

HPCTOOLKIT User's Manual

Version 2020.07.15

John Mellor-Crummey

Laksono Adhianto, Mike Fagan, Mark Krentel, Nathan Tallent, Keren Zhou

Rice University

July 15, 2020

Contents

1	Introduction	1
2	HPCToolkit Overview	6
2.1	Asynchronous Sampling	6
2.2	Call Path Profiling	7
2.3	Recovering Static Program Structure	8
2.4	Presenting Performance Measurements	9
3	Quick Start	10
3.1	Guided Tour	10
3.1.1	Compiling an Application	11
3.1.2	Measuring Application Performance	11
3.1.3	Recovering Program Structure	13
3.1.4	Analyzing Measurements & Attributing Them to Source Code	13
3.1.5	Presenting Performance Measurements for Interactive Analysis	14
3.1.6	Effective Performance Analysis Techniques	14
3.2	Additional Guidance	14
4	Effective Strategies for Analyzing Program Performance	16
4.1	Monitoring High-Latency Penalty Events	16
4.2	Computing Derived Metrics	16
4.3	Pinpointing and Quantifying Inefficiencies	19
4.4	Pinpointing and Quantifying Scalability Bottlenecks	22
4.4.1	Scalability Analysis Using Expectations	23
5	Running Applications with <code>hpcrun</code> and <code>hpclink</code>	27
5.1	Using <code>hpcrun</code>	27
5.2	Using <code>hpclink</code>	28
5.3	Hardware Counter Event Names	28
5.4	Sample Sources	29
5.4.1	Linux <code>perf_events</code>	29
5.4.2	PAPI	32
5.4.3	WALLCLOCK, REALTIME and CPUTIME	34
5.4.4	IO	35
5.4.5	MEMLEAK	36

5.5	Process Fraction	37
5.6	Starting and Stopping Sampling	37
5.7	Environment Variables for <code>hpcrun</code>	39
5.8	Platform-Specific Notes	39
5.8.1	Cray Systems	39
5.8.2	Blue Gene/Q Systems	40
5.8.3	ARM Systems	40
6	Measurement and Analysis of GPU-accelerated Applications	41
6.1	NVIDIA GPUs	41
6.1.1	PC Sampling on NVIDIA GPUs	43
6.1.2	Attributing Measurements to Source Code for NVIDIA GPUs	45
6.1.3	GPU Calling Context Tree Reconstruction	47
6.2	AMD GPUs	51
6.3	Intel GPUs	51
7	Measurement and Analysis of OpenMP Multithreading	52
8	The <code>hpcviewer</code> User Interface	53
8.1	Launching	53
8.2	Profile View	54
8.3	Panes	56
8.3.1	Source Pane	56
8.3.2	Navigation Pane	56
8.3.3	Metric Pane	58
8.4	Trace view	59
8.4.1	Trace View	62
8.4.2	Depth View	63
8.4.3	Summary View	63
8.4.4	Call Path View	63
8.4.5	Mini Map View	64
8.5	GPU Threads	64
8.6	Understanding Metrics	65
8.6.1	How Metrics are Computed	65
8.6.2	Example	65
8.7	Derived Metrics	67
8.7.1	Formulae	68
8.7.2	Examples	68
8.7.3	Creating Derived Metrics	68
8.8	Thread-level Metric Values	70
8.8.1	Plotting Graphs	70
8.8.2	Thread View	71
8.9	Filtering Tree Nodes	72
8.10	Convenience Features	74
8.10.1	Editor Pane	74
8.10.2	Metric Pane	74

8.11 Menus	75
8.11.1 File	75
8.11.2 Filter	76
8.11.3 View	77
8.11.4 Help	77
8.12 Limitations	77
9 Monitoring MPI Applications	78
9.1 Running and Analyzing MPI Programs	78
9.2 Building and Installing HPCTOOLKIT	80
10 Monitoring Statically Linked Applications	81
10.1 Linking with <code>hpmlink</code>	81
10.2 Running a Statically Linked Binary	82
10.3 Troubleshooting	82
11 Known Issues	84
11.1 A confusing label for GPU theoretical occupancy	84
11.2 Deadlock when using Darshan	84
11.3 Deadlock when using UCX	85
12 FAQ and Troubleshooting	87
12.1 Why do I see partial unwinds?	87
12.2 How do I choose <code>hpcrun</code> sampling periods?	88
12.3 <code>hpcrun</code> incurs high overhead! Why?	88
12.4 Fail to run <code>hpcviewer</code> : executable launcher was unable to locate its companion shared library	89
12.5 Mac only: <code>hpcviewer</code> runs on Java X instead of “Java 8”	89
12.6 When executing <code>hpcviewer</code> , it complains cannot create “Java Virtual Machine”	89
12.7 <code>hpcviewer</code> fails to launch due to <code>java.lang.NoSuchMethodError</code> exception.	90
12.8 <code>hpcviewer</code> fails due to <code>java.lang.OutOfMemoryError</code> exception.	90
12.9 <code>hpcviewer</code> writes a long list of Java error messages to the terminal!	90
12.10 <code>hpcviewer</code> attributes performance information only to functions and not to source code loops and lines! Why?	91
12.11 <code>hpcviewer</code> hangs trying to open a large database! Why?	91
12.12 <code>hpcviewer</code> runs glacially slowly! Why?	92
12.13 <code>hpcviewer</code> does not show my source code! Why?	92
12.13.1 Follow ‘Best Practices’	92
12.13.2 Additional Background	93
12.14 <code>hpcviewer</code> ’s reported line numbers do not exactly correspond to what I see in my source code! Why?	94
12.15 <code>hpcviewer</code> claims that there are several calls to a function within a particular source code scope, but my source code only has one! Why?	94
12.16 Trace view shows lots of white space on the left. Why?	95
12.17 I get a message about “Unable to find HPCTOOLKIT root directory”	95
12.18 Some of my syscalls return EINTR when run under <code>hpcrun</code>	95

12.19 How do I debug HPCTOOLKIT’s measurement?	96
12.19.1 Tracing <code>libmonitor</code>	96
12.19.2 Tracing HPCTOOLKIT’s Measurement Subsystem	96
12.19.3 Using a debugger to inspect an execution being monitored by HPCTOOLKIT	97
12.19.4 Using <code>hpclink</code> with <code>cmake</code>	98
A Environment Variables	101
A.1 Environment Variables for Users	101
A.2 Environment Variables for Developers	103

Chapter 1

Introduction

HPCTOOLKIT [1, 13] is an integrated suite of tools for measurement and analysis of program performance on computers ranging from multicore desktop systems to the world’s largest supercomputers. HPCTOOLKIT provides accurate measurements of a program’s work, resource consumption, and inefficiency, correlates these metrics with the program’s source code, works with multilingual, fully optimized binaries, has low measurement overhead, and scales to large parallel systems. HPCTOOLKIT’s measurements provide support for analyzing a program execution cost, inefficiency, and scaling characteristics both within and across nodes of a parallel system.

HPCTOOLKIT principally monitors an execution of a multithreaded and/or multiprocess program using asynchronous sampling, unwinding thread call stacks, and attributing the metric value associated with a sample event in a thread to the calling context of the thread/process in which the event occurred. HPCTOOLKIT’s asynchronous sampling is typically triggered by the expiration of a Linux timer or a hardware performance monitoring unit event, such reaching a threshold value for a hardware performance counter. Sampling has several advantages over instrumentation for measuring program performance: it requires no modification of source code, it avoids potential blind spots (such as code available in only binary form), and it has lower overhead. HPCTOOLKIT typically adds measurement overhead of only a few percent to an execution for reasonable sampling rates [17]. Sampling enables fine-grain measurement and attribution of costs in both serial and parallel programs.

For parallel programs, one can use HPCToolkit to measure the fraction of time threads are idle, working, or communicating. To obtain detailed information about a program’s computation performance, one can collect samples using a processor’s built-in performance monitoring units to measure metrics such as operation counts, pipeline stalls, cache misses, and data movement between processor sockets. Such detailed measurements are essential to understand the performance characteristics of applications on modern multicore microprocessors that employ instruction-level parallelism, out-of-order execution, and complex memory hierarchies. With HPCTOOLKIT, one can also easily compute derived metrics such as cycles per instruction, waste, and relative efficiency to provide insight into a program’s shortcomings.

A unique capability of HPCTOOLKIT is its ability to unwind the call stack of a thread executing highly optimized code to attribute time, hardware counter metrics, as well as

software metrics (e.g., context switches) to a full calling context. Call stack unwinding is often difficult for highly optimized code [17]. For accurate call stack unwinding, HPCToolkit employs two strategies: interpreting compiler-recorded information in DWARF Frame Descriptor Entries (FDEs) and binary analysis to compute unwind recipes directly from an application’s machine instructions. On ARM processors, HPCToolkit uses `libunwind` exclusively. On Power processors, HPCToolkit uses binary analysis exclusively. On x86_64 processors, HPCToolkit employs both strategies in an integrated fashion.

HPCTOOLKIT assembles performance measurements into a call path profile that associates the costs of each function call with its full calling context. In addition, HPCTOOLKIT uses binary analysis to attribute program performance metrics with uniquely detailed precision – full dynamic calling contexts augmented with information about call sites, inlined functions and templates, loops, and source lines. Measurements can be analyzed in a variety of ways: top-down in a calling context tree, which associates costs with the full calling context in which they are incurred; bottom-up in a view that apportions costs associated with a function to each of the contexts in which the function is called; and in a flat view that aggregates all costs associated with a function independent of calling context. This multiplicity of code-centric perspectives is essential to understanding a program’s performance for tuning under various circumstances. HPCTOOLKIT also supports a thread-centric perspective, which enables one to see how a performance metric for a calling context differs across threads, and a time-centric perspective, which enables a user to see how an execution unfolds over time. Figures 1.1–1.3 show samples of HPCToolkit’s code-centric, thread-centric, and time-centric views.

By working at the machine-code level, HPCTOOLKIT accurately measures and attributes costs in executions of multilingual programs, even if they are linked with libraries available only in binary form. HPCTOOLKIT supports performance analysis of fully optimized code – the only form of a program worth measuring; it even measures and attributes performance metrics to shared libraries that are dynamically loaded at run time. The low overhead of HPCTOOLKIT’s sampling-based measurement is particularly important for parallel programs because measurement overhead can distort program behavior.

HPCTOOLKIT is also especially good at pinpointing scaling losses in parallel codes, both within multicore nodes and across the nodes in a parallel system. Using differential analysis of call path profiles collected on different numbers of threads or processes enables one to quantify scalability losses and pinpoint their causes to individual lines of code executed in particular calling contexts [4]. We have used this technique to quantify scaling losses in leading science applications across thousands of processor cores on Cray and IBM Blue Gene systems, associate them with individual lines of source code in full calling context [15, 18], and quantify scaling losses in science applications within compute nodes at the loop nest level due to competition for memory bandwidth in multicore processors [14]. We have also developed techniques for efficiently attributing the idleness in one thread to its cause in another thread [16, 20].

HPCTOOLKIT is deployed on many DOE supercomputers, including the Sierra supercomputer (IBM Power9 + NVIDIA V100 GPUs) at Lawrence Livermore National Laboratory, Cray XC40 systems at Argonne’s Leadership Computing Facility and the National Energy Research Scientific Computing Center; the Summit supercomputer (IBM Power9 + NVIDIA V100 GPUs) at Oak Ridge Leadership Computing Facility, Blue Gene/Q systems



Figure 1.1: A code-centric view of an execution of the University of Chicago’s FLASH code executing on 8192 cores of a Blue Gene/P. This bottom-up view shows that 16% of the execution time was spent in IBM’s DCMF messaging layer. By tracking these costs up the call chain, we can see that most of this time was spent on behalf of calls to `pmmpi_allreduce` on line 419 of `amr_comm_setup`.

at Argonne Leadership Computing Facility, as well as other clusters and supercomputers based on x86_64, Power, and ARM processors.



Figure 1.2: A thread-centric view of the performance of a parallel radix sort application executing on 960 cores of a Cray XE6. The bottom pane shows a calling context for `usort` in the execution. The top pane shows a graph of how much time each thread spent executing calls to `usort` from the highlighted context. On a Cray XE6, there is one MPI helper thread for each compute node in the system; these helper threads spent no time executing `usort`. The graph shows that some of the MPI ranks spent twice as much time in `usort` as others. This happens because the radix sort divides up the work into 1024 buckets. In an execution on 960 cores, 896 cores work on one bucket and 64 cores work on two. The middle pane shows an alternate view of the thread-centric data as a histogram.



Figure 1.3: A time-centric view of part of an execution of the University of Chicago’s FLASH code on 256 cores of a Blue Gene/P. The figure shows a detail from the end of the initialization phase and part of the first iteration of the solve phase. The largest pane in the figure shows the activity of cores 2–95 in the execution during a time interval ranging from 69.376s–85.58s during the execution. Time lines for threads are arranged from top to bottom and time flows from left to right. The color at any point in time for a thread indicates the procedure that the thread is executing at that time. The right pane shows the full call stack of thread 85 at 84.82s into the execution, corresponding to the selection shown by the white crosshair; the outermost procedure frame of the call stack is shown at the top of the pane and the innermost frame is shown at the bottom. This view highlights that even though FLASH is an SPMD program, the behavior of threads over time can be quite different. The purple region highlighted by the cursor, which represents a call by all processors to `mpi_allreduce`, shows that the time spent in this call varies across the processors. The variation in time spent waiting in `mpi_allreduce` is readily explained by an imbalance in the time processes spend a prior prolongation step, shown in yellow. Further left in the figure, one can see differences among ranks executing on different cores in each node as they await the completion of an `mpi_allreduce`. A rank executing on one core of each node waits in `DCMF_Messager_advance` (which appears as blue stripes) while ranks executing on other cores in each node wait in a helper function (shown in green). In this phase, ranks await the delayed arrival of a few of their peers who have extra work to do inside `simulation_initblock` before they call `mpi_allreduce`.

Chapter 2

HPCToolkit Overview

HPCTOOLKIT’s work flow is organized around four principal capabilities, as shown in Figure 2.1:

1. *measurement* of context-sensitive performance metrics using call-stack unwinding while an application executes;
2. *binary analysis* to recover program structure from application binaries;
3. *attribution* of performance metrics by correlating dynamic performance metrics with static program structure; and
4. *presentation* of performance metrics and associated source code.

