



U.S. DEPARTMENT OF
ENERGY



**UNIVERSITY OF
CALIFORNIA**



BERKELEY LAB
LAWRENCE BERKELEY NATIONAL LABORATORY



Benchmarking C++ Code

Bryce Adelstein Lelbach aka wash <balelbach@lbl.gov>

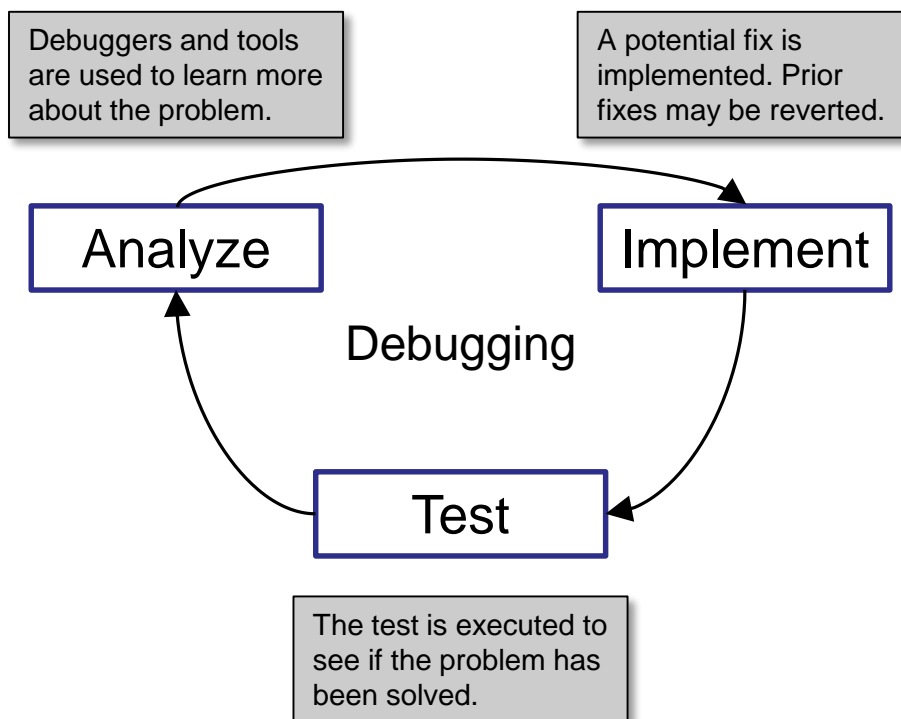
Computer Architecture Group, Computing Research Division

CppCon 2015

The Problem with Performance

Problem: **Code seg faults**

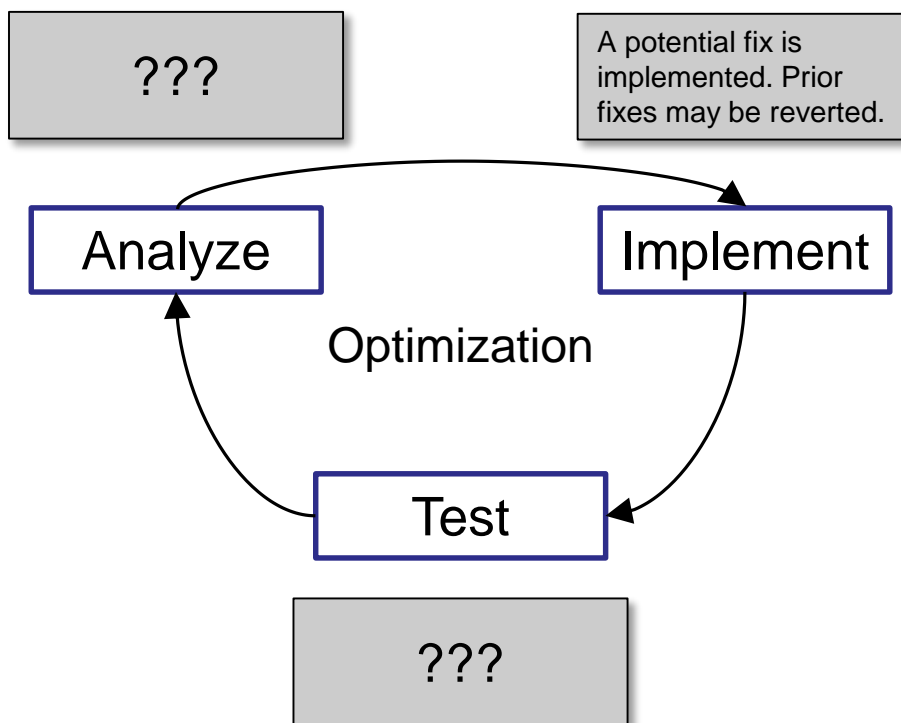
- We solve this type of problem with an iterative workflow.
- We know when we're done; we can easily get a "yes/no" answer during the testing phase.
 - Usually, there's no random error when testing for this type of problem (excluding race conditions).



The Problem with Performance

Problem: **Code is slow**

- Producing a “yes/no” answer during the testing phase is more difficult.
 - Performance is not a Boolean quantity.
 - It is often unclear when the problem is “fixed”.
 - You never really finish optimizing.
 - Performance data is subject to **random error** due to natural variability.



What is Performance?

How do we define performance, anyways?

- Not “fast”, but “fast enough”.
- Real-world metrics:
 - Ex: simulation-years/day
- Roofline:
 - Ex: FLOP/s
- Deadline:
 - Ex: takes 50 milliseconds

You need to be able to come up with meaningful definitions for performance.

Sources of Error

Observational Error: The difference between what you measure and the true result.

- Random Error: Errors caused by natural variance.
- Systemic Error: Errors caused by an inaccuracy – usually constant or proportional to the true result.

Observational error is unavoidable. Meaningful performance analysis **must** account for error.

- E.g. statistical testing approach

Variance

Computers can reproduce answers, not performance.

- Hardware jitter
 - Instruction pipelines: The pipeline fill level has an effect on the execution time for one instruction.
 - Difference in CPU/memory bus clock cycles: The CPU clock cycle is different from the memory bus clock speed. Your CPU sometimes has to wait for the synchronization of memory accesses.
 - CPU frequency scaling and power management: These features cause heterogeneities in processing power.
 - Shared hardware caches: Caches shared between multiple cores/threads are subject to variance due to concurrent use.

Source: <http://www.chronox.de/jent/doc/CPU-Jitter-NPTRNG.html>

Variance

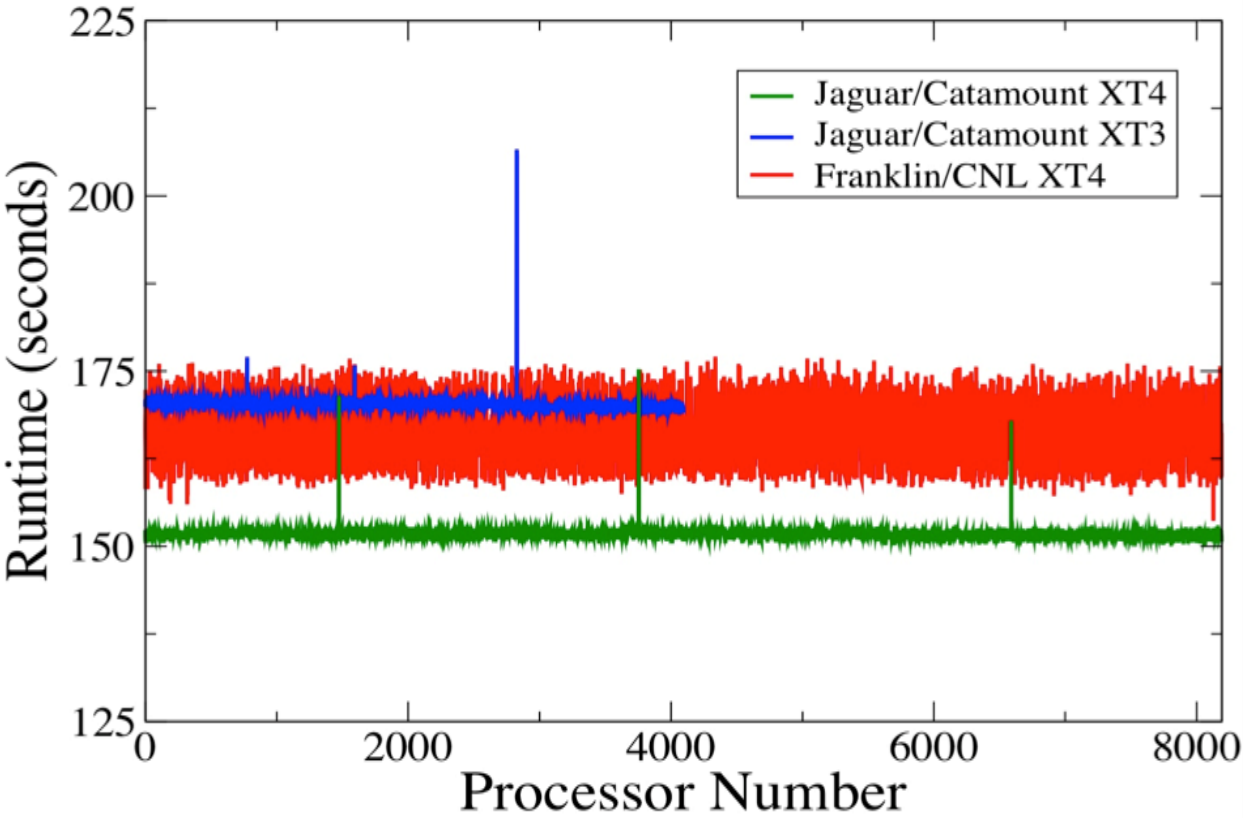
Computers can reproduce answers, not performance.

- Larger memory segments may have variance in access times due to physical distance from the CPU.
- Additionally, OS activities can cause non-determinism.
 - Some hardware interrupts require OS handling immediately after delivery.
 - Migration of non-pinned processes can affect the performance of CPU heuristics.

Observer Effect: all forms of instrumentation change the results.

Source: <http://www.chronox.de/jent/doc/CPU-Jitter-NPTRNG.html>

Variance



Source: Cy Chan, John Bachan



Statistical Best Practices

Statistical Best Practices

Statistics: A great way to lie to yourself.

