



SC21

St. Louis, MO | science & beyond.

Sampling-based Performance Analysis with HPCToolkit

Measurement and Analysis of Unmodified, Optimized Applications

John Mellor-Crummey

Department of Computer Science
Rice University

15 November 2021



Outline

- **Sampling-based call path profiling**
- **Using HPCToolkit on a single node**
- **Demo data collection and profile analysis**
- **Demo trace analysis**
- **Pointers to additional information about measurement**

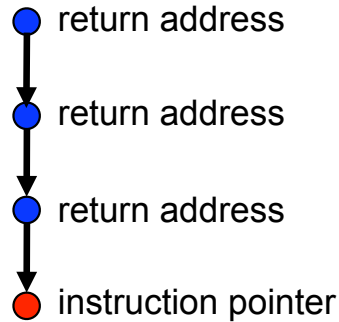
Sampling-based Performance Measurement

- **Periodically interrupt each thread in an application**
 - interrupts are triggered by a “sample source” as a metric reaches some pre-determined threshold
 - example sample sources
 - timer
 - a thousandth of a second has passed
 - hardware counters
 - five million instructions have completed
 - a million cache misses have occurred
- **Why sampling?**
 - controllable overhead
 - avoids blind spots and minimizes systematic error

Attribute Metrics to Call Paths

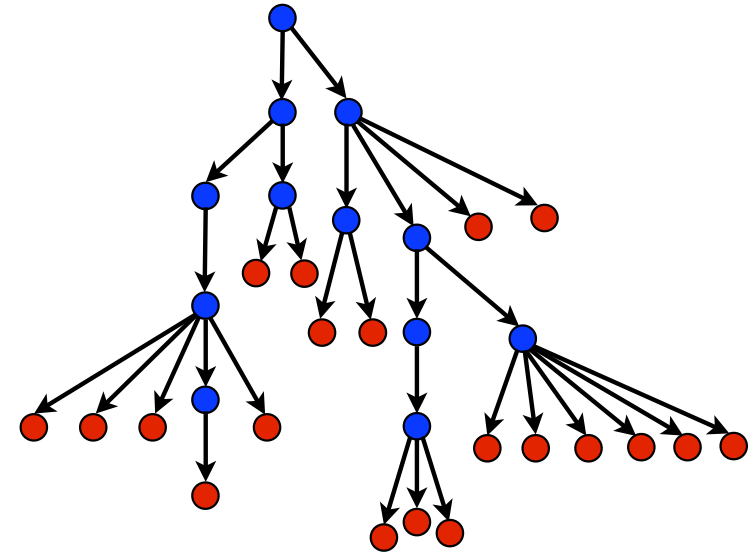
- When a thread's timer or HW counter reaches some predetermined threshold T
 - interrupt a thread
 - unwind its call stack

Call path sample



- charge T to the thread's current calling context

Calling context tree



Overhead proportional to sampling frequency, not call frequency

HPCToolkit Quickstart

```
% hpcrun myapp  
# profile CPU TIME of application and deposit results in hpctoolkit-myapp-measurements  
  
% hpcstruct hpctoolkit-myapp-measurements  
# analyze application binary and all dynamically loaded libraries involved in the execution  
  
% hpcprof hpctoolkit-myapp-measurements  
# analyze measurement data and correlate it to source using program structure from hpcstruct  
  
% hpcviewer hpctoolkit-myapp-database
```

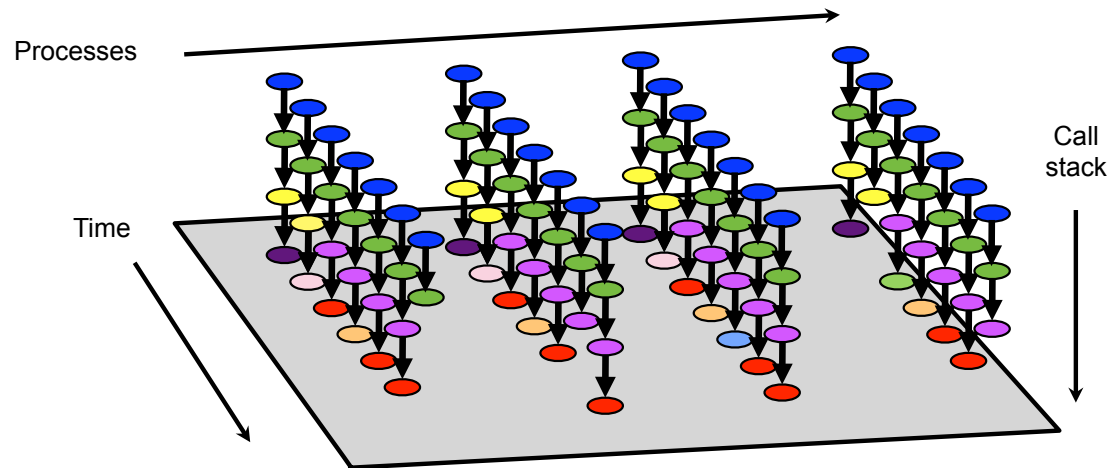
Video: Using HPCToolkit to Measure an OpenMP Program