To use HPCTOOLKIT to measure and analyze an application’s performance, one first compiles and links the application for a production run, using *full* optimization and including debugging symbols.¹ Second, one launches an application with HPCTOOLKIT’s measurement tool, `hpcrun`, which uses statistical sampling to collect a performance profile. Third, one invokes `hpcstruct`, HPCTOOLKIT’s tool for analyzing an application binary to recover information about files, procedures, loops, and inlined code. Fourth, one uses `hpcprof` to combine information about an application’s structure with dynamic performance measurements to produce a performance database. Finally, one explores a performance database with HPCTOOLKIT’s `hpcviewer` and/or Trace view graphical presentation tools.

The rest of this chapter briefly discusses unique aspects of HPCTOOLKIT’s measurement, analysis and presentation capabilities.

2.1 Asynchronous Sampling

Without accurate measurement, performance analysis results may be of questionable value. As a result, a principal focus of work on HPCTOOLKIT has been the design and

¹For the most detailed attribution of application performance data using HPCTOOLKIT, one should ensure that the compiler includes line map information in the object code it generates. While HPCTOOLKIT does not need this information to function, it can be helpful to users trying to interpret the results. Since compilers can usually provide line map information for fully optimized code, this requirement need not require a special build process. For instance, with the Intel compiler we recommend using `-g -debug inline_debug_info`.



Figure 2.1: Overview of HPCTOOLKIT’s tool work flow.

implementation of techniques to provide accurate fine-grain measurements of production applications running at scale. For tools to be useful on production applications on large-scale parallel systems, large measurement overhead is unacceptable. For measurements to be accurate, performance tools must avoid introducing measurement error. Both source-level and binary instrumentation can distort application performance through a variety of mechanisms [11]. Frequent calls to small instrumented procedures can lead to considerable measurement overhead. Furthermore, source-level instrumentation can distort application performance by interfering with inlining and template optimization. To avoid these effects, many instrumentation-based tools intentionally refrain from instrumenting certain procedures. Ironically, the more this approach reduces overhead, the more it introduces *blind spots*, i.e., intervals of unmonitored execution. For example, a common selective instrumentation technique is to ignore small frequently executed procedures — but these may be just the thread synchronization library routines that are critical. Sometimes, a tool unintentionally introduces a blind spot. A typical example is that source code instrumentation suffers from blind spots when source code is unavailable, a common condition for math and communication libraries.

To avoid these problems, HPCTOOLKIT eschews instrumentation and favors the use of *asynchronous sampling* to measure and attribute performance metrics. During a program execution, sample events are triggered by periodic interrupts induced by an interval timer or overflow of hardware performance counters. One can sample metrics that reflect work (e.g., instructions, floating-point operations), consumption of resources (e.g., cycles, bandwidth consumed in the memory hierarchy by data transfers in response to cache misses), or inefficiency (e.g., stall cycles). For reasonable sampling frequencies, the overhead and distortion introduced by sampling-based measurement is typically much lower than that introduced by instrumentation [6].

2.2 Call Path Profiling

For all but the most trivially structured programs, it is important to associate the costs incurred by each procedure with the contexts in which the procedure is called. Know-

ing the context in which each cost is incurred is essential for understanding why the code performs as it does. This is particularly important for code based on application frameworks and libraries. For instance, costs incurred for calls to communication primitives (e.g., `MPI_Wait`) or code that results from instantiating C++ templates for data structures can vary widely depending how they are used in a particular context. Because there are often layered implementations within applications and libraries, it is insufficient either to insert instrumentation at any one level or to distinguish costs based only upon the immediate caller. For this reason, HPCTOOLKIT uses call path profiling to attribute costs to the full calling contexts in which they are incurred.

HPCTOOLKIT’s `hpcrun` call path profiler uses call stack unwinding to attribute execution costs of optimized executables to the full calling context in which they occur. Unlike other tools, to support asynchronous call stack unwinding during execution of optimized code, `hpcrun` uses on-line binary analysis to locate procedure bounds and compute an unwind recipe for each code range within each procedure [17]. These analyses enable `hpcrun` to unwind call stacks for optimized code with little or no information other than an application’s machine code.

2.3 Recovering Static Program Structure

To enable effective analysis, call path profiles for executions of optimized programs must be correlated with important source code abstractions. Since measurements refer only to instruction addresses within an executable, it is necessary to map measurements back to the program source. To associate measurement data with the static structure of fully-optimized executables, we need a mapping between object code and its associated source code structure.² HPCTOOLKIT constructs this mapping using binary analysis; we call this process *recovering program structure* [17].

HPCTOOLKIT focuses its efforts on recovering procedures, inlined functions and templates, as well as loop nests, the most important elements of source code structure. To recover program structure, HPCTOOLKIT’s `hpcstruct` utility parses a load module’s machine instructions, reconstructs a control flow graph, combines line map and DWARF information about inlining with interval analysis on the control flow graph in a way that enables it to relate machine code after optimization back to the original source [17].

Two important benefits naturally accrue from this approach. First, HPCTOOLKIT can expose the structure of and assign metrics to what is actually executed, *even if source code is unavailable*. For example, `hpcstruct`’s program structure naturally reveals transformations such as loop fusion and scalarization loops that arise from compilation of Fortran 90 array notation. Similarly, it exposes calls to compiler support routines and wait loops in communication libraries of which one would otherwise be unaware. Second, we combine (post-mortem) the recovered static program structure with dynamic call paths to expose inlined frames and loop nests. This enables us to attribute the performance of samples in their full static and dynamic context and correlate it with source code.

²This object to source code mapping should be contrasted with the binary’s line map, which (if present) is typically fundamentally line based.

2.4 Presenting Performance Measurements

To enable an analyst to rapidly pinpoint and quantify performance bottlenecks, tools must present the performance measurements in a way that engages the analyst, focuses attention on what is important, and automates common analysis subtasks to reduce the mental effort and frustration of sifting through a sea of measurement details.

To enable rapid analysis of an execution’s performance bottlenecks, we have carefully designed the **hpcviewer** - a code-centric presentation tool [2] and Trace view - a time-centric presentation tool [19].

hpcviewer combines a relatively small set of complementary presentation techniques that, taken together, rapidly focus an analyst’s attention on performance bottlenecks rather than on unimportant information. To facilitate the goal of rapidly focusing an analyst’s attention on performance bottlenecks **hpcviewer** extends several existing presentation techniques. In particular, **hpcviewer** (1) synthesizes and presents three complementary views of calling-context-sensitive metrics; (2) treats a procedure’s static structure as first-class information with respect to both performance metrics and constructing views; (3) enables a large variety of user-defined metrics to describe performance inefficiency; and (4) automatically expands hot paths based on arbitrary performance metrics — through calling contexts and static structure — to rapidly highlight important performance data.

Trace view enables an application developer to visualize how a parallel execution unfolds over time. This view facilitates identification of important inefficiencies such as serialization and load imbalance, among others.

Chapter 3

Quick Start

This chapter provides a rapid overview of analyzing the performance of an application using HPCTOOLKIT. It assumes an operational installation of HPCTOOLKIT.

3.1 Guided Tour

HPCTOOLKIT’s work flow is summarized in Figure 3.1 (on page 11) and is organized around four principal capabilities:

1. *measurement* of context-sensitive performance metrics while an application executes;
2. *binary analysis* to recover program structure from application binaries;
3. *attribution* of performance metrics by correlating dynamic performance metrics with static program structure; and
4. *presentation* of performance metrics and associated source code.

To use HPCTOOLKIT to measure and analyze an application’s performance, one first compiles and links the application for a production run, using *full* optimization. Second, one launches an application with HPCTOOLKIT’s measurement tool, `hpcrun`, which uses statistical sampling to collect a performance profile. Third, one invokes `hpcstruct`, HPCTOOLKIT’s tool for analyzing an application binary to recover information about files, procedures, loops, and inlined code. Fourth, one uses `hpcprof` to combine information about an application’s structure with dynamic performance measurements to produce a performance database. For large executions where analyzing performance measurements serially would be imprudent, HPCTOOLKIT provides `hpcprof-mpi` - an MPI program that can be launched to analyze performance data from many MPI ranks and threads in parallel. Finally, one explores a performance database with one of HPCTOOLKIT’s graphical user interfaces: `hpcviewer` for code-centric analysis of performance metrics or Trace view for time-centric analysis of an execution.

The following subsections explain HPCTOOLKIT’s work flow in more detail.



Figure 3.1: Overview of HPCTOOLKIT tool’s work flow.

3.1.1 Compiling an Application

For the most detailed attribution of application performance data using HPCTOOLKIT, one should compile so as to include with line map information in the generated object code. This usually means compiling with options similar to ‘`-g -O3`’. Check your compiler’s documentation for information about the right set of options to have the compiler record information about inlining and the mapping of machine instructions to source lines. We advise picking options that indicate they will record information that relates machine instructions to source code without compromising optimization. For instance, the Portland Group (PGI) compilers, use `-gopt` in place of `-g` to collect information without interfering with optimization.

While HPCTOOLKIT does not need information about the mapping between machine instructions and source code to function, having such information included in the binary code by the compiler can be helpful to users trying to interpret performance measurements. Since compilers can usually provide information about line mappings and inlining for fully-optimized code, this requirement usually involves a one-time trivial adjustment to the an application’s build scripts to provide a better experience with tools. Such mapping information enables tools such as HPCTOOLKIT, race detectors, and memory analysis tools to attribute information more precisely.

3.1.2 Measuring Application Performance

Measurement of application performance takes two different forms depending on whether your application is dynamically or statically linked. To monitor a dynamically linked application, simply use `hpcrun` to launch the application. To monitor a statically linked application, such as those typically used on Blue Gene and Cray supercomputers, link your application using `hpclink`. In either case, the application may be sequential, multithreaded or based on MPI. The commands below give examples for an application named `app`.

- Dynamically linked applications:

Simply launch your application with `hpcrun`:

```
[<mpi-launcher>] hpcrun [hpcrun-options] app [app-arguments]
```

Of course, `<mpi-launcher>` is only needed for MPI programs and is sometimes a program like `mpixexec` or `mpirun`, or a workload manager's utilities such as Slurm's `srun` or IBM's Job Step Manager utility `jsrun`.

- Statically linked applications:

First, link `hpcrun`'s monitoring code into `app`, using `hpmlink`:

```
hpmlink <linker> -o app <linker-arguments>
```

Then monitor `app` by passing `hpcrun` options through environment variables. For instance:

```
export HPCRUN_EVENT_LIST="CYCLES"  
[<mpi-launcher>] app [app-arguments]
```

`hpmlink`'s `--help` option gives a list of environment variables that affect monitoring. See Chapter 10 for more information.

Any of these commands will produce a measurements database that contains separate measurement information for each MPI rank and thread in the application. The database is named according the form:

```
hpctoolkit-app-measurements [-<jobid>]
```

If the application `app` is run under control of a recognized batch job scheduler (such as Slurm, Cobalt, or IBM's Job Manager), the name of the measurements directory will contain the corresponding job identifier `<jobid>`. Currently, the database contains measurements files for each thread that are named using the following templates:

```
app-<mpi-rank>-<thread-id>-<host-id>-<process-id>.<generation-id>.hpcrun  
app-<mpi-rank>-<thread-id>-<host-id>-<process-id>.<generation-id>.hpctrace
```

Specifying Sample Sources

HPCTOOLKIT primarily monitors an application using asynchronous sampling. Consequently, the most common option to `hpcrun` is a list of sample sources that define how samples are generated. A sample source takes the form of an event name `e` and `howoften`, specified as `e@howoften`. The specifier `howoften` may be a number, indicating a period, e.g. `CYCLES@4000001` or it may be `f` followed by a number, `CYCLES@f200` indicating a frequency in samples/second. For a sample source with event `e` and period `p`, after every `p` instances of `e`, a sample is generated that causes `hpcrun` to inspect the and record information about the monitored application.

To configure `hpcrun` with two samples sources, `e1@howoften1` and `e2@howoften2`, use the following options:

```
--event e1@howoften1 --event e2@howoften2
```

To use the same sample sources with an `hpmlink`-ed application, use a command similar to:

```
export HPCRUN_EVENT_LIST="e1@howoften1 e2@howoften2"
```

3.1.3 Recovering Program Structure

To recover static program structure for the application `app`, use the command:

```
hpcstruct app
```

This command analyzes `app`'s binary and computes a representation of its static source code structure, including its loop nesting structure. One can accelerate analysis of a large application binary by using the '`-j n`' option to `hpcstruct`, which indicates that analysis should be performed using n threads. The command saves this information in a file named `app.hpcstruct` that should be passed to `hpcprof` with the `-S/--structure` argument.

Typically, `hpcstruct` is launched without any options.

3.1.4 Analyzing Measurements & Attributing Them to Source Code

To analyze HPCTOOLKIT's measurements and attribute them to the application's source code, use either `hpcprof` or `hpcprof-mpi`. In most respects, `hpcprof` and `hpcprof-mpi` are semantically identical. For convenience, we use the notation `hpcprof/mpi` to refer to both of these tools. Both generate the same set of summary metrics over all threads and processes in an execution. The difference between the two is that the latter is designed to process (in parallel) measurements from large-scale executions. Consequently, while the former can optionally generate separate metrics for each thread (see the `--metric/-M` option), the latter only generates summary metrics. However, the latter can also generate additional information for plotting thread-level metric values (see Section 8.8.1).

`hpcprof` is typically used as follows:

```
hpcprof -S app.hpcstruct -I <app-src>/+ \
    hpctoolkit-app-measurements1 [hpctoolkit-app-measurements2 ...]
```

and `hpcprof-mpi` is analogous:

```
<mpi-launcher> hpcprof-mpi \
    -S app.hpcstruct -I <app-src>/+ \
    hpctoolkit-app-measurements1 [hpctoolkit-app-measurements2 ...]
```

Either command will produce an HPCTOOLKIT performance database with the name `hpctoolkit-app-database`. If this database directory already exists, `hpcprof/mpi` will form a unique name using a numerical qualifier.

Both `hpcprof/mpi` can collate multiple measurement databases, as long as they are gathered against the same binary. This capability is useful for (a) combining event sets gathered over multiple executions and (b) performing scalability studies (see Section 4.4).

The above commands use two important options. The `-S/--structure` option takes a program structure file. The `-I/--include` option takes a directory `<app-src>` to application source code; the optional '+' suffix requests that the directory be searched recursively for source code. Either option can be passed multiple times to specify multiple structure files (e.g., for the application and each of the key libraries it uses) or multiple include paths that indicate roots of source trees for the application and/or of its libraries.