Statistical Best Practices

Statistics: A way to extract conclusions from your data

Statistical Best Practices

Statistics: The science of data...

collection

analysis

interpretation

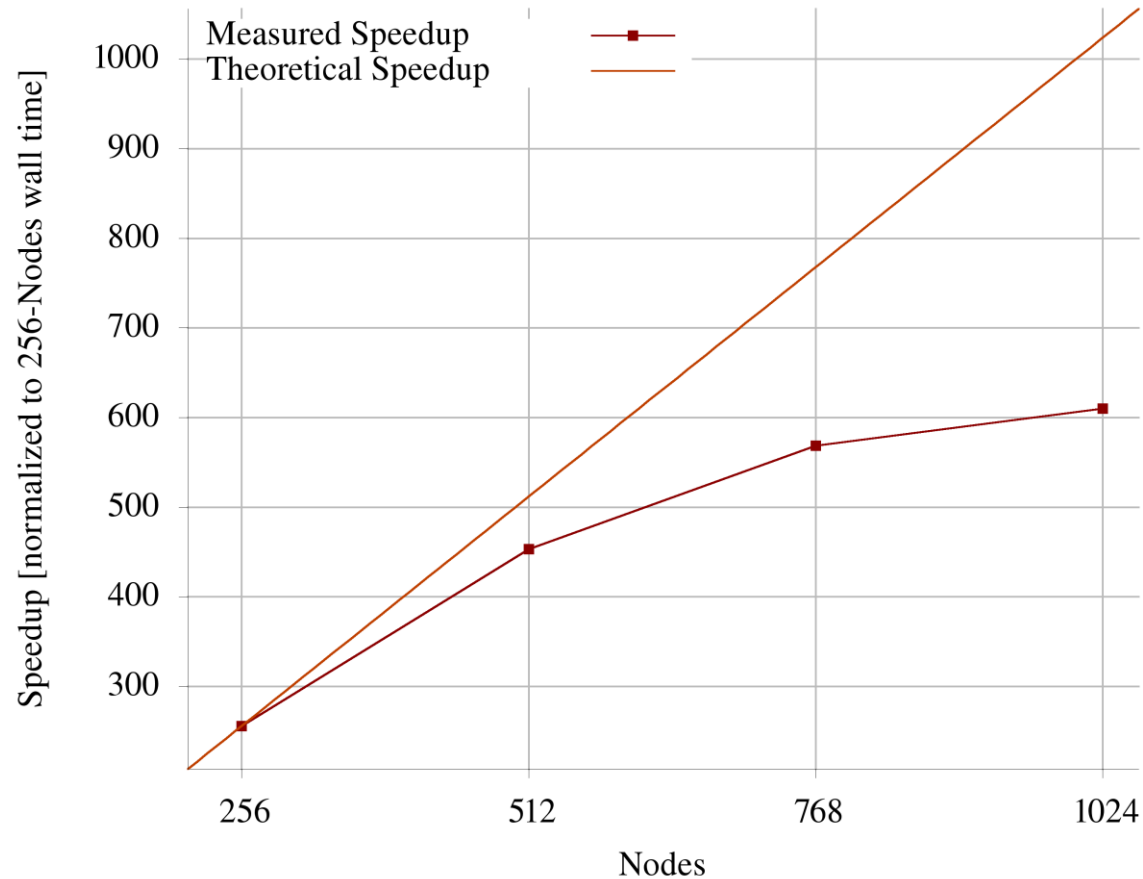
presentation

Case Study: CFD AMR Scaling

Statistical Best Practices

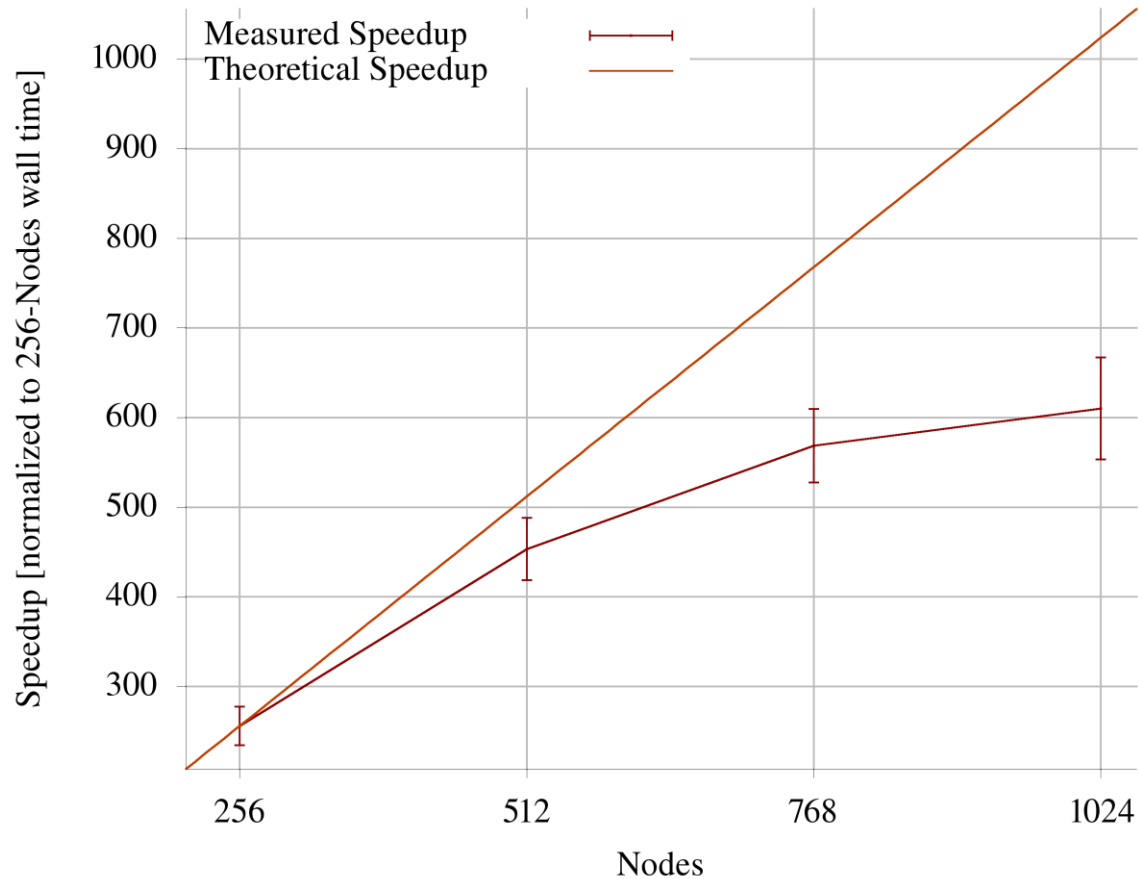
Case Study: CFD AMR Scaling

AMR Test, Strong-Scaling



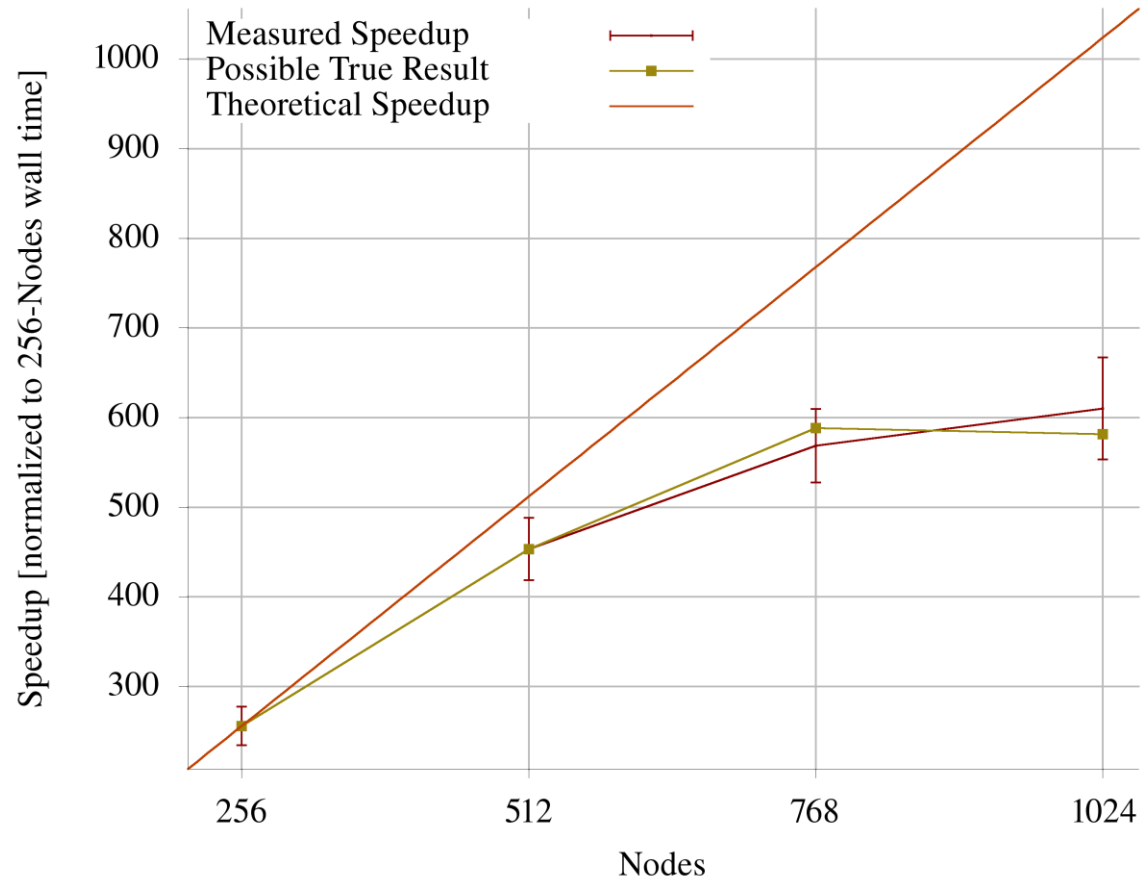
Case Study: CFD AMR Scaling

AMR Test, Strong-Scaling (with uncertainty)



Case Study: CFD AMR Scaling

AMR Test, Strong-Scaling (with uncertainty)



Statistical Best Practices

Process:

- Form a hypothesis: how do you expect performance to change?
- Come up with a test to determine if your hypothesis is right.
- Gather data.
- Statistically analyze data.
- Draw conclusions.