The terminal window shows the execution of an OpenMP program named 'lulesh2.0'. The code includes OpenMP pragmas for parallelization and a loop for calculating hourglass modes. The hpcviewer tool is used to analyze the performance of the program, displaying a table of execution metrics.

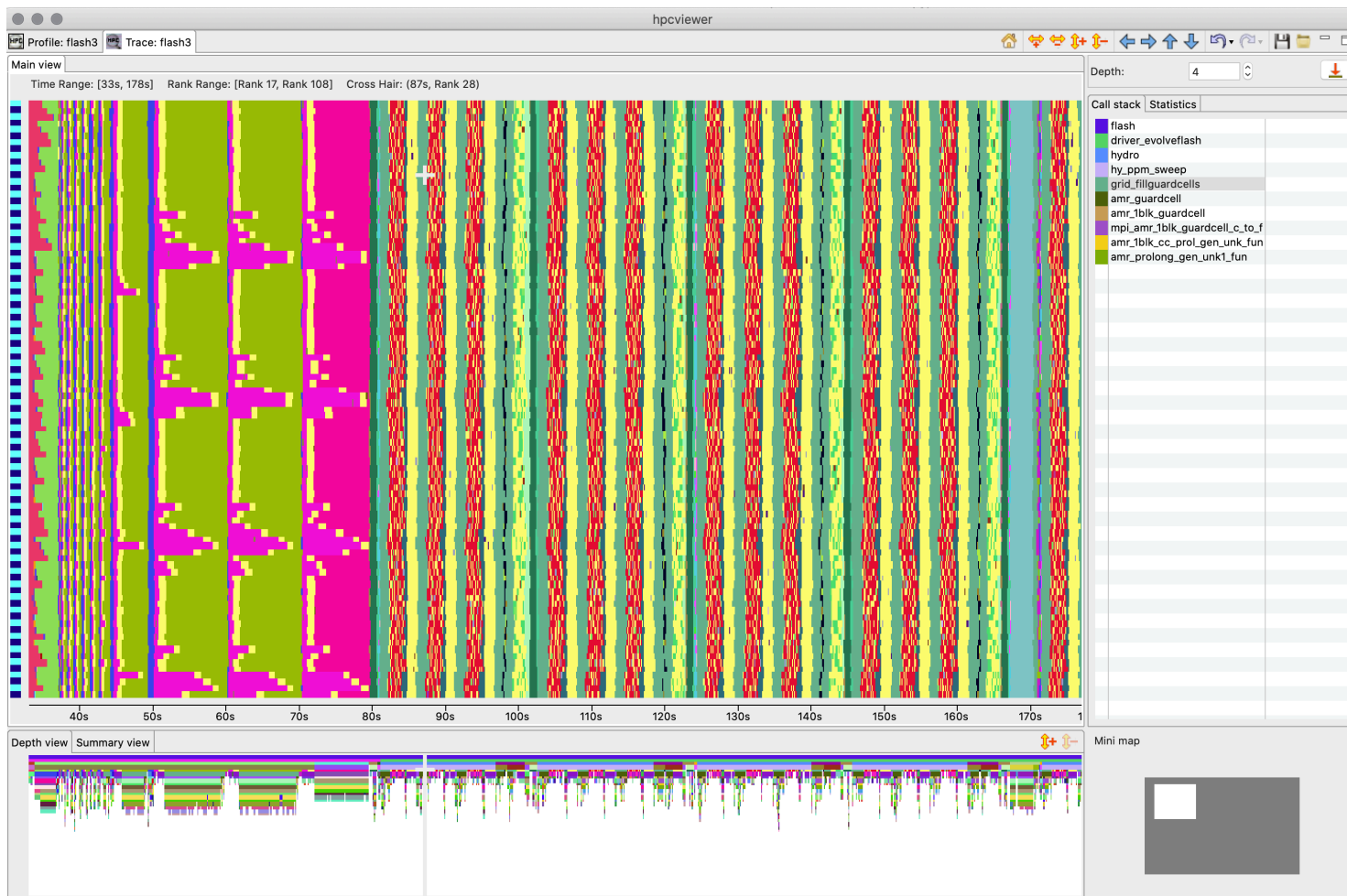
Scope	CPUTIME (sec):Sum (I)	CPUTIME (sec):Sum (E)
Experiment Aggregate Metrics	1.11e+02 100.0%	1.11e+02 100.0%
[I] do spin	1.62e+01 14.7%	1.62e+01 14.7%
CalcFBHourglassForceForElems(Domain&, double*, double*, double*, double*, double*, double*, double*, double*, double, ...)	2.16e+01 19.6%	1.30e+01 11.7%
gomp_thread_start [libgomp.so.1.0.0]	1.90e+01 17.2%	1.14e+01 10.4%
GOMP_parallel [libgomp.so.1.0.0]	2.66e+00 2.4%	1.53e+00 1.4%
[I] CalcFBHourglassForceForElems	2.66e+00 2.4%	1.53e+00 1.4%
[I] CalcHourglassControlForElems	2.66e+00 2.4%	1.53e+00 1.4%
[I] CalcVolumeForceForElems	2.66e+00 2.4%	1.53e+00 1.4%
[I] CalcForceForNodes	2.66e+00 2.4%	1.53e+00 1.4%
LagrangeNodal	2.66e+00 2.4%	1.53e+00 1.4%
LagrangeLeapFrog	2.66e+00 2.4%	1.53e+00 1.4%
main	2.66e+00 2.4%	1.53e+00 1.4%
<program root>	2.66e+00 2.4%	1.53e+00 1.4%
EvalEOSForElems(Domain&, double*, int, int*, int) [clone ._omp_fn.0]	8.37e+00 7.6%	8.04e+00 7.3%
CalcHourglassControlForElems(Domain&, double*, double) [clone ._omp_fn.0]	1.32e+01 12.0%	6.63e+00 6.0%
CalcElemFBHourglassForce	5.66e+00 5.1%	5.66e+00 5.1%
VolUder	5.57e+00 5.0%	5.57e+00 5.0%
CalcElemShapeFunctionDerivatives(double const*, double const*, double const*, double (*) [8], double*)	4.46e+00 4.0%	4.46e+00 4.0%
CalcMonotonicGradientsForElems(Domain&) [clone ._omp_fn.0]	4.70e+00 4.2%	4.30e+00 3.9%
SumElemFaceNormal	3.78e+00 3.4%	3.78e+00 3.4%
CalcMonotonicRegionForElems(Domain&, int, double) [clone ._omp_fn.0]	3.43e+00 3.1%	3.35e+00 3.0%
CalcPressureForElems(double*, double*, double*, double*, double*, double*, double, double, double, int, int*)...	3.28e+00 3.0%	2.95e+00 2.7%