Another potentially important option, especially for machines that require executing from special file systems, is the `-R/--replace-path` option for substituting instances of *old-path* with *new-path*: `-R 'old-path=new-path'`.

A possibly important detail about the above command is that source code should be considered an `hpcprof/mpi` input. This is critical when using a machine that restricts executions to a scratch parallel file system. In such cases, not only must you copy `hpcprof-mpi` into the scratch file system, but also all source code that you want `hpcprof-mpi` to find and copy into the resulting Experiment database.

3.1.5 Presenting Performance Measurements for Interactive Analysis

To interactively view and analyze an HPCTOOLKIT performance database, use `hpcviewer`. `hpcviewer` may be launched from the command line or by double-clicking on its icon on MacOS or Windows. The following is an example of launching from a command line:

```
hpcviewer hpctoolkit-app-database
```

Additional help for `hpcviewer` can be found in a help pane available from `hpcviewer`'s *Help* menu.

3.1.6 Effective Performance Analysis Techniques

To effectively analyze application performance, consider using one of the following strategies, which are described in more detail in Chapter 4.

- A waste metric, which represents the difference between achieved performance and potential peak performance is a good way of understanding the potential for tuning the node performance of codes (Section 4.3). `hpcviewer` supports synthesis of derived metrics to aid analysis. Derived metrics are specified within `hpcviewer` using spreadsheet-like formula. See the `hpcviewer` help pane for details about how to specify derived metrics.
- Scalability bottlenecks in parallel codes can be pinpointed by differential analysis of two profiles with different degrees of parallelism (Section 4.4).

The following sketches the mechanics of performing a simple scalability study between executions *x* and *y* of an application *app*:

```
hpcrun [options-x] app [app-arguments-x]           (execution x)
hpcrun [options-y] app [app-arguments-y]           (execution y)
hpcstruct app
hpcprof/mpi -S ... -I ... measurements-x measurements-y
hpcviewer hpctoolkit-database          (compute a scaling-loss metric)
```

3.2 Additional Guidance

For additional information, consult the rest of this manual and other documentation: First, we summarize the available documentation and command-line help:

Command-line help.

Each of HPCTOOLKIT's command-line tools can generate a help message summarizing the tool's usage, arguments and options. To generate this help message, invoke the tool with `-h` or `--help`.

Man pages.

Man pages are available either via the Internet (<http://hpctoolkit.org/documentation.html>) or from a local HPCTOOLKIT installation (<hpctoolkit-installation>/share/man).

Manuals.

Manuals are available either via the Internet (<http://hpctoolkit.org/documentation.html>) or from a local HPCTOOLKIT installation (<hpctoolkit-installation>/share/doc/hpctoolkit/documentation.html).

Articles and Papers.

There are a number of articles and papers that describe various aspects of HPCTOOLKIT's measurement, analysis, attribution and presentation technology. They can be found at <http://hpctoolkit.org/publications.html>.

Chapter 4

Effective Strategies for Analyzing Program Performance

This chapter describes some proven strategies for using performance measurements to identify performance bottlenecks in both serial and parallel codes.

4.1 Monitoring High-Latency Penalty Events

A very simple and often effective methodology is to profile with respect to cycles and high-latency penalty events. If HPCTOOLKIT attributes a large number of penalty events with a particular source-code statement, there is an extremely high likelihood of significant exposed stalling. This is true even though (1) modern out-of-order processors can overlap the stall latency of one instruction with nearby independent instructions and (2) some penalty events “over count”.¹ If a source-code statement incurs a large number of penalty events and it also consumes a non-trivial amount of cycles, then this region of code is an opportunity for optimization. Examples of good penalty events are last-level cache misses and TLB misses.

4.2 Computing Derived Metrics

Modern computer systems provide access to a rich set of hardware performance counters that can directly measure various aspects of a program’s performance. Counters in the processor core and memory hierarchy enable one to collect measures of work (e.g., operations performed), resource consumption (e.g., cycles), and inefficiency (e.g., stall cycles). One can also measure time using system timers.

Values of individual metrics are of limited use by themselves. For instance, knowing the count of cache misses for a loop or routine is of little value by itself; only when combined with other information such as the number of instructions executed or the total number of cache accesses does the data become informative. While a developer might not mind using mental arithmetic to evaluate the relationship between a pair of metrics for a particular program scope (e.g., a loop or a procedure), doing this for many program scopes is exhausting.

¹For example, performance monitoring units often categorize a prefetch as a cache miss.



Figure 4.1: Computing a derived metric (cycles per instruction) in hpcviewer.

To address this problem, `hpcviewer` supports calculation of derived metrics. `hpcviewer` provides an interface that enables a user to specify spreadsheet-like formula that can be used to calculate a derived metric for every program scope.

Figure 4.1 shows how to use `hpcviewer` to compute a *cycles/instruction* derived metric from measured metrics `PAPI_TOT_CYC` and `PAPI_TOT_INS`; these metrics correspond to *cycles* and *total instructions executed* measured with the PAPI hardware counter interface. To compute a derived metric, one first depresses the button marked $f(x)$ above the metric pane; that will cause the pane for computing a derived metric to appear. Next, one types in the formula for the metric of interest. When specifying a formula, existing columns of metric data are referred to using a positional name $\$n$ to refer to the n^{th} column, where the first column is written as $\$0$. The metric pane shows the formula $\$1/\3 . Here, $\$1$ refers to the column of data representing the exclusive value for `PAPI_TOT_CYC` and $\$3$ refers to the column of data representing the exclusive value for `PAPI_TOT_INS`.² Positional names for

²An *exclusive* metric for a scope refers to the quantity of the metric measured for that scope alone; an *inclusive* metric for a scope represents the value measured for that scope as well as costs incurred by



Figure 4.2: Displaying the new *cycles/instruction* derived metric in hpcviewer.

metrics you use in your formula can be determined using the *Metric* pull-down menu in the pane. If you select your metric of choice using the pull-down, you can insert its positional name into the formula using the *insert metric* button, or you can simply type the positional name directly into the formula.

At the bottom of the derived metric pane, one can specify a name for the new metric. One also has the option to indicate that the derived metric column should report for each scope what percent of the total its quantity represents; for a metric that is a ratio, computing a percent of the total is not meaningful, so we leave the box unchecked. After clicking the OK button, the derived metric pane will disappear and the new metric will appear as the rightmost column in the metric pane. If the metric pane is already filled with other columns of metric, you may need to scroll right in the pane to see the new metric. Alternatively, you can use the metric check-box pane (selected by depressing the button to the right of $f(x)$ above the metric pane) to hide some of the existing metrics so that there will be enough

any functions it calls. In hpcviewer, inclusive metric columns are marked with "(I)" and exclusive metric columns are marked with "(E)."

room on the screen to display the new metric. Figure 4.2 shows the resulting `hpcviewer` display after clicking OK to add the derived metric.

The following sections describe several types of derived metrics that are of particular use to gain insight into performance bottlenecks and opportunities for tuning.

4.3 Pinpointing and Quantifying Inefficiencies

While knowing where a program spends most of its time or executes most of its floating point operations may be interesting, such information may not suffice to identify the biggest targets of opportunity for improving program performance. For program tuning, it is less important to know how much resources (e.g., time, instructions) were consumed in each program context than knowing where resources were consumed *inefficiently*.

To identify performance problems, it might initially seem appealing to compute ratios to see how many events per cycle occur in each program context. For instance, one might compute ratios such as FLOPs/cycle, instructions/cycle, or cache miss ratios. However, using such ratios as a sorting key to identify inefficient program contexts can misdirect a user’s attention. There may be program contexts (e.g., loops) in which computation is terribly inefficient (e.g., with low operation counts per cycle); however, some or all of the least efficient contexts may not account for a significant amount of execution time. Just because a loop is inefficient doesn’t mean that it is important for tuning.

The best opportunities for tuning are where the aggregate performance losses are greatest. For instance, consider a program with two loops. The first loop might account for 90% of the execution time and run at 50% of peak performance. The second loop might account for 10% of the execution time, but only achieve 12% of peak performance. In this case, the total performance loss in the first loop accounts for 50% of the first loop’s execution time, which corresponds to 45% of the total program execution time. The 88% performance loss in the second loop would account for only 8.8% of the program’s execution time. In this case, tuning the first loop has a greater potential for improving the program performance even though the second loop is less efficient.

A good way to focus on inefficiency directly is with a derived *waste* metric. Fortunately, it is easy to compute such useful metrics. However, there is no one *right* measure of waste for all codes. Depending upon what one expects as the rate-limiting resource (e.g., floating-point computation, memory bandwidth, etc.), one can define an appropriate waste metric (e.g., FLOP opportunities missed, bandwidth not consumed) and sort by that.

For instance, in a floating-point intensive code, one might consider keeping the floating point pipeline full as a metric of success. One can directly quantify and pinpoint losses from failing to keep the floating point pipeline full *regardless of why this occurs*. One can pinpoint and quantify losses of this nature by computing a *floating-point waste* metric that is calculated as the difference between the potential number of calculations that could have been performed if the computation was running at its peak rate minus the actual number that were performed. To compute the number of calculations that could have been completed in each scope, multiply the total number of cycles spent in the scope by the peak rate of operations per cycle. Using `hpcviewer`, one can specify a formula to compute

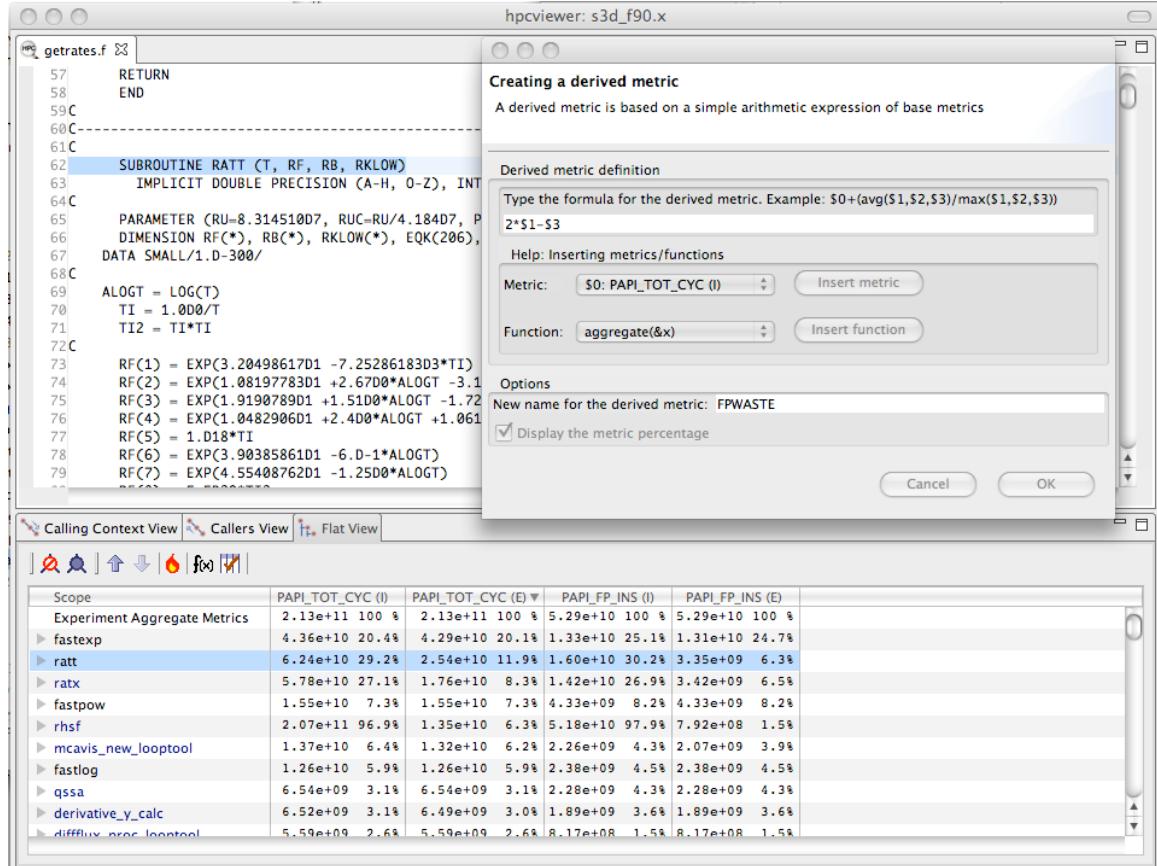


Figure 4.3: Computing a floating point waste metric in hpcviewer.

such a derived metric and it will compute the value of the derived metric for every scope. Figure 4.3 shows the specification of this floating-point waste metric for a code.³

Sorting by a waste metric will rank order scopes to show the scopes with the greatest waste. Such scopes correspond directly to those that contain the greatest opportunities for improving overall program performance. A waste metric will typically highlight loops where

- a lot of time is spent computing efficiently, but the aggregate inefficiencies accumulate,
- less time is spent computing, but the computation is rather inefficient, and
- scopes such as copy loops that contain no computation at all, which represent a complete waste according to a metric such as floating point waste.

Beyond identifying and quantifying opportunities for tuning with a waste metric, one can compute a companion derived metric *relative efficiency* metric to help understand how easy it might be to improve performance. A scope running at very high efficiency will typically be much harder to tune than running at low efficiency. For our floating-point

³Many recent processors have trouble accurately counting floating-point operations accurately, which is unfortunate. If your processor can't accurately count floating-point operations, a floating-point waste metric will be less useful.