Statistical Best Practices

Come up with a test to determine if your hypothesis is right.

- Identify independent/dependent/control variables.
- Determine what relevant metric you'll use (metric will be derived from dependent variables).
- Consider the assumptions you're making:
 - Assumptions about independence of variables.
 - Assumptions about distribution of samples.
 - Usually we assume a normal distribution.

Gathering Data

Amortizing: When measuring “small” events, we often measure by amortization to reduce the observer effect.

- E.g. time an N-iteration for loop and divide by N to get the amortized time per iteration.

```
high_resolution_timer t; // Start timing.  
  
for (std::size_t i = 0; i < N; ++i)  
    A[i] = A[i] + B[i] * C[i];  
  
double time_per_iteration = t.elapsed() / N;
```

- We treat this as one sample, not N samples.

Gathering Data

Sampling: Each independent measurement we take is a sample.

- Samples are representative of the “population” (AKA the true performance).
- Our goal is to gather samples in sufficient quantity and quality to be representative of the population.

It's crucial to both sample within one execution of the test and across multiple executions of the test.

- Gathering data across multiple executions gives a better representation of system noise.

Gathering Data

Running “hot” vs “cold”.

- Often, you need to make sure that both your test as a whole (e.g. each execution), and the particular region your measuring (e.g. each sample) are not running “cold” on the CPU.
 - I/O, caching and branch prediction may be off if you’re running cold.

You can do this by doing some warmup executions/runs before you start measuring.

- E.g. don’t measure first execution or first few iterations.

Uncertainty

Uncertainty: representation of the amount of error in a certain measurement.

- Instrument uncertainty: the inherent amount of uncertainty in an instrument.
 - Ex: if your clock ticks in microseconds, it would have an instrument uncertainty of +/- 500 nanoseconds (1/2th the unit of measurement).
- The sample standard deviation of a set of samples is a frequently used method for estimating the uncertainty of the average of the samples.

When dealing with derived metrics that use averaged data, you can formulate a derived uncertainty based on the uncertainties of the averaged data.

Uncertainty

Given **uncorrelated** averaged data A and B with standard deviations σ_A and σ_B , and constants a and b .

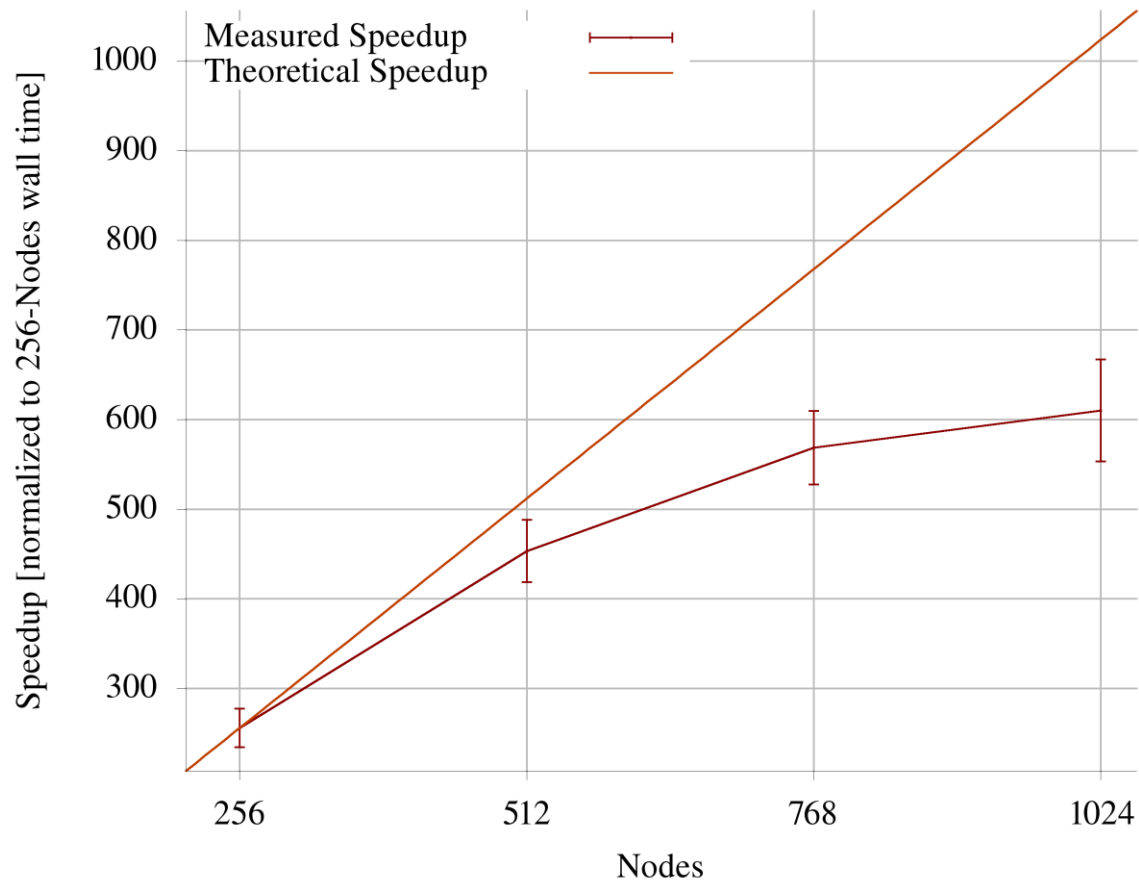
Function	Standard Deviation
$f = aA$	$\sigma_f = a\sigma_A$
$f = aA \pm bB$	$\sigma_f = \sqrt{a^2\sigma_A^2 + b^2\sigma_B^2}$
$f = AB$ or $f = A/B$	$\sigma_f \approx f \sqrt{\left(\frac{\sigma_A}{A}\right)^2 + \left(\frac{\sigma_B}{B}\right)^2}$

Case Study: CFD AMR Scaling

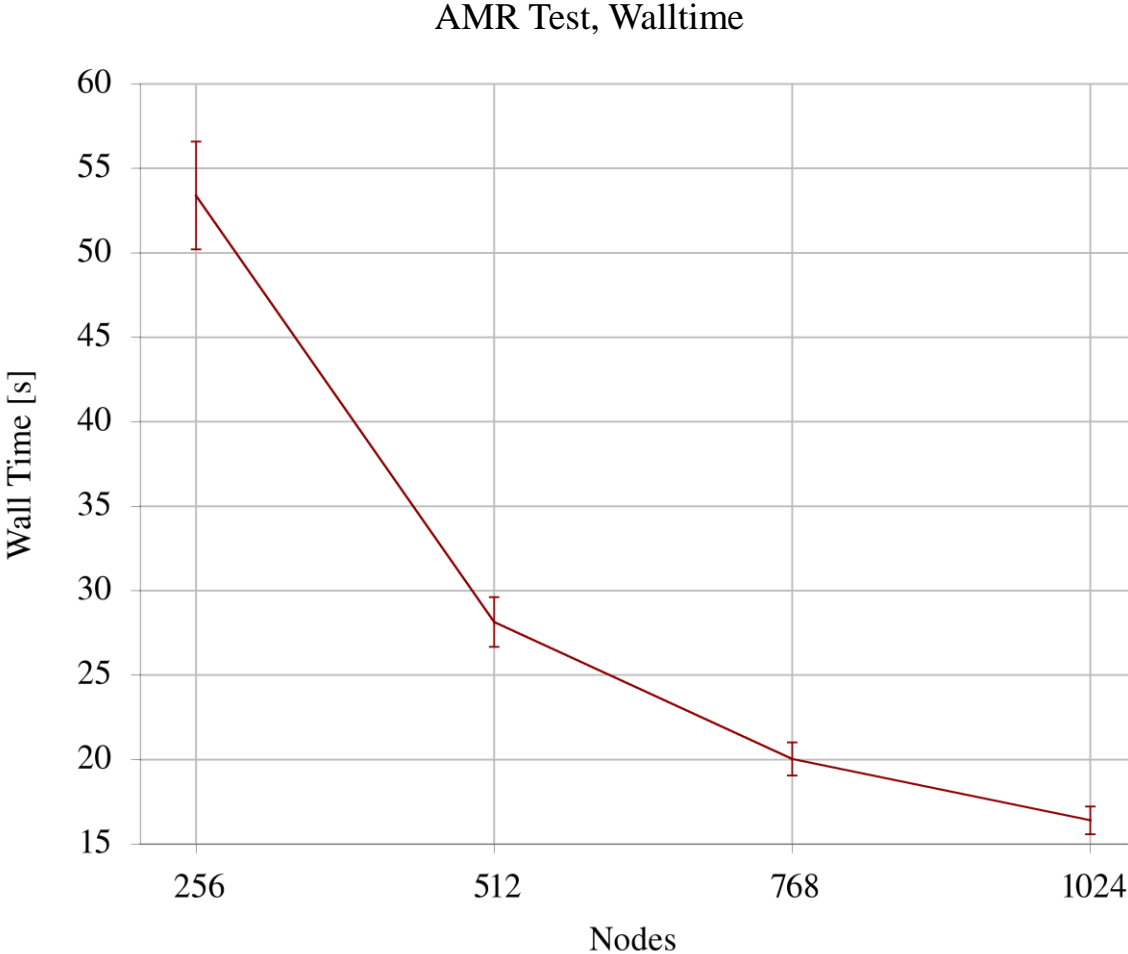
Statistical Best Practices

Case Study: CFD AMR Scaling

AMR Test, Strong-Scaling (with uncertainty)



Case Study: CFD AMR Scaling



Example: Boost.Accumulators

Statistical Best Practices

Example: Boost.Accumulators

“Boost.Accumulators provides accumulators to which numbers can be added to get, for example, the mean or the **standard deviation**.”

The Boost C++ Libraries, Boris Schäling

Example: Boost.Accumulators

```
using namespace boost::accumulators;

int main()
{
    accumulator_set<
        double, stats<tag::count, tag::mean, tag::median, tag::variance>
    > acc;

    acc(42);

    // ... Accumulate data ...

    auto stdev = std::sqrt(variance(acc));

    std::cout << "Mean:    " << mean(acc) << "\n"
               << "Median:  " << median(acc) << "\n"
               << "Stdev:   " << stdev << "\n";
}
```

Example: Boost.Accumulators

On Jan 17, 2012, at 3:38 PM, Victor Yankee wrote:

- > How can I calculate the Sample Standard Deviation over a `std::vector` of doubles using accumulators?
- >
- > Or is there a faster way in boost math or some such?