Understanding Temporal Behavior

- **Profiling compresses out the temporal dimension**
 - Temporal patterns, e.g. serial sections and dynamic load imbalance are invisible in profiles
- **What can we do? Trace call path samples**
 - N times per second, take a call path sample of each thread
 - Organize the samples for each thread along a time line
 - View how the execution evolves left to right
 - What do we view? assign each procedure a color; view a depth slice of an execution



Video: Using HPCToolkit to Analyze the Trace of an MPI Program



See Slide Deck for Additional Details about Measurement

- Measuring applications when using a job launcher
- Specifying sample sources
 - timers
 - hardware counters
- Controlling measurement frequency
 - automatic
 - frequency-based sampling
 - period-based sampling

Measuring Performance with hpcrun

- Profile a dynamic binary (sequential or multithreaded)
 - `hpcrun [measurement options] myapp`
- Use hpcrun with example job launcher commands
 - `jsrun -n 32 -g 1 -a 1 hpcrun [measurement options] myapp`
 - `srun -n 1 -G 1 hpcrun [measurement options] myapp`
 - `aprun -n 16 -N 8 -d 8 hpcrun [measurement options] myapp`
- Specifying CPU events to measure
 - `hpcrun -e <event1>[@<howoften1>] -e <event2>[@<howoften2>] myapp`

Note: To profile statically-linked applications, you must link your application with HPCToolkit's measurement subsystem using `hplink`. See the HPCToolkit manual for details.

hpcrun - Tracing

- Specify tracing simply by adding “-t” as an argument to hpcrun
- Requirements
 - must be measuring execution with a time-based metric
 - Linux timer
 - “cycles” measured with `perf_event`

Sample Sources: Linux Timers

- **CPUTIME** (DEFAULT if no sample source is specified)

- -e CPUTIME@<period>: interrupt each thread every <period> microseconds it executes
- does not include time blocked in the kernel
 - disadvantage: misses time a thread is blocked
 - advantage: a blocked thread is never woken to take a sample