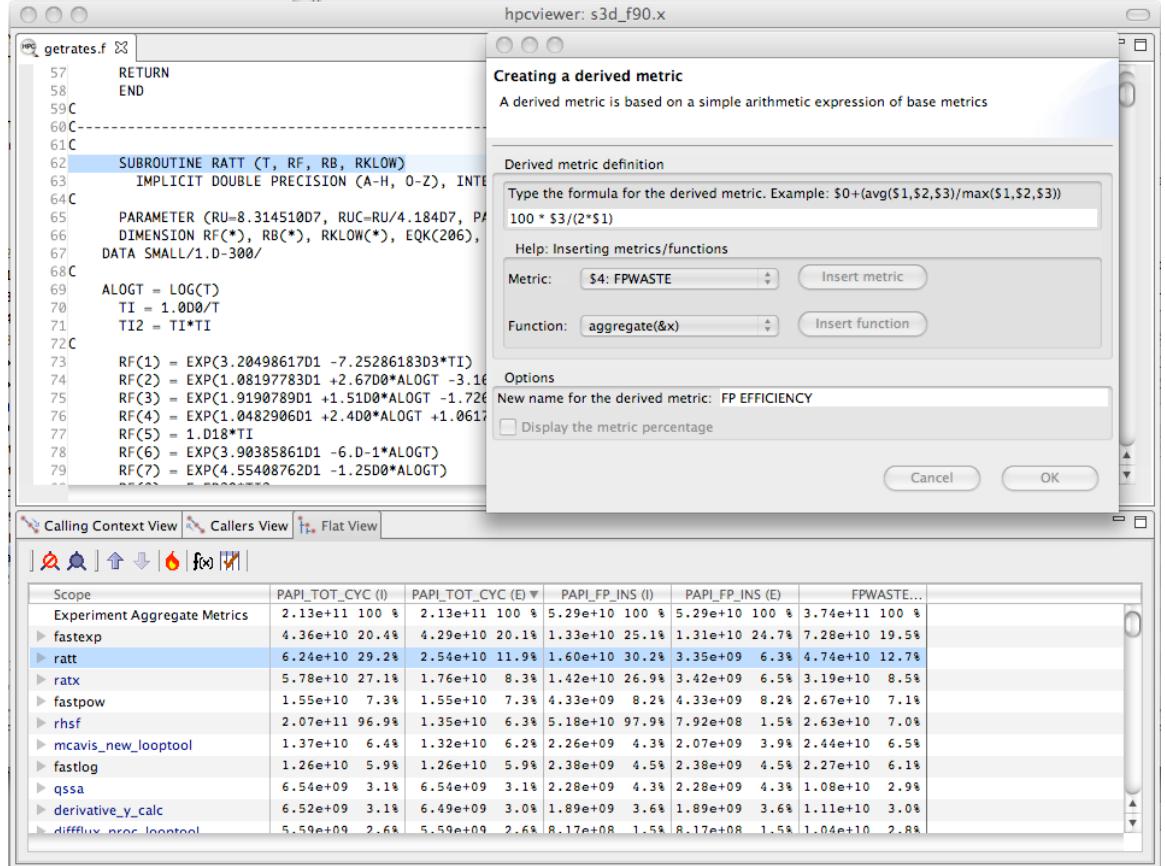


Figure 4.4: Computing floating point efficiency in percent using `hpcviewer`.

waste metric, we one can compute the floating point efficiency metric by dividing measured FLOPs by potential peak FLOPs and multiplying the quantity by 100. Figure 4.4 shows the specification of this floating-point efficiency metric for a code.

Scopes that rank high according to a waste metric and low according to a companion relative efficiency metric often make the best targets for optimization. Figure 4.5 shows the specification of this floating-point efficiency metric for a code. Figure 4.5 shows an `hpcviewer` display that shows the top two routines that collectively account for 32.2% of the floating point waste in a reactive turbulent combustion code. The second routine (`ratt`) is expanded to show the loops and statements within. While the overall floating point efficiency for `ratt` is at 6.6% of peak (shown in scientific notation in the `hpcviewer` display), the most costly loop in `ratt` that accounts for 7.3% of the floating point waste is executing at only 0.114% efficiency. Identifying such sources of inefficiency is the first step towards improving performance via tuning.



Figure 4.5: Using floating point waste and the percent of floating point efficiency to evaluate opportunities for optimization.

4.4 Pinpointing and Quantifying Scalability Bottlenecks

On large-scale parallel systems, identifying impediments to scalability is of paramount importance. On today’s systems fashioned out of multicore processors, two kinds of scalability are of particular interest:

- scaling within nodes, and
- scaling across the entire system.

HPCTOOLKIT can be used to readily pinpoint both kinds of bottlenecks. Using call path profiles collected by `hpcrun`, it is possible to quantify and pinpoint scalability bottlenecks of any kind, *regardless of cause*.

To pinpoint scalability bottlenecks in parallel programs, we use *differential profiling* — mathematically combining corresponding buckets of two or more execution profiles. Differential profiling was first described by McKenney [10]; he used differential profiling to

compare two *flat* execution profiles. Differencing of flat profiles is useful for identifying what parts of a program incur different costs in two executions. Building upon McKenney’s idea of differential profiling, we compare call path profiles of parallel executions at different scales to pinpoint scalability bottlenecks. Differential analysis of call path profiles pinpoints not only differences between two executions (in this case scalability losses), but the contexts in which those differences occur. Associating changes in cost with full calling contexts is particularly important for pinpointing context-dependent behavior. Context-dependent behavior is common in parallel programs. For instance, in message passing programs, the time spent by a call to `MPI_Wait` depends upon the context in which it is called. Similarly, how the performance of a communication event scales as the number of processors in a parallel execution increases depends upon a variety of factors such as whether the size of the data transferred increases and whether the communication is collective or not.

4.4.1 Scalability Analysis Using Expectations

Application developers have expectations about how the performance of their code should scale as the number of processors in a parallel execution increases. Namely,

- when different numbers of processors are used to solve the same problem (strong scaling), one expects an execution’s speedup to increase linearly with the number of processors employed;
- when different numbers of processors are used but the amount of computation per processor is held constant (weak scaling), one expects the execution time on a different number of processors to be the same.

In both of these situations, a code developer can express their expectations for how performance will scale as a formula that can be used to predict execution performance on a different number of processors. One’s expectations about how overall application performance should scale can be applied to each context in a program to pinpoint and quantify deviations from expected scaling. Specifically, one can scale and difference the performance of an application on different numbers of processors to pinpoint contexts that are not scaling ideally.

To pinpoint and quantify scalability bottlenecks in a parallel application, we first use `hpcrun` to collect call path profile for an application on two different numbers of processors. Let E_p be an execution on p processors and E_q be an execution on q processors. Without loss of generality, assume that $q > p$.

In our analysis, we consider both *inclusive* and *exclusive* costs for CCT nodes. The inclusive cost at n represents the sum of all costs attributed to n and any of its descendants in the CCT, and is denoted by $I(n)$. The exclusive cost at n represents the sum of all costs attributed strictly to n , and we denote it by $E(n)$. If n is an interior node in a CCT, it represents an invocation of a procedure. If n is a leaf in a CCT, it represents a statement inside some procedure. For leaves, their inclusive and exclusive costs are equal.

It is useful to perform scalability analysis for both inclusive and exclusive costs; if the loss of scalability attributed to the inclusive costs of a function invocation is roughly equal to the loss of scalability due to its exclusive costs, then we know that the computation in that function invocation does not scale. However, if the loss of scalability attributed



Figure 4.6: Computing the scaling loss when weak scaling a white dwarf detonation simulation with FLASH3 from 256 to 8192 cores. For weak scaling, the time on an MPI rank in each of the simulations will be the same. In the figure, column 0 represents the inclusive cost for one MPI rank in a 256-core simulation; column 2 represents the inclusive cost for one MPI rank in an 8192-core simulation. The difference between these two columns, computed as $\$2 - \0 , represents the excess work present in the larger simulation for each unique program context in the calling context tree. Dividing that by the total time in the 8192-core execution $@2$ gives the fraction of wasted time. Multiplying through by 100 gives the percent of the time wasted in the 8192-core execution, which corresponds to the % scalability loss.

to a function invocation's inclusive costs outweighs the loss of scalability accounted for by exclusive costs, we need to explore the scalability of the function's callees.

Given CCTs for an ensemble of executions, the next step to analyzing the scalability of their performance is to clearly define our expectations. Next, we describe performance expectations for weak scaling and intuitive metrics that represent how much performance deviates from our expectations. More information about our scalability analysis technique can be found elsewhere [4, 18].

Weak Scaling

Consider two weak scaling experiments executed on p and q processors, respectively, $p < q$. In Figure 4.6 shows how we can use a derived metric to compute and attribute scalability losses. Here, we compute the difference in inclusive cycles spent on one core of a 8192-core run and one core in a 256-core run in a weak scaling experiment. If the code had perfect weak scaling, the time for an MPI rank in each of the executions would be identical.

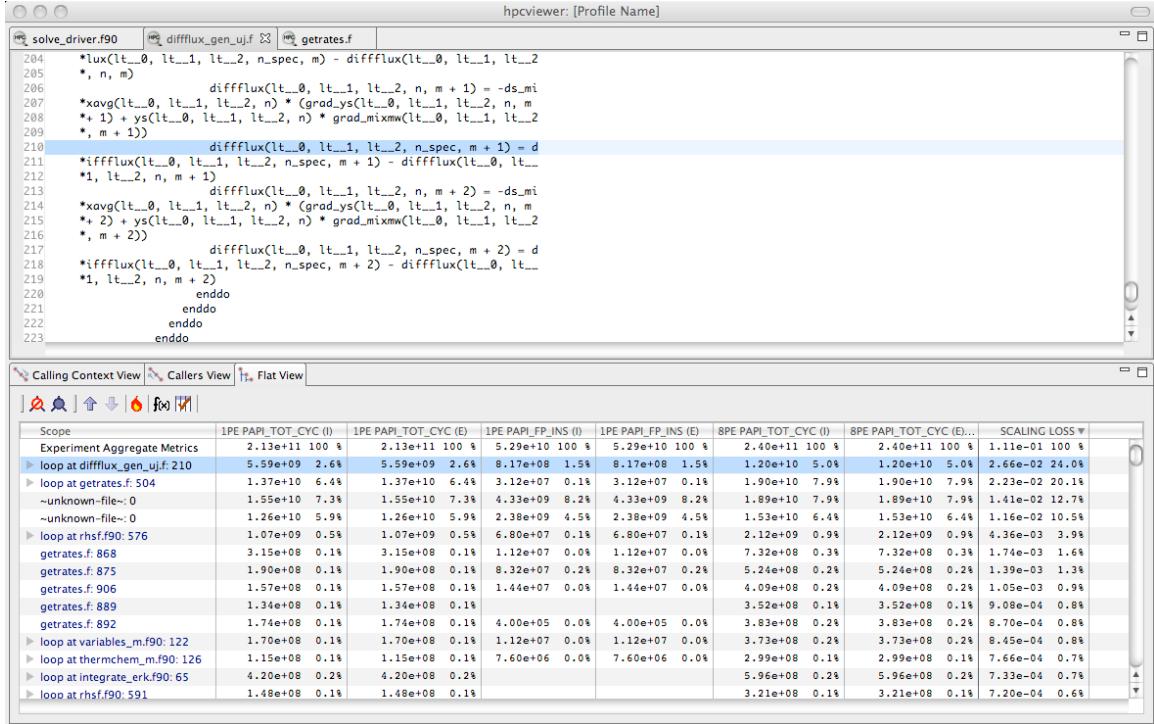


Figure 4.7: Using the fraction the scalability loss metric of Figure 4.6 to rank order loop nests by their scaling loss.

In this case, they are not. We compute the excess work by computing the difference for each scope between the time on the 8192-core run and the time on the 256-core core run. We normalize the differences of the time spent in the two runs by dividing then by the total time spent on the 8192-core run. This yields the fraction of wasted effort for each scope when scaling from 256 to 8192 cores. Finally, we multiply these results by 100 to compute the % scalability loss. This example shows how one can compute a derived metric to that pinpoints and quantifies scaling losses across different node counts of a Blue Gene/P system.

A similar analysis can be applied to compute scaling losses between jobs that use different numbers of core counts on individual processors. Figure 4.7 shows the result of computing the scaling loss for each loop nest when scaling from one to eight cores on a multicore node and rank order loop nests by their scaling loss metric. Here, we simply compute the scaling loss as the difference between the cycle counts of the eight-core and the one-core runs, divided through by the aggregate cost of the process executing on eight cores. This figure shows the scaling lost written in scientific notation as a fraction rather than multiplying through by 100 to yield a percent. In this figure, we examine scaling losses in the flat view, showing them for each loop nest. The source pane shows the loop nest responsible for the greatest scaling loss when scaling from one to eight cores. Unsurprisingly, the loop with the worst scaling loss is very memory intensive. Memory bandwidth is a precious commodity on multicore processors.

While we have shown how to compute and attribute the fraction of excess work in a weak scaling experiment, one can compute a similar quantity for experiments with strong scaling.

When differencing the costs summed across all of the threads in a pair of strong-scaling experiments, one uses exactly the same approach as shown in Figure 4.6. If comparing weak scaling costs summed across all ranks in p and q core executions, one can simply scale the aggregate costs by $1/p$ and $1/q$ respectively before differencing them.

Exploring Scaling Losses

Scaling losses can be explored in `hpcviewer` using any of its three views.

- *Top-down view.* This view represents the dynamic calling contexts (call paths) in which costs were incurred.
- *Bottom-up view.* This view enables one to look upward along call paths. This view is particularly useful for understanding the performance of software components or procedures that are used in more than one context, such as communication library routines.
- *Flat view.* This view organizes performance measurement data according to the static structure of an application. All costs incurred in *any* calling context by a procedure are aggregated together in the flat view.

`hpcviewer` enables developers to explore top-down, bottom-up, and flat views of CCTs annotated with costs, helping to quickly pinpoint performance bottlenecks. Typically, one begins analyzing an application’s scalability and performance using the top-down calling context tree view. Using this view, one can readily see how costs and scalability losses are associated with different calling contexts. If costs or scalability losses are associated with only a few calling contexts, then this view suffices for identifying the bottlenecks. When scalability losses are spread among many calling contexts, e.g., among different invocations of `MPI_Wait`, often it is useful to switch to the bottom-up of the data to see if many losses are due to the same underlying cause. In the bottom-up view, one can sort routines by their exclusive scalability losses and then look upward to see how these losses accumulate from the different calling contexts in which the routine was invoked.

Scaling loss based on excess work is intuitive; perfect scaling corresponds to a excess work value of 0, sublinear scaling yields positive values, and superlinear scaling yields negative values. Typically, CCTs for SPMD programs have similar structure. If CCTs for different executions diverge, using `hpcviewer` to compute and report excess work will highlight these program regions.

Inclusive excess work and exclusive excess work serve as useful measures of scalability associated with nodes in a calling context tree (CCT). By computing both metrics, one can determine whether the application scales well or not at a CCT node and also pinpoint the cause of any lack of scaling. If a node for a function in the CCT has comparable positive values for both inclusive excess work and exclusive excess work, then the loss of scaling is due to computation in the function itself. However, if the inclusive excess work for the function outweighs that accounted for by its exclusive costs, then one should explore the scalability of its callees. To isolate code that is an impediment to scalable performance, one can use the *hot path* button in `hpcviewer` to trace a path down through the CCT to see where the cost is incurred.