Google is your friend (second hit for "boost accumulator standard deviation"):

<http://stackoverflow.com/questions/7616511/calculate-mean-and-standard-deviation-from-a-vector-of-samples-in-c-using-boos>
and <http://stackoverflow.com/questions/4316716/is-it-possible-to-use-boost-accumulators-with-vectors>

```
accumulator_set<double, stats<tag::variance> > acc;  
for_each(a_vec.begin(), a_vec.end(), bind<void>(ref(acc), _1));  
  
cout << mean(acc) << endl;  
cout << sqrt(variance(acc)) << endl;
```

-- Marshall

Marshall Clow Idio Software <mailto:mclow.lists_at_[hidden]>

A.D. 1517: Martin Luther nails his 95 Theses to the church door and is promptly moderated down to (-1, Flamebait).

-- Yu Suzuki

Example: Boost.Accumulators

But this is the POPULATION standard deviation (variance divided by N) I think.

What I was asking for was how to calculate the SAMPLE standard deviation (variance divided by $N-1$).

Thanks,
Vic

Example: Boost.Accumulators

Two different forms of standard deviation

- Uncorrected, takes the standard deviation of an entire population:

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

- Corrected, takes the standard deviation of a sample of a population:

$$\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2}$$

Example: Boost.Accumulators

```
using namespace boost::accumulators;

int main()
{
    accumulator_set<
        double, stats<tag::count, tag::mean, tag::median, tag::variance>
    > acc;

    acc(42);

    // ... Accumulate data ...

    auto n = count(acc);
    auto stdev = std::sqrt(variance(acc) * (n / (n - 1.0)));

    std::cout << "Mean:    " << mean(acc) << "\n"
               << "Median:  " << median(acc) << "\n"
               << "Stdev:   " << stdev << "\n";
}
```

Gathering Data

Process for collecting good data:

- Take individual measurements in your code. Use amortization if relevant.
- Accumulate multiple measurements and uncertainty estimations in code.
- Gather results from multiple executions of the test, and recompute uncertainty estimations.
 - Given two averages, μ_1 and μ_2 (and a combined average μ), of n_1 and n_2 data points, with sample standard deviations σ_1 and σ_2 , the combined sample standard deviation of both datasets is:

$$\sigma = \sqrt{\frac{n_1^2\sigma_1^2 + n_2^2\sigma_2^2 - n_2\sigma_1^2 - n_1\sigma_2^2 - n_1\sigma_1^2 - n_1\sigma_2^2 + n_1n_2(\mu_1 - \mu_2)^2}{(n_1 + n_2 - 1)(n_1 + n_2)}}$$

Confidence Intervals

Confidence Interval: a way to describe the amount of uncertainty associated with a sample of a population.

- Constructed from three pieces of information:
 - Confidence level (r) - e.g. 90%, 95%, 99%.
 - Statistical data, including sample size (n).
 - Uncertainty for the data (σ).

$$CI = \frac{z\sigma}{\sqrt{n}}$$

- z is the critical value. For large sample sizes, you can look this up in a table. For small sample sizes, use the Student's t inverse cumulative distribution function:

$$z = T_{inv}(1 - r, n - 1)$$

Confidence Intervals

One of the useful things you can do with confidence intervals is determine the correct sample size, based on an initial “pilot” set of samples.

- Given a margin of error e_m , a critical value z , an uncertainty σ , and a mean μ :

$$n = \left(\frac{z\sigma}{\frac{e_m}{2}\mu} \right)^2$$

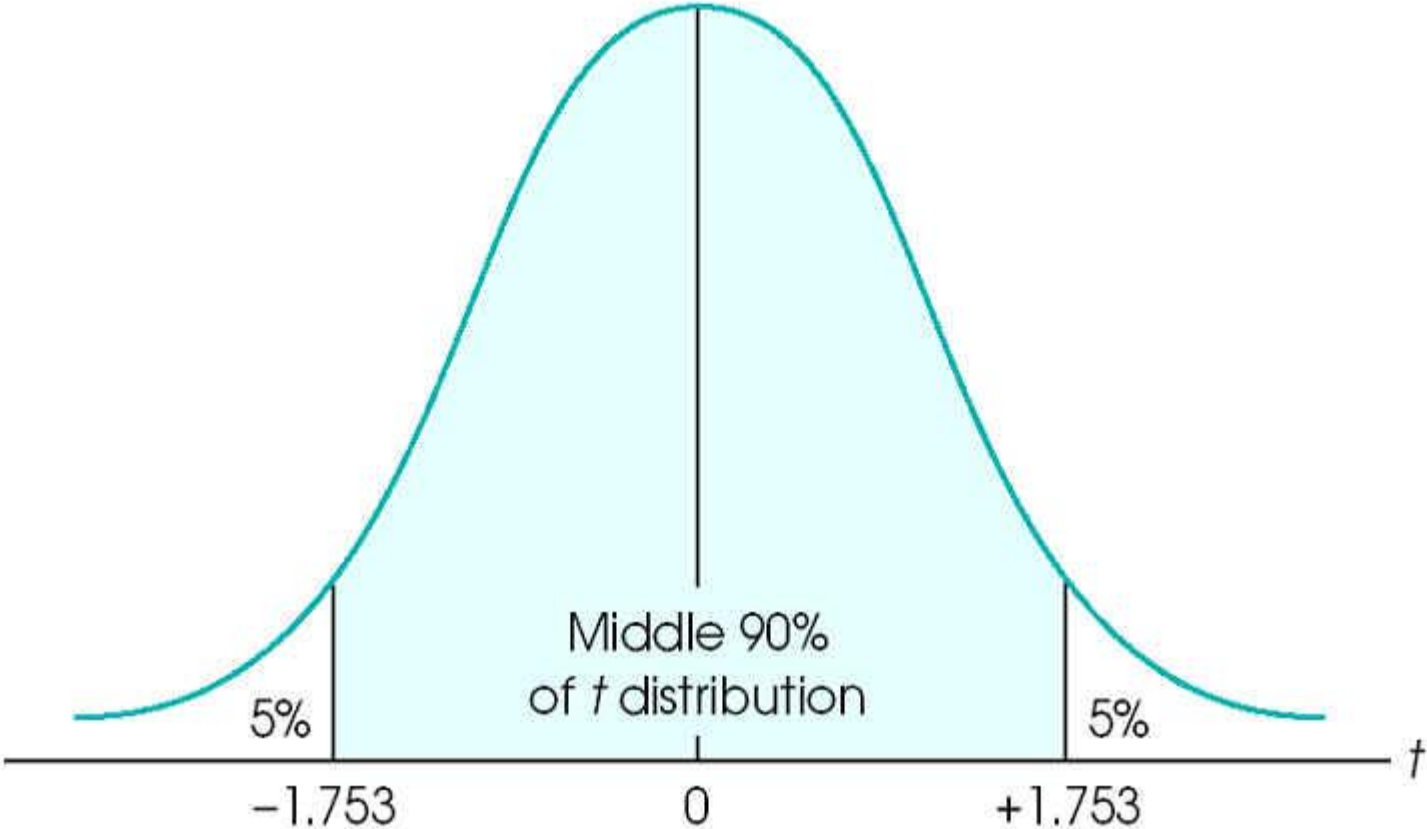
- If this calculation indicates an unreasonable large sample size is needed, the experiment may need to be redesigned.
- Typically, if your uncertainties are big relative to your data (mean and standard deviation have the same magnitude), there is too much noise to get meaningful results from your data.

Confidence Intervals

Meaning of confidence intervals

- If the true performance lies outside of the 95% confidence interval, then an event occurred which had a probability of 5% or less of happening.
- A 95% confidence interval does **not** mean that 95% of the data lies within the interval.
- A confidence interval isn't a range of plausible values for a sample mean. It can be interpreted as an estimate of plausible values for the population.

