Receives



Approximated Timeline



Process 1 Sends; Process 2 Receives





SC 2007







Related Work in Noise Measurement

□ Petrini et al. SC'03 - Microbenchmark, simulation, modeling to *close the loop* • Specific to application (SAGE). Accurate model difficult to produce. Gioiosa et al. ISSIPIT'04 - Microbenchmark & measurement • Identify noise sources using OProfile sampling. Quantifies only *local-noise*. □ Agarwal et al. HiPC'05 - Theoretical Modeling under different distributions • Assumptions include: Balanced Load, Stationary, Balanced Noise, Identical noise □ Beckman et al. CLUSTER'06 (Emulation), Sottile et al. IPDPS'06 (Simulation) • Injected artificial noise at runtime into micro-benchmarks to understand effects • Modify application traces by adding artificial noise.

Profiling LU Application using KTAU OS Metrics





and pinning No irq balance

std. dev 0.19 mean 0 | n,c,t 0,C,O n, i, i 1, C, O 4F-6 n,r,t 2,C,O 0 n,c,t 3,C,O n,c,t 4,C,O ú n,c,t 5,C,O n,c,t 6,C,O 0 7 7E-6 n, i, 17, 0, 0 0.758 n, c, t 8, C, O 0 n,c,t 9,C,O 2.2E-6 n,c,t 10,0,0 5.9E-6 n,c,t 11,0,0 n, i, i 12, 0, 0 1.CE-5 | n,c,t 13,0.0 2.1E-6 n,c,t 14,0,0 0 n,c,t 15,0,0 0.78 n,c,t 16,0,0 0 n, c, t 17, 0, 0 4 6E-6 n,c,t 18,0,0 6.7E-6 n,c,t 19,0.0 0.755 n,c,t 20,0,0 2.1E-6 n,c,t 21,0,0 n, c, t 22, 0, 0 n,c,t 23,0,0 0 n,c,t 24,0.0 0.757 n,c,t 25,0,0 n,c,t 26,0,0 0 n, c, t 27, 0, 0 Ú. n,c,t 28,0,0 0 5.5E-6 | n,c,t 29,0.0 0 n,c,t 30,0,0 n,c,t 31,0,0

global timer interrupt



Appendix C

Profiling LU Application using KTAU OS Metrics



global timer interrupt







and pinning No irq balance

Appendix C













Appendix C



Appendix C

application events.

Acknowledgments

- □ San Diego Supercomputing Center (SDSC) & Don Thorp O Access to Opteron cluster
- Argonne National Laboratory
 - \circ Access to BG/L and other compute resources
- **DOE FAST-OS Project**
 - Funding as part of the joint ZeptoOS project between ANL and UO