Chapter 5

Running Applications with `hpcrun` and `hpclink`

This chapter describes the mechanics of using `hpcrun` and `hpclink` to profile an application and collect performance data. For advice on how to choose events, perform scaling studies, etc., see Chapter 4 *Effective Strategies for Analyzing Program Performance*.

5.1 Using `hpcrun`

The `hpcrun` launch script is used to run an application and collect performance data for *dynamically linked* binaries. For dynamically linked programs, this requires no change to the program source and no change to the build procedure. You should build your application natively at full optimization. `hpcrun` inserts its profiling code into the application at runtime via `LD_PRELOAD`.

The basic options for `hpcrun` are `-e` (or `--event`) to specify a sampling source and rate and `-t` (or `--trace`) to turn on tracing. Sample sources are specified as ‘`event@howoften`’ where `event` is the name of the source and `howoften` is either a number specifying the period (threshold) for that event, or `f` followed by a number, e.g., `@f100` specifying a target sampling frequency for the event in samples/second.¹ Note that a higher period implies a lower rate of sampling. The `-e` option may be used multiple times to specify that multiple sample sources be used for measuring an execution. The basic syntax for profiling an application with `hpcrun` is:

```
hpcrun -t -e event@howoften ... app arg ...
```

For example, to profile an application using hardware counter sample sources provided by Linux `perf_events` and sample cycles at 300 times/second (the default sampling frequency) and sample every 4,000,000 instructions, you would use:

```
hpcrun -e CYCLES -e INSTRUCTIONS@4000000 app arg ...
```

¹Frequency-based sampling and the frequency-based notation for `howoften` is only available for sample sources managed by Linux `perf_events`. For Linux `perf_events`, HPCTOOLKIT uses a default sampling frequency of 300 samples/second.

The units for timer-based sample sources (`CPUTIME`, `REALTIME`, and `WALLCLOCK`) are microseconds, so to sample an application with tracing every 5,000 microseconds (200 times/second), you would use:

```
hpcrun -t -e CPUTIME@5000 app arg ...
```

`hpcrun` stores its raw performance data in a *measurements* directory with the program name in the directory name. On systems with a batch job scheduler (eg, PBS) the name of the job is appended to the directory name.

```
hpctoolkit-app-measurements[-jobid]
```

It is best to use a different measurements directory for each run. So, if you're using `hpcrun` on a local workstation without a job launcher, you can use the '`-o dirname`' option to specify an alternate directory name.

For programs that use their own launch script (eg, `mpirun` or `mpiexec` for MPI), put the application's run script on the outside (first) and `hpcrun` on the inside (second) on the command line. For example,

```
mpirun -n 4 hpcrun -e CYCLES mpiapp arg ...
```

Note that `hpcrun` is intended for profiling dynamically linked *binaries*. It will not work well if used to profile a shell script. At best, you would be profiling the shell interpreter, not the script commands, and sometimes this will fail outright.

It is possible to use `hpcrun` to launch a statically linked binary, but there are two problems with this. First, it is still necessary to build the binary with `hpclink`. Second, static binaries are commonly used on parallel clusters that require running the binary directly and do not accept a launch script. However, if your system allows it, and if the binary was produced with `hpclink`, then `hpcrun` will set the correct environment variables for profiling statically or dynamically linked binaries. All that `hpcrun` really does is set some environment variables (including `LD_PRELOAD`) and `exec` the binary.

5.2 Using `hpclink`

For now, see Chapter 10 on *Monitoring Statically Linked Applications*.

5.3 Harware Counter Event Names

HPCToolkit uses libpfm4 [9] to translate from an event name string to an event code recognized by the kernel. An event name is case insensitive and is defined as followed:

```
[pmu:] [event_name] [:unit_mask] [:modifier|:modifier=val]
```

- **pmu.** Optional name of the PMU (group of events) to which the event belongs to. This is useful to disambiguate events in case events from difference sources have the same name. If no pmu is specified, the first match event is used.

- **event_name**. The name of the event. It must be the complete name, partial matches are not accepted.
 - **unit_mask**. Some events can be refined using sub-events. A **unit_mask** designates an optional sub-event. An event may have multiple unit masks and it is possible to combine them (for some events) by repeating **:unit_mask** pattern.
 - **modifier**. A modifier is an optional filter that restricts when an event counts. The form of a modifier may be either **:modifier** or **:modifier=val**. For modifiers without a value, the presence of the modifier is interpreted as a restriction. Events may allow use of multiple modifiers at the same time.
 - **hardware event modifiers**. Some hardware events support one or more modifiers that restrict counting to a subset of events. For instance, on an Intel Broadwell EP, one can add a modifier to **MEM_LOAD_UOPS_RETIRED** to count only load operations that are an **L2_HIT** or an **L2_MISS**. For information about all modifiers for hardware events, one can direct HPCTOOLKIT’s measurement subsystem to list all native events and their modifiers as described in Section 5.4.
 - **precise_ip**. For some events, it is possible to control the amount of skid. Skid is a measure of how many instructions may execute between an event and the PC where the event is reported. Smaller skid enables more accurate attribution of events to instructions. Without a skid modifier, **hpcrun** allows arbitrary skid because some architectures don’t support anything more precise. One may optionally specify one of the following as a skid modifier:
 - * **:p** : a sample must have constant skid.
 - * **:pp** : a sample is requested to have 0 skid.
 - * **:ppp** : a sample must have 0 skid.
 - * **:P** : autodetect the least skid possible.
- NOTE: If the kernel or the hardware does not support the specified value of the skid, no error message will be reported but no samples will be recorded.

5.4 Sample Sources

This section provides an overview of how to use sample sources supported by HPCToolkit. To see a list of the available sample sources and events that **hpcrun** supports, use ‘**hpcrun -L**’ (dynamic) or set ‘**HPCRUN_EVENT_LIST=LIST**’ (static). Note that on systems with separate compute nodes, it is best to run this on a compute node.

5.4.1 Linux perf_events

Linux **perf_events** provides a powerful interface that supports measurement of both application execution and kernel activity. Using **perf_events**, one can measure both hardware and software events. Using a processor’s hardware performance monitoring unit (PMU), the **perf_events** interface can measure an execution using any hardware counter supported by the PMU. Examples of hardware events include cycles, instructions completed, cache misses,

and stall cycles. Using instrumentation built in to the Linux kernel, the `perf_events` interface can measure software events. Examples of software events include page faults, context switches, and CPU migrations.

Capabilities of HPCToolkit’s `perf_events` Interface

Frequency-based sampling. The Linux `perf_events` interface supports frequency-based sampling. With frequency-based sampling, the kernel automatically selects and adjusts an event period with the aim of delivering samples for that event at a target sampling frequency.² Unless a user explicitly specifies an event count threshold for an event, HPCToolkit’s measurement interface will use frequency-based sampling by default. HPCToolkit’s default sampling frequency is $\min(300, M - 1)$, where M is the value specified in the system configuration file `/proc/sys/kernel/perf_event_max_sample_rate`.

For circumstances where the user wants to use frequency-based sampling but HPCToolkit’s default sampling frequency is inappropriate, one can specify the target sampling frequency for a particular event using the notation `event@f` when specifying an event or change the default sampling frequency. When measuring a dynamically-linked executable using `hpcrun`, one can change the default sampling frequency using `hpcrun`’s `-c` option. To set a new default sampling frequency for a statically-linked executable instrumented with `hpclink`, set the `HPCRUN_PERF_COUNT` environment variable. The section below entitled *Launching* provides examples of how to monitor an execution using frequency-based sampling.

Multiplexing. Using multiplexing enables one to monitor more events in a single execution than the number of hardware counters a processor can support for each thread. The number of events that can be monitored in a single execution is only limited by the maximum number of concurrent events that the kernel will allow a user to multiplex using the `perf_events` interface.

When more events are specified than can be monitored simultaneously using a thread’s hardware counters,³ the kernel will employ multiplexing and divide the set of events to be monitored into groups, monitor only one group of events at a time, and cycle repeatedly through the groups as a program executes.

For applications that have very regular, steady state behavior, e.g., an iterative code with lots of iterations, multiplexing will yield results that are suitably representative of execution behavior. However, for executions that consist of unique short phases, measurements collected using multiplexing may not accurately represent the execution behavior. To obtain more accurate measurements, one can run an application multiple times and in each run collect a subset of events that can be measured without multiplexing. Results from several such executions can be imported into HPCToolkit’s `hpcviewer` and analyzed together.

²The kernel may be unable to deliver the desired frequency if there are fewer events per second than the desired frequency.

³How many events can be monitored simultaneously on a particular processor may depend on the events specified.

Thread blocking. When a program executes, a thread may block waiting for the kernel to complete some operation on its behalf. For instance, a thread may block waiting for data to become available so that a `read` operation can complete. On systems running Linux 4.3 or newer, one can use the `perf_events` sample source to monitor how much time a thread is blocked and where the blocking occurs. To measure the time a thread spends blocked, one can profile with `BLOCKTIME` event and another time-based event, such as `CYCLES`. The `BLOCKTIME` event shouldn't have any frequency or period specified, whereas `CYCLES` may have a frequency or period specified.

Launching

When sampling with native events, by default `hpcrun` will profile using `perf_events`. To force HPCToolkit to use PAPI rather than `perf_events` to oversee monitoring of a PMU event (assuming that HPCToolkit has been configured to include support for PAPI), one must prefix the event with ‘`papi:::`’ as follows:

```
hpcrun -e papi::CYCLES
```

For PAPI presets, there is no need to prefix the event with ‘`papi:::`’. For instance it is sufficient to specify `PAPI_TOT_CYC` event without any prefix to profile using PAPI. For more information about using PAPI, see Section 5.4.2.

Below, we provide some examples of various ways to measure `CYCLES` and `INSTRUCTIONS` using HPCTOOLKIT’s `perf_events` measurement substrate:

To sample an execution 100 times per second (frequency-based sampling) counting `CYCLES` and 100 times a second counting `INSTRUCTIONS`:

```
hpcrun -e CYCLES@f100 -e INSTRUCTIONS@f100 ...
```

To sample an execution every 1,000,000 cycles and every 1,000,000 instructions using period-based sampling:

```
hpcrun -e CYCLES@1000000 -e INSTRUCTIONS@1000000
```

By default, `hpcrun` uses frequency-based sampling with the rate 300 samples per second per event type. Hence the following command causes HPCTOOLKIT to sample `CYCLES` at 300 samples per second and `INSTRUCTIONS` at 300 samples per second:

```
hpcrun -e CYCLES -e INSTRUCTIONS ...
```

One can specify a different default sampling period or frequency using the `-c` option. The command below will sample `CYCLES` and `INSTRUCTIONS` at 200 samples per second each:

```
hpcrun -c f200 -e CYCLES -e INSTRUCTIONS ...
```

Notes

- Linux `perf_events` uses one file descriptor for each event to be monitored. Furthermore, since `hpcrun` generates one `hpcrun` file for each thread, and an additional `hpctrace` file if traces is enabled. Hence for e events and t threads, the required number of file descriptors is:

$$t \times e + t + t \text{ (if trace is enabled)}$$

For instance, if one profiles a multi-threaded program that executes with 500 threads using 4 events, then the required number of file descriptors is

$$\begin{aligned} & 500 \text{ threads} \times 4 \text{ events} + 500 \text{ hpcrun files} + 500 \text{ hpctrace files} \\ & = 3000 \text{ file descriptors} \end{aligned}$$

If the number of file descriptors exceeds the number of maximum number of open files, then the program will crash. To remedy this issue, one needs to increase the number of maximum number of open files allowed.

- When a system is configured with suitable permissions, HPC Toolkit will sample call stacks within the Linux kernel in addition to application-level call stacks. This feature can be useful to measure kernel activity on behalf of a thread (e.g., zero-filling allocated pages when they are first touched) or to observe where, why, and how long a thread blocks. For a user to be able to sample kernel call stacks, the configuration file `/proc/sys/kernel/perf_event_paranoid` must have a value ≤ 1 . To associate addresses in kernel call paths with function names, the value of `/proc/sys/kernel/kptr_restrict` must be 0 (number zero). If these settings are not configured in this way on your system, you will need someone with administrator privileges to change them for you to be able to sample call stacks within the kernel.
- Due to a limitation present in all Linux kernel versions currently available, HPC Toolkit's measurement subsystem can only approximate a thread's blocking time. At present, Linux reports when a thread blocks but does not report when a thread resumes execution. For that reason, HPC Toolkit's measurement subsystem approximates the time a thread spends blocked using sampling as the time between when the thread blocks and when the thread receives its first sample after resuming execution.
- Users need to be cautious when considering measured counts of events that have been collected using hardware counter multiplexing. Currently, it is not obvious to a user if a metric was measured using a multiplexed counter. Information about whether a counter was multiplexed is only available in the `experiment.xml` file produced when post-processing measurement data with `hpcprof` or `hpcprof-mpi`, but is not visible in `hpcviewer`.

5.4.2 PAPI

PAPI, the Performance API, is a library for providing access to the hardware performance counters. PAPI aims to provide a consistent, high-level interface that consists of a universal set of event names that can be used to measure performance on any processor,

PAPI_BR_INS	Branch instructions
PAPI_BR_MSP	Conditional branch instructions mispredicted
PAPI_FP_INS	Floating point instructions
PAPI_FP_OPS	Floating point operations
PAPI_L1_DCA	Level 1 data cache accesses
PAPI_L1_DCM	Level 1 data cache misses
PAPI_L1_ICH	Level 1 instruction cache hits
PAPI_L1_ICM	Level 1 instruction cache misses
PAPI_L2_DCA	Level 2 data cache accesses
PAPI_L2_ICM	Level 2 instruction cache misses
PAPI_L2_TCM	Level 2 cache misses
PAPI_LD_INS	Load instructions
PAPI_SR_INS	Store instructions
PAPI_TLB_DM	Data translation lookaside buffer misses
PAPI_TOT_CYC	Total cycles
PAPI_TOT_IIS	Instructions issued
PAPI_TOT_INS	Instructions completed

Table 5.1: Some commonly available PAPI events. The exact set of available events is system dependent.

independent of any processor-specific event names. In some cases, PAPI event names represent quantities synthesized by combining measurements based on multiple native events available on a particular processor. For instance, in some cases PAPI reports total cache misses by measuring and combining data misses and instruction misses. PAPI is available from the University of Tennessee at <http://icl.cs.utk.edu/papi>.