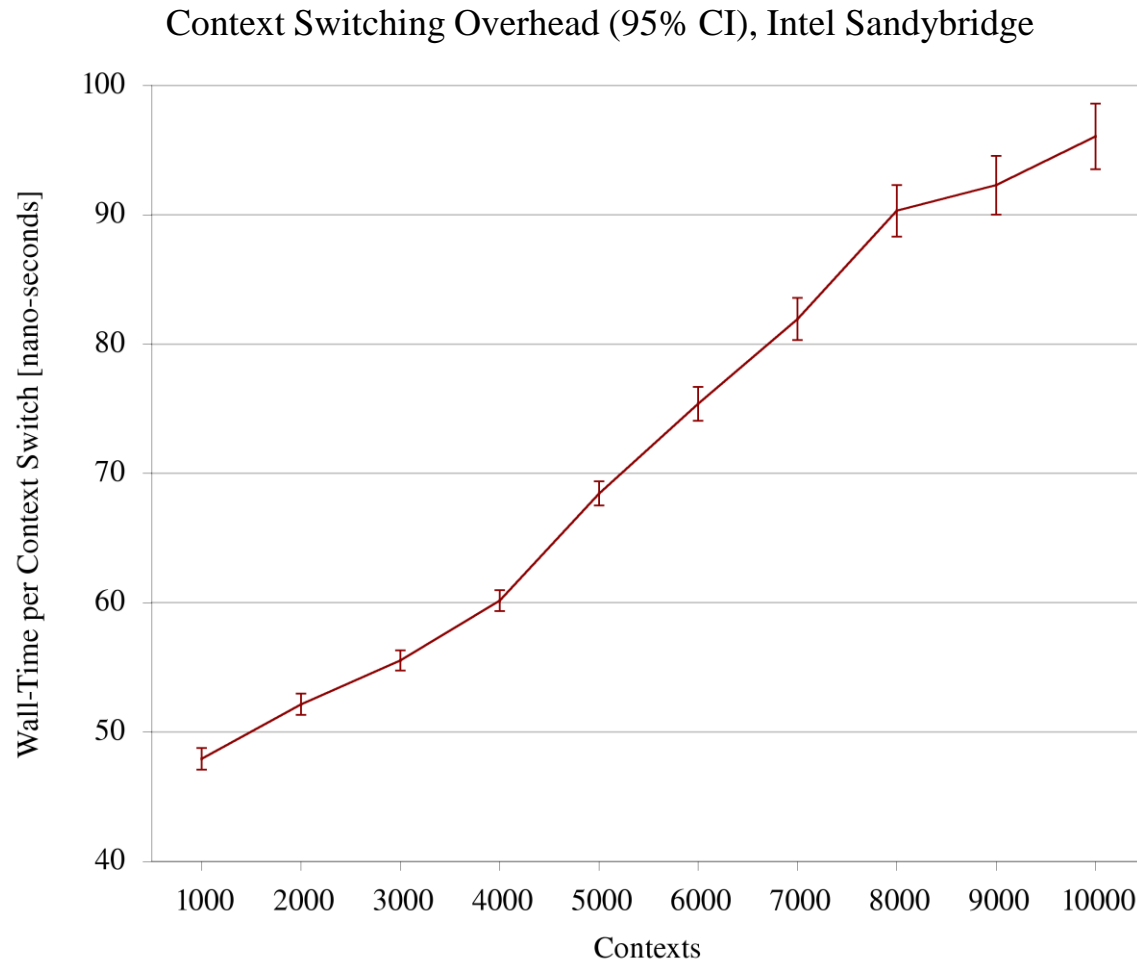
Confidence Intervals



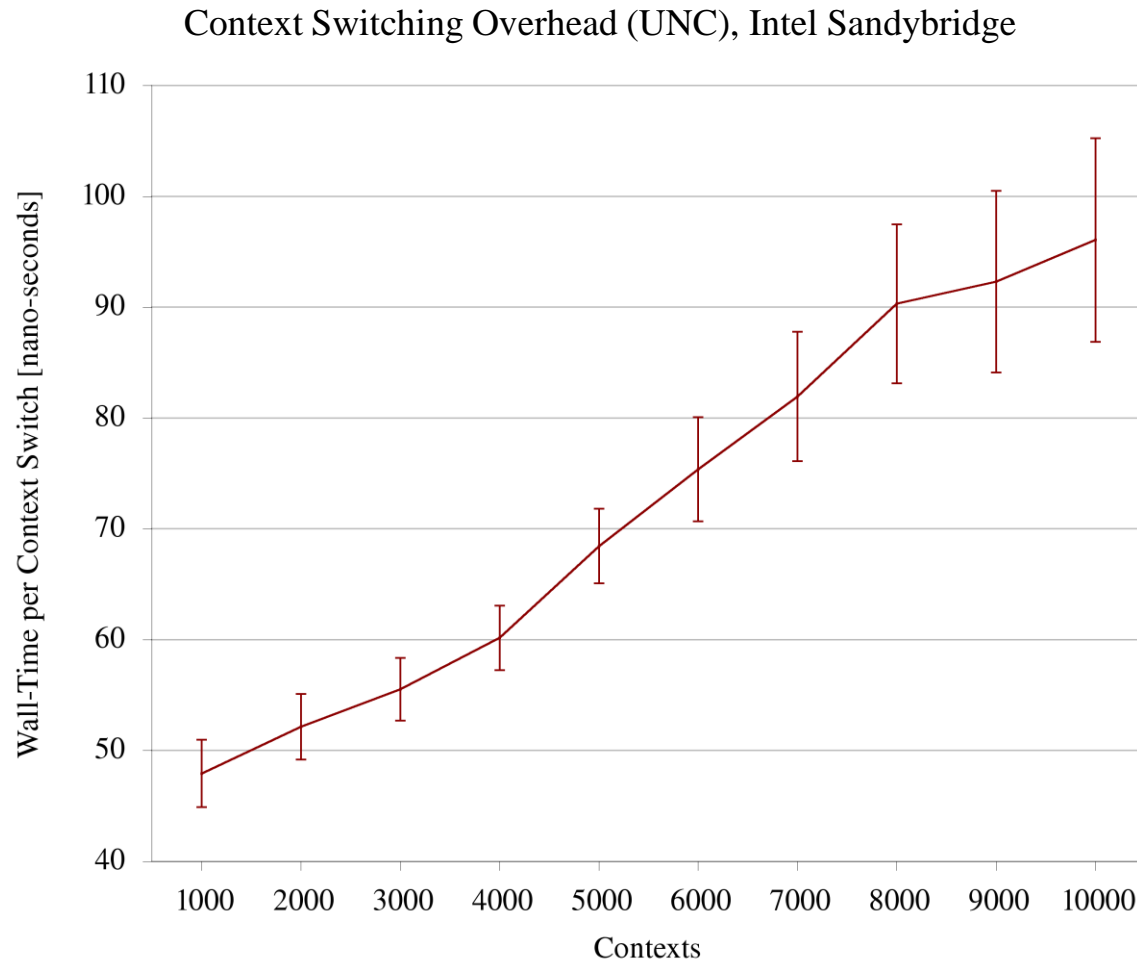
Case Study: HPX CS Overhead

Statistical Best Practices

Case Study: HPX CS Overhead



Case Study: HPX CS Overhead



Mean-Median Test

Normality test: Tests to determine if a data-set fits a normal distribution well.

- There are graphical (QQ plot), informal/back-of-the-envelope and rigorous normality tests.

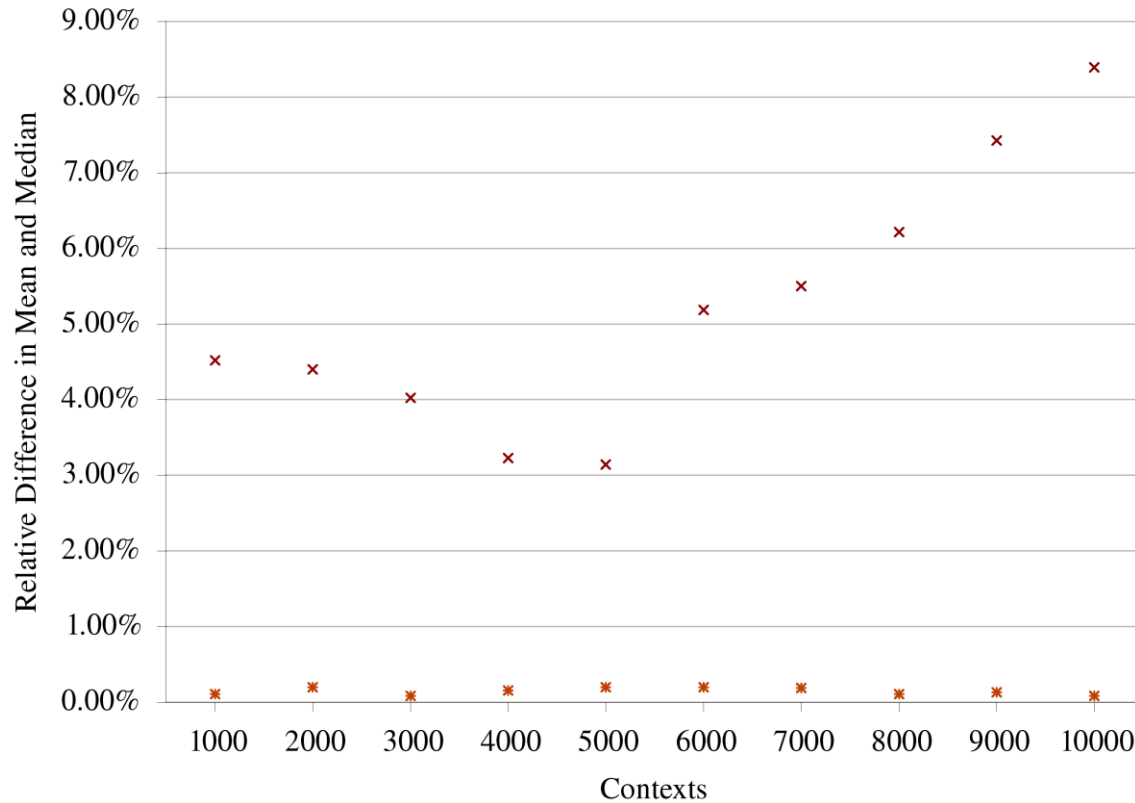
The mean (μ), median (m) and mode of normally distributed data should be the same, so...

$$\frac{|\mu - m|}{\max(\mu, m)}$$

- This will give you the relative difference between the mean and median (a percentage represented as a decimal). If this is larger than 1%, your data is probably not normally distributed.