Best for analysis of
profile data

- **REALTIME**

- -e REALTIME@<period>: interrupt each thread every <period> microseconds
- includes time blocked in the kernel
 - advantage: shows where a thread spends its time, even when blocked
 - disadvantages
 - activates a blocked thread to take a sample
 - a blocked thread appears active even when blocked

May produce more
intuitive traces

Sample Sources: Hardware Counters

- Each core in a modern processor has a performance monitoring unit with counters for HW events
 - each thread has a small number of HW counters
- Linux kernel: perf_event subsystem for performance monitoring
 - access and manipulate
 - hardware counters: cycles, instructions, ...
 - software counters: context switches, page faults, ...
 - available in Linux kernels 2.6.31+

A useful explanation about events available through perf
<https://sites.google.com/site/lbathen/research/perf>

Sample Sources: perf_event Hardware Event Counters

- PERF_COUNT_HW_CPU_CYCLES
- PERF_COUNT_HW_INSTRUCTIONS
- PERF_COUNT_HW_CACHE_REFERENCES
- PERF_COUNT_HW_CACHE_MISSES
- PERF_COUNT_HW_BRANCH_INSTRUCTIONS
- PERF_COUNT_HW_BRANCH_MISSES
- PERF_COUNT_HW_BUS_CYCLES
- PERF_COUNT_HW_STALLED_CYCLES_FRONTEND
- PERF_COUNT_HW_STALLED_CYCLES_BACKEND
- PERF_COUNT_HW_REF_CPU_CYCLES

Sample Sources: `perf_event` Hardware Cache Events

- **Hardware cache**

- `PERF_COUNT_HW_CACHE_L1D`
- `PERF_COUNT_HW_CACHE_L1I`
- `PERF_COUNT_HW_CACHE_LL`
- `PERF_COUNT_HW_CACHE_DTLB`
- `PERF_COUNT_HW_CACHE_ITLB`
- `PERF_COUNT_HW_CACHE_BPU`

- **Operations**

- `PERF_COUNT_HW_CACHE_OP_READ`
- `PERF_COUNT_HW_CACHE_OP_WRITE`
- `PERF_COUNT_HW_CACHE_OP_PREFETCH`

- **Results**

- `PERF_COUNT_HW_CACHE_RESULT_ACCESS`
- `PERF_COUNT_HW_CACHE_RESULT_MISS`

Sample Sources: `perf_event` Software Events

- `PERF_COUNT_SW_CPU_CLOCK`
- `PERF_COUNT_SW_TASK_CLOCK`
- `PERF_COUNT_SW_PAGE_FAULTS`
- `PERF_COUNT_SW_CONTEXT_SWITCHES`
- `PERF_COUNT_SW_CPU_MIGRATIONS`
- `PERF_COUNT_SW_PAGE_FAULTS_MIN`
- `PERF_COUNT_SW_PAGE_FAULTS_MAJ`
- `PERF_COUNT_SW_ALIGNMENT_FAULTS`
- `PERF_COUNT_SW_EMULATION_FAULTS`

useful when monitoring data-intensive codes

Sample Sources: Measuring Other HW Events with `perf_event`

- See the full list of available events with
 - `hpcrun -L`
- Perf events are grouped by categories indicated by a prefix
 - `ix86arch::` // Intel architecture
 - `perf::` // `perf_event` builtin
 - `bdw_ep::` // Broadwell EP specific
 - ...
- For convenience
 - you may omit the category prefix, e.g. “`perf::`”
 - you may specify `perf_event` counter names using lower case

Controlling perf_event Sampling Frequency

- **Automatic**

Recommended

- HPCToolkit samples `perf_event` counters `min(300x/second, maximum Linux allows)`
 - may be higher than necessary for long executions
 - reducing the frequency will reduce measurement overhead

- **Specify frequency**

- use the `@f<freq>` suffix for an event to specify frequency
 - `hpcrun -e cycles@f100 -e instructions@f200 ...`
- specify a different default frequency using the `-c` option
 - example: sample both cycles and instructions 200x per second
 - `hpcrun -c f200 -e cycles -e instructions ...`

- **Specify period**

- use the `@<period>` suffix for an event to specify a period
 - `hpcrun -e cycles@1000000 -e instructions@5000000 ...`

Sample Sources: Multiplexing Events

- A single execution can measure more HW events than the number of counters available per thread
- If you specify more events than counters available
 - `perf_event` will automatically multiplex them
- How multiplexing works with Linux `perf_event` subsystem
 - at any time, the number of events being collected will not exceed the number of HW counters per thread
 - kernel will partition events into sets that can be monitored simultaneously using counter resources
 - monitors one set of events for a while then switches to another
 - kernel uses schedules event sets round-robin
 - multiplexing is convenient but there is some loss of accuracy
 - advice: multiplexing is fine for casual execution analysis