PAPI focuses mostly on in-core CPU events: cycles, cache misses, floating point operations, mispredicted branches, etc. For example, the following command samples total cycles and L2 cache misses.

```
hpcrun -e PAPI_TOT_CYC@15000000 -e PAPI_L2_TCM@400000 app arg ...
```

The precise set of PAPI preset and native events is highly system dependent. Commonly, there are events for machine cycles, cache misses, floating point operations and other more system specific events. However, there are restrictions both on how many events can be sampled at one time and on what events may be sampled together and both restrictions are system dependent. Table 5.1 contains a list of commonly available PAPI events.

To see what PAPI events are available on your system, use the `papi_avail` command from the `bin` directory in your PAPI installation. The event must be both available and not derived to be usable for sampling. The command `papi_native_avail` displays the machine's native events. Note that on systems with separate compute nodes, you normally need to run `papi_avail` on one of the compute nodes.

When selecting the period for PAPI events, aim for a rate of approximately a few hundred samples per second. So, roughly several million or tens of million for total cycles

or a few hundred thousand for cache misses. PAPI and `hpcrun` will tolerate sampling rates as high as 1,000 or even 10,000 samples per second (or more). However, rates higher than a few hundred samples per second will only increase measurement overhead and distort the execution of your program; they won't yield more accurate results.

Beginning with Linux kernel version 2.6.32, support for accessing performance counters using the Linux `perf_events` performance monitoring subsystem is built into the kernel. `perf_events` provides a measurement substrate for PAPI on Linux.

On Blue Gene platforms that are not based on Linux, PAPI is highly recommended as it provides an essential substrate for accessing hardware performance counters. On modern Linux systems that include support for `perf_events`, PAPI is only recommended for monitoring events outside the scope of the `perf_events` interface.

Proxy Sampling HPCTOOLKIT supports proxy sampling for derived PAPI events. For HPCTOOLKIT to sample a PAPI event directly, the event must not be derived and must trigger hardware interrupts when a threshold is exceeded. For events that cannot trigger interrupts directly, HPCToolkit's proxy sampling sample on another event that is supported directly and then reads the counter for the derived event. In this case, a native event can serve as a proxy for one or more derived events.

To use proxy sampling, specify the `hpcrun` command line as usual and be sure to include at least one non-derived PAPI event. The derived events will be accumulated automatically when processing a sample trigger for a native event. We recommend adding `PAPI_TOT_CYC` as a native event when using proxy sampling, but proxy sampling will gather data as long as the event set contains at least one non-derived PAPI event. Proxy sampling requires one non-derived PAPI event to serve as the proxy; a Linux timer can't serve as the proxy for a PAPI derived event.

For example, on newer Intel CPUs, often PAPI floating point events are all derived and cannot be sampled directly. In that case, you could count FLOPs by using cycles a proxy event with a command line such as the following. The period for derived events is ignored and may be omitted.

```
hpcrun -e PAPI_TOT_CYC@6000000 -e PAPI_FP_OPS app arg ...
```

Attribution of proxy samples is not as accurate as regular samples. The problem, of course, is that the event that triggered the sample may not be related to the derived counter. The total count of events should be accurate, but their location at the leaves in the Calling Context tree may not be very accurate. However, the higher up the CCT, the more accurate the attribution becomes. For example, suppose you profile a loop of mixed integer and floating point operations and sample on `PAPI_TOT_CYC` directly and count `PAPI_FP_OPS` via proxy sampling. The attribution of flops to individual statements within the loop is likely to be off. But as long as the loop is long enough, the count for the loop as a whole (and up the tree) should be accurate.

5.4.3 WALLCLOCK, REALTIME and CPUTIME

HPCTOOLKIT supports three timer-based sample sources: CPUTIME, REALTIME and WALLCLOCK. The unit for periods of these timers is microseconds.

Before describing this capability further, it is worth noting that the CYCLES event supported by Linux `perf_events` or PAPI’s `PAPI_TOT_CYC` are generally superior to any of the timer-based sampling sources.

The `CPUTIME` and `REALTIME` sample sources are based on the POSIX timers `CLOCK_THREAD_CPUTIME_ID` and `CLOCK_REALTIME` with the Linux `SIGEV_THREAD_ID` extension. `CPUTIME` only counts time when the CPU is running; `REALTIME` counts real (wall clock) time, whether the process is running or not. Signal delivery for these timers is thread-specific, so these timers are suitable for profiling multithreaded programs. Sampling using the `REALTIME` sample source may break some applications that don’t handle interrupted syscalls well. In that case, consider using `CPUTIME` instead. It is worth noting that `REALTIME` and `CPUTIME` are not available on Blue Gene, where compute nodes run a custom microkernel instead of Linux.

The `WALLCLOCK` sample source is based on the `ITIMER_PROF` interval timer. `WALLCLOCK` counts time when an application is running or the kernel is running on behalf of the application. Unlike the POSIX timers that support `CPUTIME` and `REALTIME`, the interval timer supporting `WALLCLOCK` does not support thread-specific signal delivery on Linux. `WALLCLOCK` is available on Blue Gene even though `REALTIME` and `CPUTIME` are not. On Linux the maximum interrupt rate for the interval timer is limited by the system’s Hz rate, commonly 1,000 cycles per second, but may be lower. That is, `WALLCLOCK@10` will not generate any higher sampling rate than `WALLCLOCK@1000`. On IBM Blue Gene, the interval timer is not bound by the Hz rate and so sampling rates faster than 1,000 per second are possible.

The following example, which specifies a period of 5000 microseconds will sample each thread in `app` at a rate of approximately 200 times per second.

```
hpcrun -e REALTIME@5000 app arg ...
```

Note: do not use more than one timer-based sample source to monitor a program execution. When using a sample source such as `CPUTIME`, `REALTIME`, or `WALLCLOCK`, we recommend not using another time-based sampling source such as Linux `perf_events` `CYCLES` or PAPI’s `PAPI_TOT_CYC`. Technically, this is feasible and `hpcrun` won’t die. However, multiple time-based sample sources would compete with one another to measure the execution and likely lead to dropped samples and possibly distorted results.

5.4.4 IO

The `IO` sample source counts the number of bytes read and written. This displays two metrics in the viewer: “IO Bytes Read” and “IO Bytes Written.” The `IO` source is a synchronous sample source. It overrides the functions `read`, `write`, `fread` and `fwrite` and records the number of bytes read or written along with their dynamic context synchronously rather than relying on data collection triggered by interrupts.

To include this source, use the `IO` event (no period). In the static case, two steps are needed. Use the `--io` option for `hpmlink` to link in the `IO` library and use the `IO` event to activate the `IO` source at runtime. For example,

```
(dynamic) hpcrun -e IO app arg ...
(static) hpmlink --io gcc -g -O -static -o app file.c ...
           export HPCRUN_EVENT_LIST=IO
           app arg ...
```

The `IO` source is mainly used to find where your program reads or writes large amounts of data. However, it is also useful for tracing a program that spends much time in `read` and `write`. The hardware performance counters do not advance while running in the kernel, so the trace viewer may misrepresent the amount of time spent in syscalls such as `read` and `write`. By adding the `IO` source, `hpcrun` overrides `read` and `write` and thus is able to more accurately count the time spent in these functions.

5.4.5 MEMLEAK

The `MEMLEAK` sample source counts the number of bytes allocated and freed. Like `IO`, `MEMLEAK` is a synchronous sample source and does not generate asynchronous interrupts. Instead, it overrides the malloc family of functions (`malloc`, `calloc`, `realloc` and `free` plus `memalign`, `posix_memalign` and `valloc`) and records the number of bytes allocated and freed along with their dynamic context.

`MEMLEAK` allows you to find locations in your program that allocate memory that is never freed. But note that failure to free a memory location does not necessarily imply that location has leaked (missing a pointer to the memory). It is common for programs to allocate memory that is used throughout the lifetime of the process and not explicitly free it.

To include this source, use the `MEMLEAK` event (no period). Again, two steps are needed in the static case. Use the `--memleak` option for `hpclink` to link in the `MEMLEAK` library and use the `MEMLEAK` event to activate it at runtime. For example,

```
(dynamic) hpcrun -e MEMLEAK app arg ...
(static) hpclink --memleak gcc -g -O -static -o app file.c ...
          export HPCRUN_EVENT_LIST=MEMLEAK
          app arg ...
```

If a program allocates and frees many small regions, the `MEMLEAK` source may result in a high overhead. In this case, you may reduce the overhead by using the `memleak` probability option to record only a fraction of the mallocs. For example, to monitor 10% of the mallocs, use:

```
(dynamic) hpcrun -e MEMLEAK --memleak-prob 0.10 app arg ...
(static) export HPCRUN_EVENT_LIST=MEMLEAK
          export HPCRUN_MEMLEAK_PROB=0.10
          app arg ...
```

It might appear that if you monitor only 10% of the program's mallocs, then you would have only a 10% chance of finding the leak. But if a program leaks memory, then it's likely that it does so many times, all from the same source location. And you only have to find that location once. So, this option can be a useful tool if the overhead of recording all mallocs is prohibitive.

Rarely, for some programs with complicated memory usage patterns, the `MEMLEAK` source can interfere with the application's memory allocation causing the program to segfault. If this happens, use the `hpcrun` debug (dd) variable `MEMLEAK_NO_HEADER` as a workaround.

```
(dynamic) hpcrun -e MEMLEAK -dd MEMLEAK_NO_HEADER app arg ...
(static) export HPCRUN_EVENT_LIST=MEMLEAK
          export HPCRUN_DEBUG_FLAGS=MEMLEAK_NO_HEADER
          app arg ...
```

The `MEMLEAK` source works by attaching a header or a footer to the application's `malloc`'d regions. Headers are faster but have a greater potential for interfering with an application. Footers have higher overhead (require an external lookup) but have almost no chance of interfering with an application. The `MEMLEAK_NO_HEADER` variable disables headers and uses only footers.

5.5 Process Fraction

Although `hpcrun` can profile parallel jobs with thousands or tens of thousands of processes, there are two scaling problems that become prohibitive beyond a few thousand cores. First, `hpcrun` writes the measurement data for all of the processes into a single directory. This results in one file per process plus one file per thread (two files per thread if using tracing). Unix file systems are not equipped to handle directories with many tens or hundreds of thousands of files. Second, the sheer volume of data can overwhelm the viewer when the size of the database far exceeds the amount of memory on the machine.

The solution is to sample only a fraction of the processes. That is, you can run an application on many thousands of cores but record data for only a few hundred processes. The other processes run the application but do not record any measurement data. This is what the process fraction option (`-f` or `--process-fraction`) does. For example, to monitor 10% of the processes, use:

```
(dynamic) hpcrun -f 0.10 -e event@howoften app arg ...
(dynamic) hpcrun -f 1/10 -e event@howoften app arg ...
(static)  export HPCRUN_EVENT_LIST='event@howoften'
          export HPCRUN_PROCESS_FRACTION=0.10
          app arg ...
```

With this option, each process generates a random number and records its measurement data with the given probability. The process fraction (probability) may be written as a decimal number (0.10) or as a fraction (1/10) between 0 and 1. So, in the above example, all three cases would record data for approximately 10% of the processes. Aim for a number of processes in the hundreds.

5.6 Starting and Stopping Sampling

HPCTOOLKIT supports an API for the application to start and stop sampling. This is useful if you want to profile only a subset of a program and ignore the rest. The API supports the following functions.

```
void hpctoolkit_sampling_start(void);
void hpctoolkit_sampling_stop(void);
```

For example, suppose that your program has three major phases: it reads input from a file, performs some numerical computation on the data and then writes the output to another file. And suppose that you want to profile only the compute phase and skip the read and write phases. In that case, you could stop sampling at the beginning of the program, restart it before the compute phase and stop it again at the end of the compute phase.

This interface is process wide, not thread specific. That is, it affects all threads of a process. Note that when you turn sampling on or off, you should do so uniformly across all processes, normally at the same point in the program. Enabling sampling in only a subset of the processes would likely produce skewed and misleading results.

And for technical reasons, when sampling is turned off in a threaded process, interrupts are disabled only for the current thread. Other threads continue to receive interrupts, but they don't unwind the call stack or record samples. So, another use for this interface is to protect syscalls that are sensitive to being interrupted with signals. For example, some Gemini interconnect (GNI) functions called from inside `gasnet_init()` or `MPI_Init()` on Cray XE systems will fail if they are interrupted by a signal. As a workaround, you could turn sampling off around those functions.

Also, you should use this interface only at the top level for major phases of your program. That is, the granularity of turning sampling on and off should be much larger than the time between samples. Turning sampling on and off down inside an inner loop will likely produce skewed and misleading results.

To use this interface, put the above function calls into your program where you want sampling to start and stop. Remember, starting and stopping apply process wide. For C/C++, include the following header file from the HPCTOOLKIT `include` directory.

```
#include <hpctoolkit.h>
```

Compile your application with `libhpctoolkit` with `-I` and `-L` options for the include and library paths. For example,

```
gcc -I /path/to/hpctoolkit/include app.c ... \
    -L /path/to/hpctoolkit/lib/hpctoolkit -lhpctoolkit ...
```

The `libhpctoolkit` library provides weak symbol no-op definitions for the start and stop functions. For dynamically linked programs, be sure to include `-lhpctoolkit` on the link line (otherwise your program won't link). For statically linked programs, `hpclink` adds strong symbol definitions for these functions. So, `-lhpctoolkit` is not necessary in the static case, but it doesn't hurt.

To run the program, set the `LD_LIBRARY_PATH` environment variable to include the HPCTOOLKIT `lib/hpctoolkit` directory. This step is only needed for dynamically linked programs.

```
export LD_LIBRARY_PATH=/path/to/hpctoolkit/lib/hpctoolkit
```

Note that sampling is initially turned on until the program turns it off. If you want it initially turned off, then use the `-ds` (or `--delay-sampling`) option for `hpcrun` (dynamic) or set the `HPCRUN_DELAY_SAMPLING` environment variable (static).

```
(dynamic) hpcrun -ds -e event@howoften app arg ...
(static) export HPCRUN_EVENT_LIST='event@howoften'
          export HPCRUN_DELAY_SAMPLING=1
          app arg ...
```

5.7 Environment Variables for `hpcrun`

For most systems, `hpcrun` requires no special environment variable settings. There are two situations, however, where `hpcrun`, to function correctly, *must* refer to environment variables. These environment variables, and corresponding situations are:

HPCTOOLKIT To function correctly, `hpcrun` must know the location of the HPCTOOLKIT top-level installation directory. The `hpcrun` script uses elements of the installation `lib` and `libexec` subdirectories. On most systems, the `hpcrun` can find the requisite components relative to its own location in the file system. However, some parallel job launchers *copy* the `hpcrun` script to a different location as they launch a job. If your system does this, you must set the `HPCTOOLKIT` environment variable to the location of the HPCTOOLKIT top-level installation directory before launching a job.