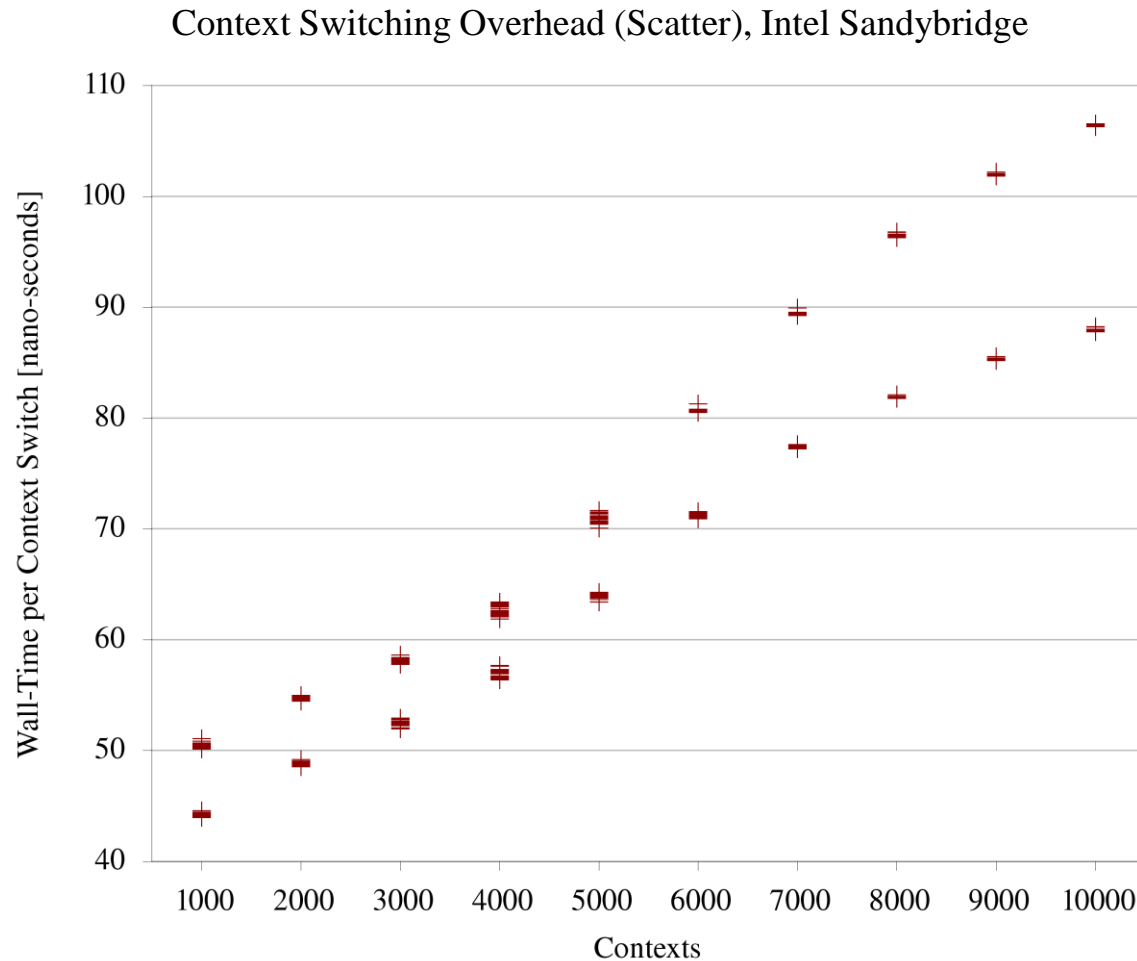
Case Study: HPX CS Overhead

Context Switching Overhead (Mean-Median Test), Intel Sandybridge



Sandybridge Xeon E5-2690, 2x8 cores, 2.90GHz ×
Ivybridge Xeon E5-2670 v2, 2x10 cores, 2.50GHz *

Case Study: HPX CS Overhead



Time-Based Benchmarking

Time-Based Benchmarking

We have access to a few different clocksources for benchmarking on modern (x86) CPUs:

- System-wide high-resolution clock:
 - Monotonic, **frequency-stable**, higher latency and overhead.
 - Resolution is in nanoseconds.
 - Times can be passed between threads.
 - *nix, this is accessed via `clock_gettime` reading `CLOCK_MONOTONIC`.
 - Windows, this is accessed via `QueryPerformanceCounter/Frequency`.
 - Suitable for measuring most events (microseconds and up).

Time-Based Benchmarking

We have access to three different clocksources for benchmarking on modern (x86) CPUs:

- Timestamp Counter (TSC):
 - Monotonic, lower latency and overhead.
 - Resolution is in CPU cycles (with caveats), tick is in base clock cycles.
 - All newer (4-5 year old) CPUs guarantee a constant TSC frequency, even if the CPU frequency changes (e.g. frequency scaling, Intel Turbo mode).
 - Constant TSC frequency == timing data is **not** representative of # of cycles executed.
 - Ticks with the base clock, which runs at 100 or 133 Mhz (depending on microarchitecture).
 - Assembly instruction(s) for reading this counter.
 - Cycle counts are thread-specific.
 - Suitable for measuring short events (cycles to minutes).

<chrono>

Standard facilities for manipulating dates and times, introduced in C++11

- Three types:
 - Duration: A span of time, defined as some number of ticks of some time unit.
 - Time Point: A duration of time that has passed since the epoch of specific clock.
 - Clocks: An object with a starting point and a tick rate, which can be queried for the current time.

<chrono> is the best way to measure durations that are microsecond magnitude or large.

Source: cppreference.com

<chrono>

Clock	Description
system_clock	Wall clock time from the system-wide realtime clock.
steady_clock	Monotonic clock that will never be adjusted.
high_resolution_clock	The clock with the shortest tick period available.

Example: `high_resolution_timer`

Time-Based Benchmarking

Example: high_resolution_timer

```
struct high_resolution_timer
{
    high_resolution_timer() : start_time_(take_time_stamp()) {}

    void restart()
    { start_time_ = take_time_stamp(); }

    double elapsed() const // Return elapsed time in seconds.
    { return double(take_time_stamp() - start_time_) * 1e-9; }

    std::uint64_t elapsed_nanoseconds() const
    { return take_time_stamp() - start_time_; }

protected:
    static std::uint64_t take_time_stamp()
    {
        return std::chrono::duration_cast<std::chrono::nanoseconds>
            (std::chrono::steady_clock::now().time_since_epoch()).count();
    }

private:
    std::uint64_t start_time_;
};
```

Non-Time-Based Benchmarking

Memory Benchmarking

Approaches to instrumenting memory allocation:

- What do we want to look at?
 - Objects (allocated/deallocated)
 - Memory (total, per object size, per object type)
- External tools:
 - googleperftools/TCMalloc (MALLOCSTATS)
 - MemTrack
- Overload operator new/delete
 - Writing a member operator new/delete is a great technique for tracking memory performance for a specific object.
 - I suggest a static member variable to store the performance data; if you need thread safety, use thread-local storage and accumulate afterwards.

Example: Instrumenting operator new

Non-Time-Based Benchmarking

Example: Instrumenting operator new

```
struct A {  
    static std::size_t allocated;  
  
    static void* operator new(std::size_t sz)  
    {  
        allocated += sz/sizeof(A);  
        return ::operator new(sz);  
    }  
    static void* operator new[](std::size_t sz)  
    {  
        allocated += sz/sizeof(A);  
        return ::operator new(sz);  
    }  
};  
  
std::size_t A::allocated = 0;
```


Counting Copies/Moves

When we started transition the HPX codebase to support move semantics a few years ago, we wrote some tests to make sure we got it right.

- We passed mock objects that count copies/moves through our framework and looked at the results.
- Once we were confident our interfaces were doing things right (minimizing the number of copies, etc), we wrote unit tests to verify the move/copy counts wouldn't change.
- Especially important for us – HPX is an asynchronous programming framework, so there are places where we duplicate data to facilitate asynchrony.
 - We wanted to ensure we only copied `async()` arguments once.

Hardware Performance Counters

X86 processors have a diverse set of hardware performance counters.

- Pros:
 - Low overhead.
 - Very diverse and descriptive information.
- Cons:
 - Microarchitecture specific.
 - Some counters are estimations, or suffer from inaccuracies (overcounting, etc).
 - You need very specialized knowledge to use these for performance analysis.
 - Fortunately, there's an awesome tool which has this knowledge baked into it.

Hardware Performance Counters

Low-level frameworks for accessing hardware counters from within your code:

- Linux: PAPI framework
- Windows: Performance Counter framework
- Mac: kpc.h

There are some external sampling-based profiling tools that provide access to this information.

- Ex: Intel VTune Amplifier.

Performance Analysis Tools

Intel VTune Amplifier

Sampling-based profiling tool: runs your application, and collects “snapshots” of performance metrics while your program is running.

- Works on Intel processors, Windows/Linux/Mac OS X/Android, not tied to any particular compiler.
- Requires no code changes to use.
- Multiple data sources: timers, hardware performance counters and operating system metrics.
- Performance data can be viewed per function or at assembly/source code granularity.
- Analyzes everything: kernel calls, sub-processes, threads, etc.

Intel VTune Amplifier

Sampling-based profiling tool: runs your application, and collects “snapshots” of performance metrics while your program is running.

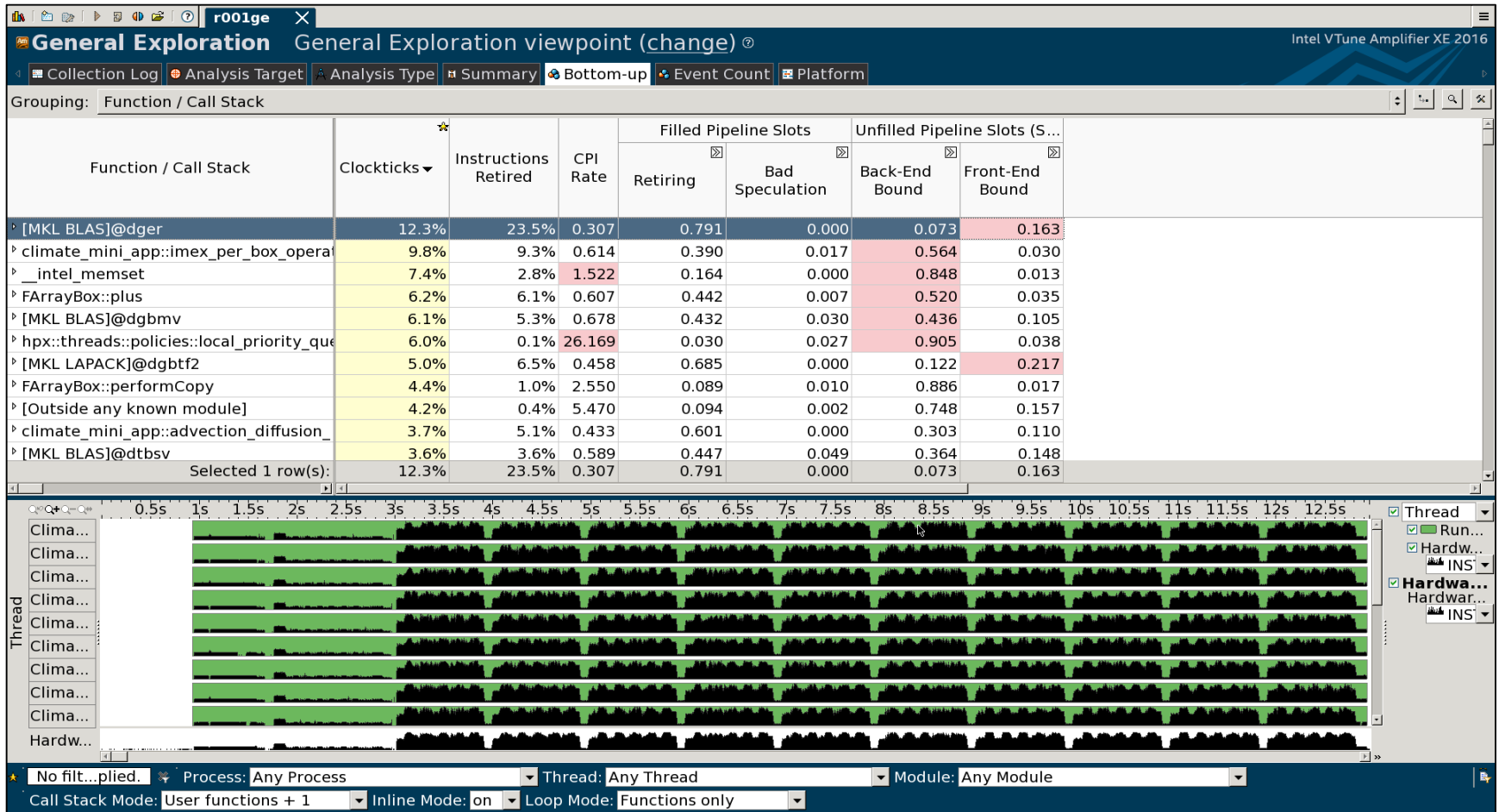
- Provides built-in analysis passes which derive useful, higher-level performance metrics from micro-architecture specific raw hardware counters.
- Also supports user-defined analysis passes.
- Support for instrumenting parallel and distributed code.
 - Built-in support for OS-threading frameworks.
 - Built-in support for OpenMP, MPI and Intel TBB.
 - Provides an instrumentation API which parallel programming frameworks can use to inform the profiler about their threading and concurrency data structures.