Note to system administrators: if your system provides a module system for configuring software packages, then constructing a module for HPCTOOLKIT to initialize these environment variables to appropriate settings would be convenient for users.

5.8 Platform-Specific Notes

5.8.1 Cray Systems

If you are trying to profile a dynamically-linked executable on a Cray that is still using the ALPS job launcher and you see an error like the following

```
/var/spool/alps/103526/hpcrun: Unable to find HPCTOOLKIT root directory.
Please set HPCTOOLKIT to the install prefix, either in this script,
or in your environment, and try again.
```

in your job's error log then read on. Otherwise, skip this section.

The problem is that the Cray job launcher copies HPCToolkit's `hpcrun` script to a directory somewhere below `/var/spool/alps/` and runs it from there. By moving `hpcrun` to a different directory, this breaks `hpcrun`'s method for finding HPCTOOLKIT's install directory.

To fix this problem, in your job script, set `HPCTOOLKIT` to the top-level HPCTOOLKIT installation directory (the directory containing the `bin`, `lib` and `libexec` subdirectories) and export it to the environment. (If launching statically-linked binaries created using `hpmlink`, this step is unnecessary, but harmless.) Figure 5.1 show a skeletal job script that sets the `HPCTOOLKIT` environment variable before monitoring a dynamically-linked executable with `hpcrun`:

Your system may have a module installed for `hpctoolkit` with the correct settings for `PATH`, `HPCTOOLKIT`, etc. In that case, the easiest solution is to load the `hpctoolkit` module. Try “`module show hpctoolkit`” to see if it sets `HPCTOOLKIT`.

```

#!/bin/sh
#PBS -l mppwidth=#nodes
#PBS -l walltime=00:30:00
#PBS -V

export HPCTOOLKIT=/path/to/hpctoolkit/install/directory
export CRAY_ROOTFS=DSL

cd $PBS_O_WORKDIR
aprun -n #nodes hpcrun -e event@howoften dynamic-app arg ...

```

Figure 5.1: A sketch of how to help HPCToolkit find its dynamic libraries when using Cray’s ALPS job launcher.

5.8.2 Blue Gene/Q Systems

Blue Gene Q systems provide the `WALLCLOCK` interval timer, but not the POSIX `CPUTIME` and `REALTIME` timers.

The Linux `perf_events` subsystem is unavailable on Blue Gene Q systems. One should use the PAPI interface to monitor executions using hardware performance counters.

5.8.3 ARM Systems

HPCTOOLKIT’s measurement infrastructure depends upon `libunwind` for call stack unwinding on ARM.

Chapter 6

Measurement and Analysis of GPU-accelerated Applications

Note: This release contains beta support for monitoring computations offloaded onto NVIDIA GPUs and alpha support for monitoring computations offloaded onto AMD GPUs.

HPCToolkit can measure both the CPU and GPU performance of GPU-accelerated applications. It can measure CPU performance using asynchronous sampling triggered by Linux timers or hardware counter events as described in Section 5.4 and it can monitor GPU performance using tool support libraries provided by GPU vendors.

The foundation of HPCToolkit’s support for measuring the performance of GPU-accelerated applications is a vendor-independent monitoring substrate. A thin software layer connects NVIDIA’s CUPTI (CUDA Performance Tools Interface) [12] and AMD’s ROC-tracer (ROCM Tracer Callback/Activity Library) [3] monitoring libraries to this substrate. HPCToolkit reports GPU performance metrics in a vendor-neutral way. For instance, rather than focusing on NVIDIA warps or AMD wavefronts, HPCToolkit presents both as fine-grain, thread-level parallelism.

In the following sections, we describe how to use HPCToolkit to measure GPU metrics for GPU-accelerated applications. We begin with a discussion of NVIDIA GPUs since HPCToolkit’s support for them is more complete than for others.

6.1 NVIDIA GPUs

HPCToolkit supports two levels of performance monitoring for NVIDIA GPUs: coarse-grain profiling and tracing of GPU activities at the operation level (e.g., kernel launches, data allocations, memory copies, ...) , and fine-grain profiling of GPU computations using PC sampling, which measures GPU computations at a granularity of individual machine instructions. Section 6.1.1 describes fine-grain GPU performance measurement using PC sampling and the metrics it measures or computes.

While performing coarse-grain GPU monitoring of kernel executions, memory copies, and other GPU activities, HPCToolkit will collect a trace of activity for each GPU stream if tracing is enabled. Table 6.1 shows the possible command-line arguments to `hpcrun` that

Argument to <code>hpcrun</code>	What is monitored
<code>-e gpu=nvidia</code>	profiling of GPU operations
<code>-e gpu=nvidia -t</code>	profiling and tracing of GPU operations
<code>-e gpu=nvidia,pc</code>	PC sampling of GPU computations in addition to profiling of GPU operations
<code>-e gpu=nvidia,pc -t</code>	PC sampling of GPU computations in addition to profiling and tracing of GPU operations

Table 6.1: Monitoring performance on NVIDIA GPUs.

Metric	Description
GKER (sec)	GPU time: kernel execution (seconds)
GMEM (sec)	GPU time: memory allocation/deallocation (seconds)
GMSET (sec)	GPU time: memory set (seconds)
GXCOPY (sec)	GPU time: explicit data copy (seconds)
GSYNC (sec)	GPU time: synchronization (seconds)

Table 6.2: GPU operation timings.

will enable different levels of monitoring for NVIDIA GPUs. When fine-grain monitoring using PC sampling is enabled, coarse-grain profiling is also performed, so tracing is available in this mode as well. However, since PC sampling dilates the CPU overhead of GPU-accelerated codes, tracing is not recommended when PC sampling is enabled.

Coarse-grain profiling attributes to each calling context the total time of all GPU operations initiated in that context. Table 6.2 shows the classes of GPU operations for which timings are collected. In addition, HPCToolkit records metrics for operations performed including memory allocation and deallocation (Table 6.3), memory set (Table 6.4), explicit memory copies (Table 6.5), and synchronization (Table 6.6). In addition, HPCToolkit reports GPU kernel characteristics, including including register usage, thread count per block, and theoretical occupancy as shown in Table 6.7. NVIDIA defines theoretical occupancy as the ratio of the active threads in a streaming multiprocessor to the maximum active threads supported by the hardware in one streaming multiprocessor.

At present, using NVIDIA’s CUPTI library adds substantial measurement overhead. Unlike CPU monitoring based on asynchronous sampling, GPU performance monitoring uses vendor-provided callback interfaces to intercept the initiation of each GPU operation. Accordingly, the overhead of GPU performance monitoring depends upon how frequently GPU operations are initiated. In our experience to date, profiling (and if requested, tracing) on NVIDIA GPUs using NVIDIA’s CUPTI interface roughly doubles the execution time of a GPU-accelerated application. In our experience, we have seen NVIDIA’s PC sampling dilate the execution time of a GPU-accelerated program by 30× using CUDA 10 or earlier. Our early experience with CUDA 11 indicates that overhead using PC sampling is much lower and less than 5×. The overhead of GPU monitoring is principally on the host side. As measured by CUPTI, the time spent in GPU operations or PC samples is expected to be relatively accurate. However, since execution as a whole is slowed while measuring GPU performance, when evaluating GPU activity reported by HPCToolkit, one must be careful.

Metric	Description
GMEM:UNK (B)	GPU memory alloc/free: unknown memory kind (bytes)
GMEM:PAG (B)	GPU memory alloc/free: pageable memory (bytes)
GMEM:PIN (B)	GPU memory alloc/free: pinned memory (bytes)
GMEM:DEV (B)	GPU memory alloc/free: device memory (bytes)
GMEM:ARY (B)	GPU memory alloc/free: array memory (bytes)
GMEM:MAN (B)	GPU memory alloc/free: managed memory (bytes)
GMEM:DST (B)	GPU memory alloc/free: device static memory (bytes)
GMEM:MST (B)	GPU memory alloc/free: managed static memory (bytes)
GMEM:COUNT	GPU memory alloc/free: count

Table 6.3: GPU memory allocation and deallocation.

Metric	Description
GMSET:UNK (B)	GPU memory set: unknown memory kind (bytes)
GMSET:PAG (B)	GPU memory set: pageable memory (bytes)
GMSET:PIN (B)	GPU memory set: pinned memory (bytes)
GMSET:DEV (B)	GPU memory set: device memory (bytes)
GMSET:ARY (B)	GPU memory set: array memory (bytes)
GMSET:MAN (B)	GPU memory set: managed memory (bytes)
GMSET:DST (B)	GPU memory set: device static memory (bytes)
GMSET:MST (B)	GPU memory set: managed static memory (bytes)
GMSET:COUNT	GPU memory set: count

Table 6.4: GPU memory set metrics.

For instance, if a GPU-accelerated program runs in 1000 seconds without HPC Toolkit monitoring GPU activity but slows to 2000 seconds when GPU profiling and tracing is enabled, then if GPU profiles and traces show that the GPU is active for 25% of the execution time, one should re-scale the accurate measurements of GPU activity by considering the $2 \times$ dilation when monitoring GPU activity. Without monitoring, one would expect the same level of GPU activity, but the host time would be twice as fast. Thus, without monitoring, the ratio of GPU activity to host activity would be roughly double.

6.1.1 PC Sampling on NVIDIA GPUs

NVIDIA’s GPUs have supported PC sampling since Maxwell [5]. Instruction samples are collected separately on each active streaming multiprocessor (SM) and merged in a buffer returned by NVIDIA’s CUPTI. In each sampling period, one warp scheduler of each active SM samples the next instruction from one of its active warps. Sampling rotates through an SM’s warp schedulers in a round robin fashion. When an instruction is sampled, its stall reason (if any) is recorded. If all warps on a scheduler are stalled when a sample is taken, the sample is marked as a latency sample, meaning no instruction will be issued by the warp scheduler in the next cycle. Figure 6.1 shows a PC sampling example on an SM with four

Metric	Description
GXCOPY:UNK (B)	GPU explicit memory copy: unknown kind (bytes)
GXCOPY:H2D (B)	GPU explicit memory copy: host to device (bytes)
GXCOPY:D2H (B)	GPU explicit memory copy: device to host (bytes)
GXCOPY:H2A (B)	GPU explicit memory copy: host to array (bytes)
GXCOPY:A2H (B)	GPU explicit memory copy: array to host (bytes)
GXCOPY:A2A (B)	GPU explicit memory copy: array to array (bytes)
GXCOPY:A2D (B)	GPU explicit memory copy: array to device (bytes)
GXCOPY:D2A (B)	GPU explicit memory copy: device to array (bytes)
GXCOPY:D2D (B)	GPU explicit memory copy: device to device (bytes)
GXCOPY:H2H (B)	GPU explicit memory copy: host to host (bytes)
GXCOPY:P2P (B)	GPU explicit memory copy: peer to peer (bytes)
GXCOPY:COUNT	GPU explicit memory copy: count

Table 6.5: GPU explicit memory copy metrics.

Metric	Description
GSYNC:UNK (sec)	GPU synchronizations: unknown kind
GSYNC:EVT (sec)	GPU synchronizations: event
GSYNC:STRE (sec)	GPU synchronizations: stream event wait
GSYNC:STR (sec)	GPU synchronizations: stream
GSYNC:CTX (sec)	GPU synchronizations: context
GSYNC:COUNT	GPU synchronizations: count

Table 6.6: GPU synchronization metrics.

schedulers. Among the six collected samples, four are latency samples, so the estimated stall ratio is 4/6.

Figure 6.8 shows the stall metrics recorded by HPCToolkit using CUPTI’s PC sampling. Figure 6.9 shows PC sampling summary statistics recorded by HPCToolkit. Of particular note is the metric **GSAMP:UTIL**. HPCToolkit computes approximate GPU utilization using information gathered using PC sampling. Given the average clock frequency and the sampling rate, if all SMs are active, then HPCToolkit knows how many instruction samples would be expected (**GSAMP:EXP**) if the GPU was fully active for the interval when it was in use. HPCToolkit approximates the percentage of GPU utilization by comparing the measured samples with the expected samples using the following formula: $100 * (\text{GSAMP : TOT}) / (\text{GSAMP : EXP})$.

For CUDA 10, measurement using PC sampling with CUPTI serializes the execution of GPU kernels. Thus, measurement of GPU kernels using PC sampling will distort the execution of a GPU-accelerated application by blocking concurrent execution of GPU kernels. For applications that rely on concurrent kernel execution to keep the GPU busy, this will significantly distort execution and PC sampling measurements will only reflect the GPU activity of kernels running in isolation.

Metric	Description
GKER:STMEM (B)	GPU kernel: static memory (bytes)
GKER:DYMEM (B)	GPU kernel: dynamic memory (bytes)
GKER:LMEM (B)	GPU kernel: local memory (bytes)
GKER:FGP_ACT	GPU kernel: fine-grain parallelism, actual
GKER:FGP_MAX	GPU kernel: fine-grain parallelism, maximum
GKER:THR_REG	GPU kernel: thread register count
GKER:BLK_THR	GPU kernel: thread count
GKER:BLK_SM (B)	GPU kernel: block local memory (bytes)
GKER:COUNT	GPU kernel: launch count
GKER:OCC_THR	GPU kernel: theoretical occupancy

Table 6.7: GPU kernel characteristic metrics.

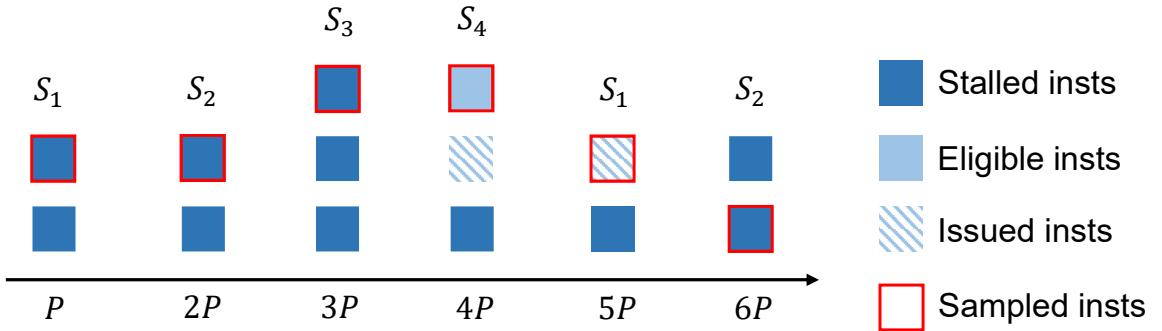


Figure 6.1: NVIDIA’s GPU PC sampling example on an SM. $P - 6P$ represent six sample periods P cycles apart. $S_1 - S_4$ represent four schedulers on an SM.