Intel VTune Amplifier

Sampling-based profiling tool: runs your application, and collects “snapshots” of performance metrics while your program is running.

- Powerful GUI.
 - Standalone Windows/Linux/Mac GUI as well as integration with Visual Studio and Eclipse.
 - Data can be collected remotely via the command line interface and then fed into the GUI.
 - Great interface for filtering data (e.g. focusing in on just one section of the program’s execution).
 - Built-in analysis passes contain a lot of information about how to interpret results.

Intel VTune Amplifier



Intel VTune Amplifier

General Exploration General Exploration viewpoint (change) Intel VTune Amplifier XE 2016

Collection Log Analysis Target Analysis Type Summary Bottom-up Event Count Platform

Grouping: Function / Call Stack

Function / Call Stack	Filled Pipeline Slots		Unfilled Pipeline Slots (Stalls)														Front-End Bound				
	Retiring	Bad Speculation	Memory Latency							Memory Replace...			Memory Reissues			Divi...		Fla... Mer...	Slow LEA...		
			LLC...	LLC ..	LLC ..	DTL...	Con..	Dat...	L1D ..	L2 ...	LLC...	Loa..	Spli...	Spli...	4K ...						
▸ [MKL BLAS]@dger	0.815	0.000	0.038	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.023	0.002	0.114	0.000	0.000	0.035	0.169
A significant proportion of cycles is spent dealing with false 4k aliasing between loads and stores. Use the source/assembly view to identify the aliasing loads and stores, and then adjust your data layout so that the loads and stores no longer alias.																					
Threshold: (((5 * LD_BLOCKS_PARTIAL.ADDRESS_ALIAS) / CPU_CLK_UNHALTED.THREAD) > 0.1) * (CPU_CLK_UNHALTED.THREAD / > 0.05))																					
▸ [MKL LAPACK]@dgbtf2	0.686	0.026	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.009	0.188	0.000	0.068	0.000	0.226
▸ climate_mini_app::imex_per_box_operat	0.384	0.025	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.003	0.000	0.000	0.000	0.000	0.071	0.000	0.081	0.000	0.061	0.000	0.157
▸ FArrayBox::performCopy	0.056	0.015	0.985	0.000	0.018	0.027	0.000	0.009	0.005	0.003	0.012	0.000	0.000	0.000	0.005	0.000	0.000	0.000	0.000	0.000	0.005
▸ [MKL BLAS]@dtbsv	0.468	0.000	0.120	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.061	0.401	0.000	0.022	0.000	0.150
▸ climate_mini_app::advection_diffusion_	0.509	0.069	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.003	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.132
▸ climate_mini_app::imex_per_box_operat	0.561	0.000	1.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.000	0.000	0.000	0.000	0.020	0.000	0.000	0.130	0.000	0.124

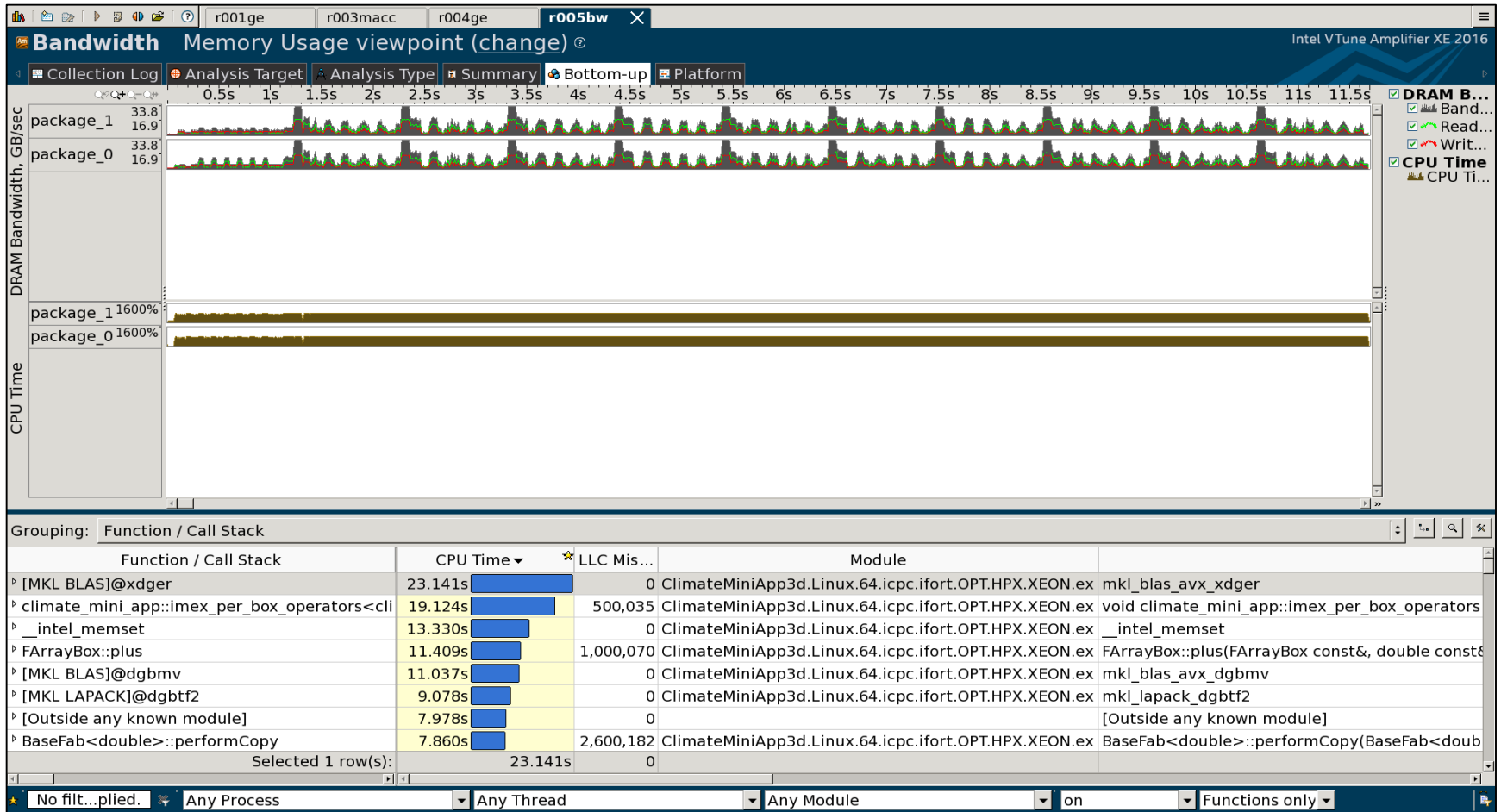
Highlighted 0 row(s)

Thread: Climax... Hardware...

Filter: ... shown Process: Any Process Thread: Any Thread Module: Any Module

Call Stack Mode: User functions + 1 Inline Mode: on Loop Mode: Functions only

Intel VTune Amplifier



Intel Vectorization Adviser

New tool in Intel Parallel Studio XE 2016: Vectorization Adviser.

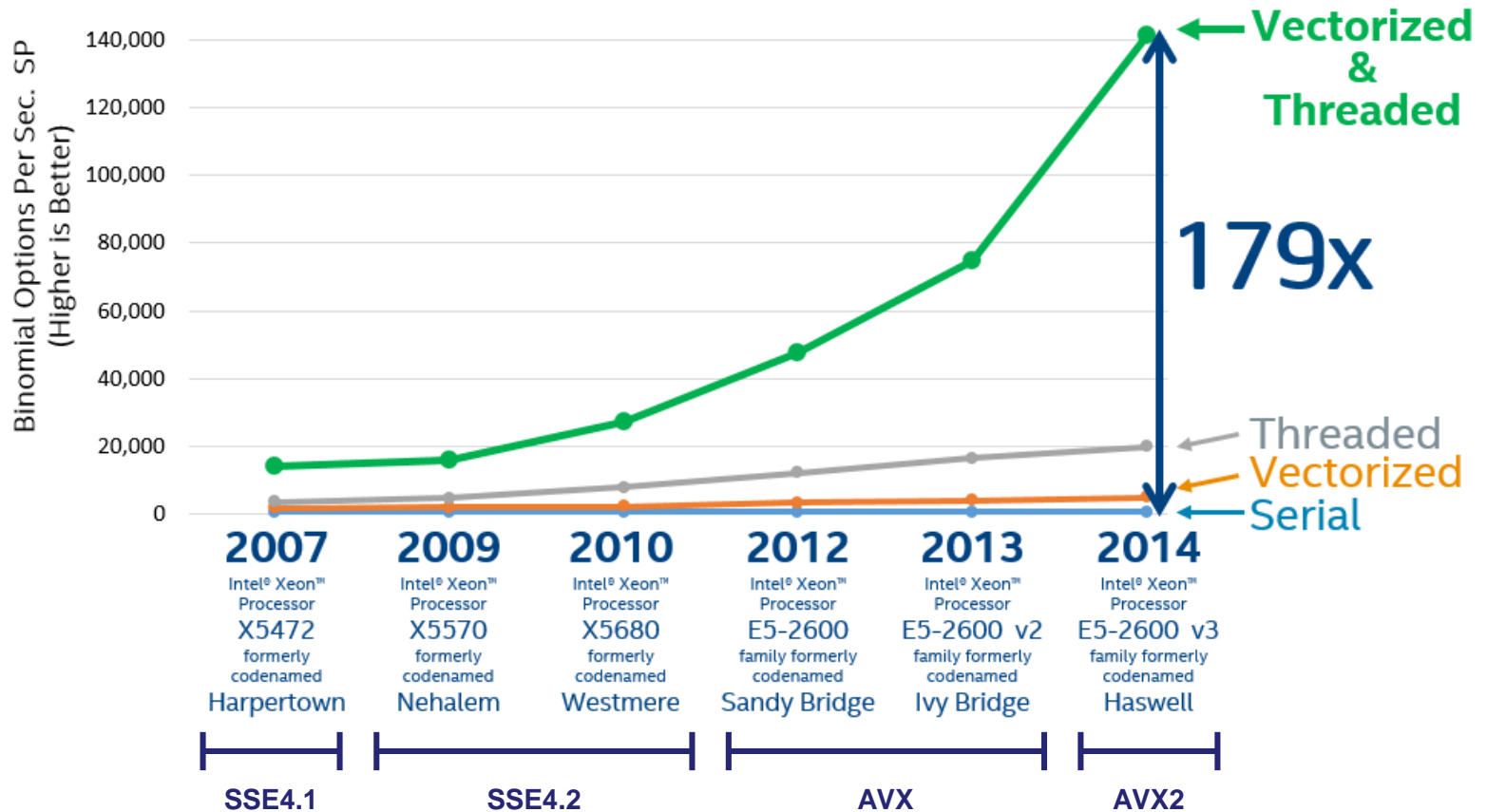
- Integrates the Intel compiler's vectorization reports into the GUI performance profiling framework.

The screenshot shows the Intel Advisor XE 2016 interface. At the top, there's a title bar with the text "Where should I add vectorization and/or threading parallelism?" and "Intel Advisor XE 2016". Below the title bar are several tabs: "Summary", "Survey Report", "Refinement Reports", "Annotation Report", and "Suitability Report". The "Survey Report" tab is active. Below the tabs, there's a filter section with "Elapsed time: 54.44s", "Vectorized", "Not Vectorized", and a search icon. The main area is a table with columns: "Function Call Sites and Loops", "Vector Issues", "Self Time", "Total Time", "Trip Counts", "Loop Type", "Why No Vectorization?", and "Vectorized Loops". The "Vectorized Loops" column is further divided into "Vecto..." and "Efficiency".

Function Call Sites and Loops	Vector Issues	Self Time	Total Time	Trip Counts	Loop Type	Why No Vectorization?	Vectorized Loops	
							Vecto...	Efficiency
[loop at stl_algo.h:4740 in std::tr...]		0.170s	0.170s		Scalar	non-vectorizable loop ins ...		
[loop at loopstl.cpp:2449 in s234_]	2 Ineffective peeled/rem ..	0.170s	0.170s	12; 4	Collapse	Collapse	AVX	~100%
[loop at loopstl.cpp:2449 in s ...]		0.150s	0.150s	12	Vectorized (B		AVX	
[loop at loopstl.cpp:2449 in s ...]		0.020s	0.020s	4	Remainder			
[loop at loopstl.cpp:7900 in vas_]		0.170s	0.170s	500	Scalar	vectorization possible but ...		
[loop at loopstl.cpp:3509 in s2 ...]	1 High vector register ...	0.160s	0.160s	12	Expand	Expand	AVX	~69%
[loop at loopstl.cpp:3891 in s279_]	2 Ineffective peeled/rem ..	0.150s	0.150s	125; 4	Expand	Expand	AVX	~96%
[loop at loopstl.cpp:6249 in s414_]		0.150s	0.150s	12	Expand	Expand	AVX	~100%
[loop at stl_numeric.h:247 in std ...]	1 Assumed dependency ...	0.150s	0.150s	49	Scalar	vector dependence preve ...		

Source: <https://software.intel.com/en-us/intel-advisor-xe/>

Intel Vectorization Adviser



Source: <https://software.intel.com/en-us/intel-advisor-xe/>

Write Performance Tests

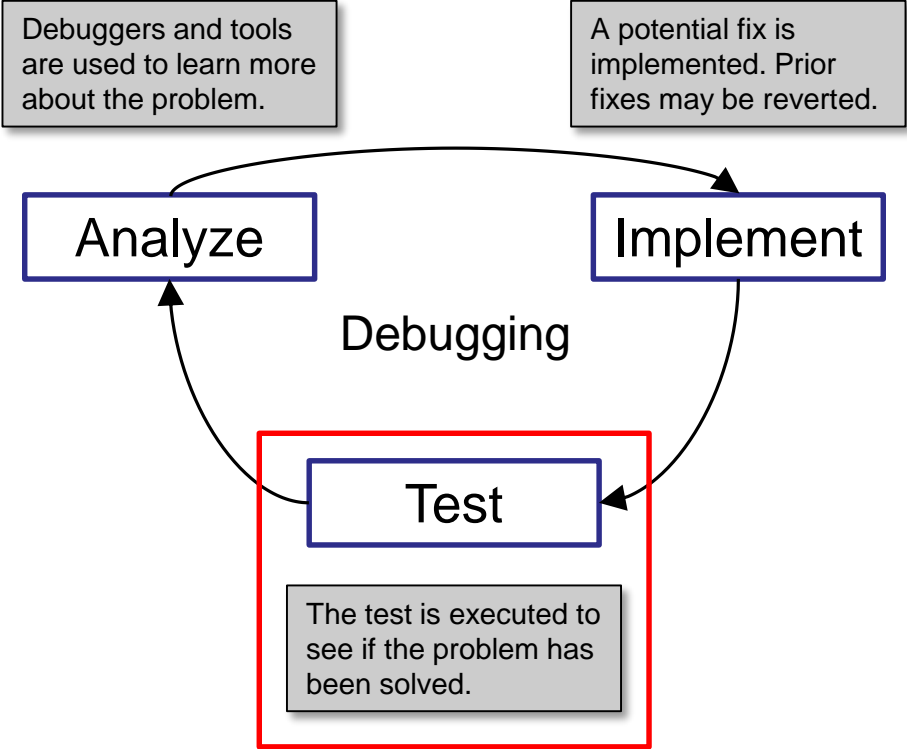
Write Performance Tests

Idea: Let's write unit and regression tests for performance, just like we do for correctness.

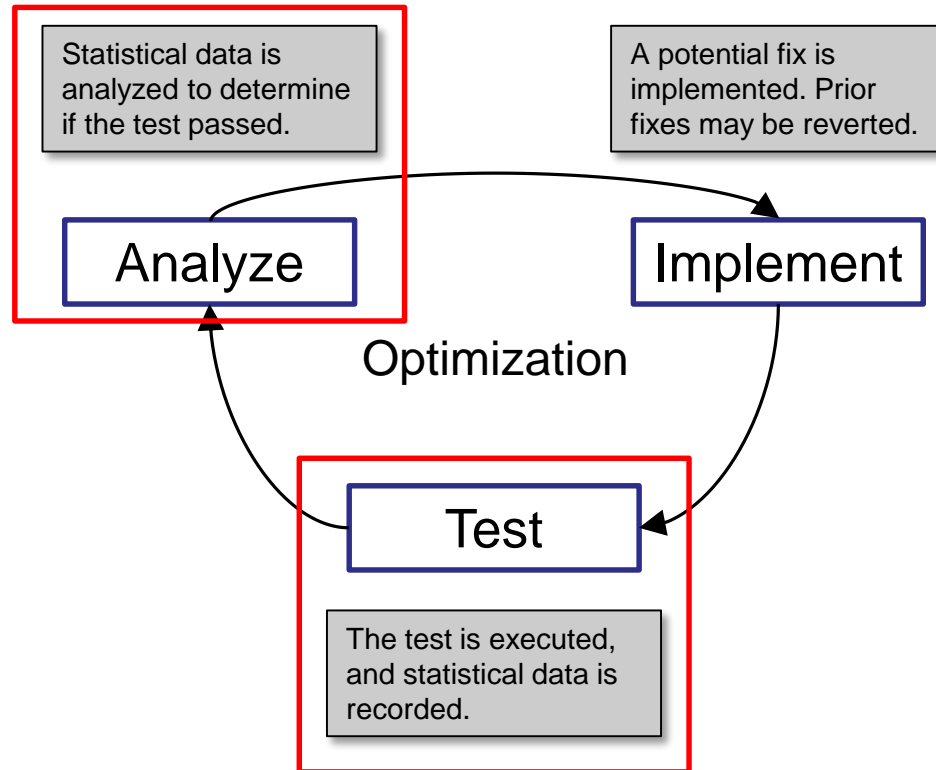
Challenges:

- Implementing automated performance testing that follow the kind of best practices we've been talking about requires a lot of machinery.
- If your performance tests are stateful (rely on the results of previous automated tests), you need even more machinery.
- You need more than just the machinery to run the tests - you also need automated analysis to determine whether the test has failed.

Write Performance Tests



Write Performance Tests



Stateful Performance Tests

Stateful performance tests: performance benchmarks which yield results that cannot be interpreted without contextual information.

- Output: absolute values.
- Most of your existing benchmarks are already stateful.
- To automate these tests, the current performance (e.g. trunk) needs to be compared against some prior results. There's two options for doing this:
 - Automated build system stores prior results for comparison. Requires more machinery, but allows you to track performance over time.
 - Automated build system checks out and builds an older version of the code to compare against.

Stateful Performance Tests



Stateless Performance Tests

Stateless performance tests: performance benchmarks which test for a performance “failure”, and can provide a “**yes/no**” answer without external data.

- Output: relative values.

The idea is to compare different implementation options which you believe to have a performance impact.

- Ex: Lockfree queue vs lock-based queue.
- Ex: Recomputing data locally vs overhead for sharing.
- Ex: Algorithmic complexity testing.

Stateless Performance Tests

Whenever you face a design trade-off
with performance implications,
write a stateless test!

Summary

- Take a scientific approach to performance benchmarking.
 - Hypothesize, design a test, run the test, analyze, draw conclusions.
- Manage and test your assumptions about your tests.
- Collect a statistically significant quantity of data.
- Measure and propagate error.
- Develop unit and regression tests for performance.



U.S. DEPARTMENT OF
ENERGY



**UNIVERSITY OF
CALIFORNIA**