6.1.2 Attributing Measurements to Source Code for NVIDIA GPUs

NVIDIA’s `nvcc` compiler doesn’t record information about how GPU machine code maps to CUDA source without proper compiler arguments. Using the `-G` compiler option to `nvcc`, one may generate NVIDIA CUBINs with full DWARF information that includes not only line maps, which map each machine instruction back to a program source line, but also detailed information about inlined code. However, the price of turning on `-G` is that optimization by `nvcc` will be disabled. For that reason, the performance of code compiled `-G` is vastly slower. While a developer of a template-based programming model may find this option useful to see how a program employs templates to instantiate GPU code, measurements of code compiled with `-G` should be viewed with skeptical eye.

One can use `nvcc`’s `-lineinfo` option to instruct `nvcc` to record line map information during compilation.¹ The `-lineinfo` option can be used in conjunction with `nvcc` optimization. Using `-lineinfo`, one can measure and interpret the performance of optimized code. However, line map information is a poor substitute for full DWARF information. When `nvcc` inlines code during optimization, the resulting line map information simply shows that source lines that were compiled into a GPU function. A developer examining

¹Line maps relate each machine instruction back to the program source line from where it came.

Metric	Description
GINST	GPU instructions executed
GINST:STL_ANY	GPU instruction stalls: any
GINST:STL_NONE	GPU instruction stalls: no stall
GINST:STL_IFET	GPU instruction stalls: await availability of next instruction (fetch or branch delay)
GINST:STL_IDEP	GPU instruction stalls: await satisfaction of instruction input dependence
GINST:STL_GMEM	GPU instruction stalls: await completion of global memory access
GINST:STL_TMEM	GPU instruction stalls: texture memory request queue full
GINST:STL_SYNC	GPU instruction stalls: await completion of thread or memory synchronization
GINST:STL_CMEM	GPU instruction stalls: await completion of constant or immediate memory access
GINST:STL_PIPE	GPU instruction stalls: await completion of required compute resources
GINST:STL_MTHR	GPU instruction stalls: global memory request queue full
GINST:STL_NSEL	GPU instruction stalls: not selected for issue but ready
GINST:STL_OTHR	GPU instruction stalls: other
GINST:STL_SLP	GPU instruction stalls: sleep

Table 6.8: GPU instruction execution and stall metrics.

Metric	Description
GSAMP:DRP	GPU PC samples: dropped
GSAMP:EXP	GPU PC samples: expected
GSAMP:TOT	GPU PC samples: measured
GSAMP:PER (cyc)	GPU PC samples: period (GPU cycles)
GSAMP:UTIL (%)	GPU utilization computed using PC sampling

Table 6.9: GPU PC sampling statistics.

performance measurements for a function must reason on their own about how any source lines from outside the function got there as the result of inlining and/or macro expansion.

When HPCToolkit uses NVIDIA’s CUPTI to monitor a GPU-accelerated application, CUPTI notifies HPCToolkit every time it loads a CUDA binary, known as a CUBIN, into a GPU. At runtime, HPCToolkit computes a cryptographic hash of a CUBIN’s contents and records the CUBIN into the execution’s measurement directory. For instance, if a GPU-accelerated application loaded CUBIN into a GPU, NVIDIA’s CUPTI informed HPCToolkit that the CUBIN was being loaded, and HPCToolkit computed its cryptographic hash as

972349aed8, then HPCToolkit would record `972349aed8.cubin` inside a `cubins` subdirectory of an HPCToolkit measurement directory.

To attribute GPU performance measurements back to source, HPCToolkit's `hpcstruct` supports analysis of NVIDIA CUBIN binaries. Since many CUBIN binaries may be loaded by a GPU-accelerated application during execution, an application's measurements directory may contain a `cubins` subdirectory populated with many CUBINs. In this case, it would be inconvenient to require a developer to apply `hpcstruct` manually to analyze each CUBIN. To simplify the analysis of an execution's CUBINs, a developer may apply HPCToolkit's `hpcstruct` directly to a measurement directory to analyze all of the CUBINs it contains. Namely, for a measurements directory `hpctoolkit-laghos-measurements` collected during an execution of the GPU-accelerated laghos mini-app [7], one can analyze all of the CUBINs collected during execution by using the following command:

```
hpcstruct hpctoolkit-laghos-measurements
```

Since there may be many CUBINs inside an HPCToolkit measurements directory for a GPU-accelerated application, it is useful to accelerate the analysis of an execution's CUBINs by employing parallelism. One can analyze large or many CUBINs in an HPCToolkit measurements directory using multiple threads by specifying `hpcstruct`'s '`-j n`' option. For instance, one can use 16 threads to analyze the CUBINs in the `hpctoolkit-laghos-measurements` directory with the following command:

```
hpcstruct -j 16 hpctoolkit-laghos-measurements
```

With the default settings of `hpcstruct`, if the '`-j n`' argument is specified, each CUBIN larger than 100MB will be analyzed using n threads. Smaller CUBINs will be analyzed in parallel by using a pool of n threads to analyze n small CUBINs concurrently.

By default, when applied to a measurements directory, `hpcstruct` performs only lightweight analysis of the GPU functions in each CUBIN. When a measurements directory contains fine-grain measurements collected using PC sampling, it is useful to perform a more detailed analysis to recover information about the loops and call sites of GPU functions in an NVIDIA CUBIN. Unfortunately, NVIDIA has refused to provide an API that would enable HPCToolkit to perform instruction-level analysis of CUBINs directly. Instead, HPCToolkit must invoke NVIDIA's `nvdismasm` command line utility to compute control flow graphs for functions in a CUBIN. The version of `nvdismasm` in CUDA 10 is slow and fails to compute control flow graphs for some GPU functions. In such cases, `hpcstruct` reverts to lightweight analysis of GPU functions that considers only line map information. Because analysis of CUBINs using `nvdismasm` is slow, it is not performed by default. To enable detailed analysis of GPU functions, use the `--gpucfg yes` option to `hpcstruct`, as shown below:

```
hpcstruct -j 16 --gpucfg yes hpctoolkit-laghos-measurements
```

6.1.3 GPU Calling Context Tree Reconstruction

The CUPTI API returns flat PC samples without any information about GPU call stacks. With complex code generated from template-based GPU programming models,

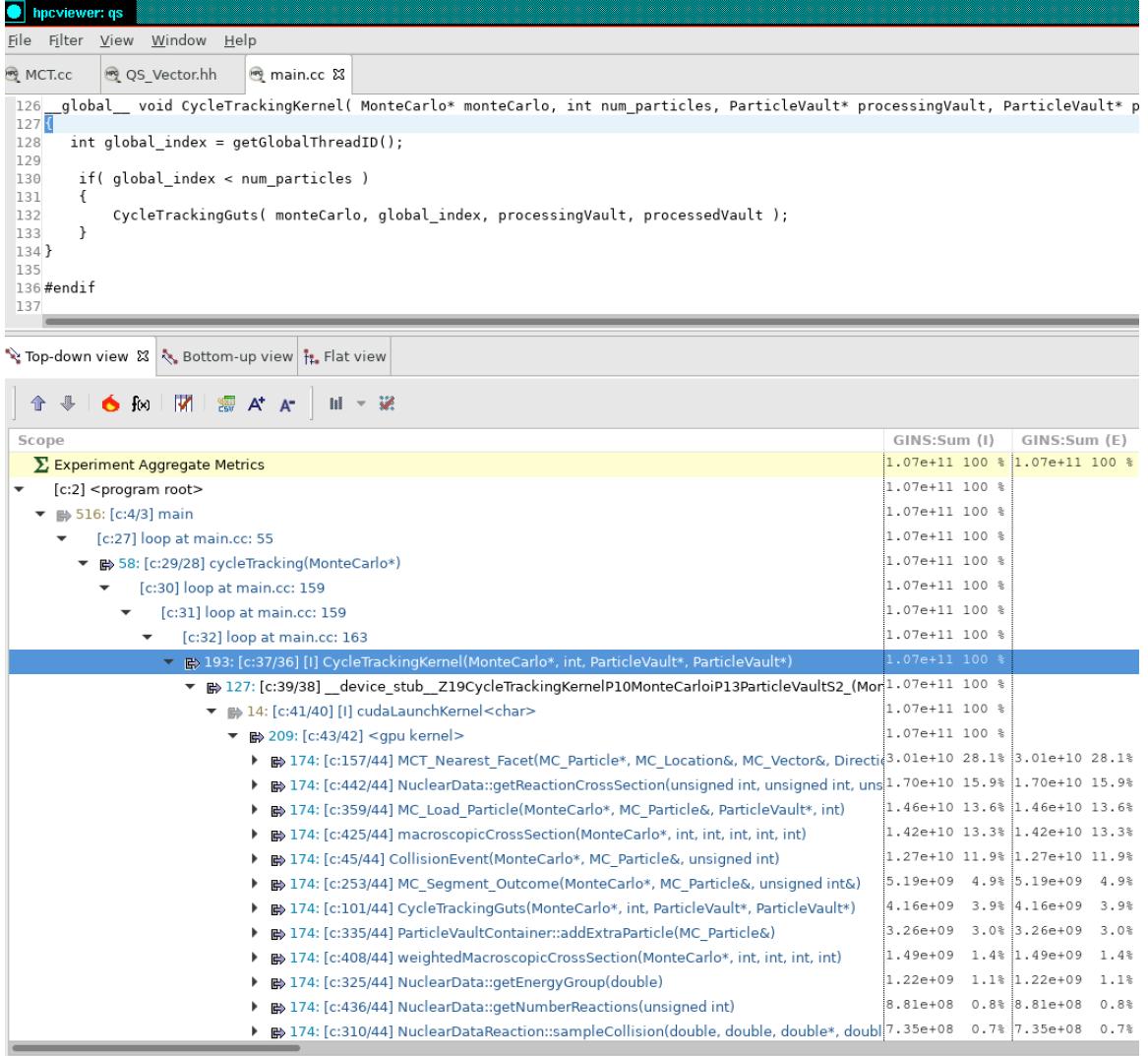


Figure 6.2: A screenshot of hpcviewer for the GPU-accelerated Quicksilver proxy app without GPU CCT reconstruction.

calling contexts on GPUs are essential for developers to understand the code and its performance. Lawrence Livermore National Laboratory's GPU-accelerated Quicksilver proxy app [8] illustrates this problem. Figure 6.2 shows a hpcviewer screenshot of Quicksilver without approximate reconstruction the GPU calling context tree. The figure shows a top-down view of heterogeneous calling contexts that span both the CPU and GPU. In the middle of the figure is a placeholder <gpu kernel> that is inserted by HPCToolkit. Above the placeholder is a CPU calling context where a GPU kernel was invoked. Below the <gpu kernel> placeholder, hpcviewer shows a dozen of the GPU functions that were executed on behalf of the GPU kernel CycleTrackingKernel.

Currently, no API is available for efficiently unwinding call stacks on NVIDIA's GPUs. To address this issue, we designed a method to reconstruct approximate GPU calling con-

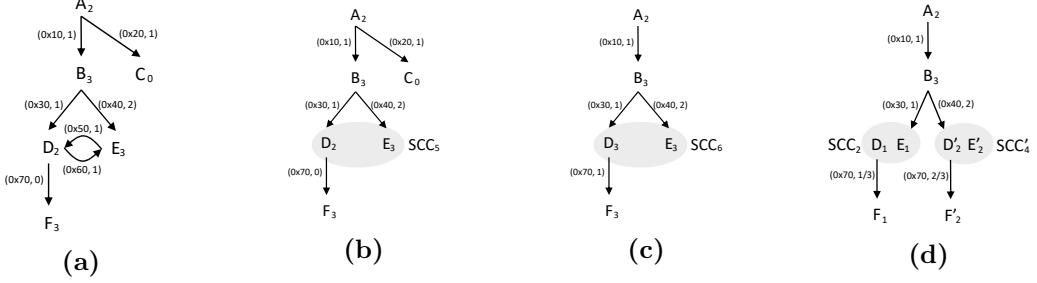


Figure 6.3: Reconstruct a GPU calling context tree. A-F represent GPU functions. Each subscript denotes the number of samples associated with the function. Each (a, c) pair indicates an edge at address a has c call instruction samples.

texts using post-mortem analysis. This analysis is only performed when (1) an execution has been monitored using PC sampling, and (2) an execution’s CUBINs have analyzed in detail using `hpcstruct` with `--gpucfg yes` option.

To reconstruct approximate calling context trees for GPU computations, HPCToolkit uses information about call sites identified by `hpcstruct` in conjunction with PC samples measured for each `call` instruction in GPU binaries.

Without the ability to measure each function invocation in detail, HPCToolkit assumes that each invocation of a particular GPU function incurs the same costs. The costs of each GPU function are apportioned among its caller or callers using the following rules:

- If a GPU function G can only be invoked from a single call site, all of the measured cost of G will be attributed to its call site.
- If a GPU function G can be called from multiple call sites and PC samples have been collected for one or more of the call instructions for G, the costs for G are proportionally divided among G’s call sites according to the distribution of PC samples for calls that invoke G. For instance, consider the case where there are three call sites where G may be invoked, 5 samples are recorded for the first call instruction, 10 samples are recorded for the second call instruction, and no samples are recorded for the third call. In this case, HPCToolkit divides the costs for G among the first two call sites, attributing $5/15$ of G’s costs to the first call site and $10/15$ of G’s costs to the second call site.
- If no call instructions for a GPU function G have been sampled, the costs of G are apportioned evenly among each of G’s call sites.

IHPCToolkit’s `hpcprof` analyzes the static call graph associated with each GPU kernel invocation. If the static call graph for the GPU kernel contains cycles, which arise from recursive or mutually-recursive calls, `hpcprof` replaces each cycle with a strongly connected component (SCC). In this case, `hpcprof` unlinks call graph edges between vertices within the SCC and adds an SCC vertex to enclose the set of vertices in each SCC. The rest of `hpcprof`’s analysis treats an SCC vertex as a normal “function” in the call graph.

Figure 6.3 illustrates the reconstruction of an approximate calling context tree for a GPU computation given the static call graph (computed by `hpcstruct` from a CU-

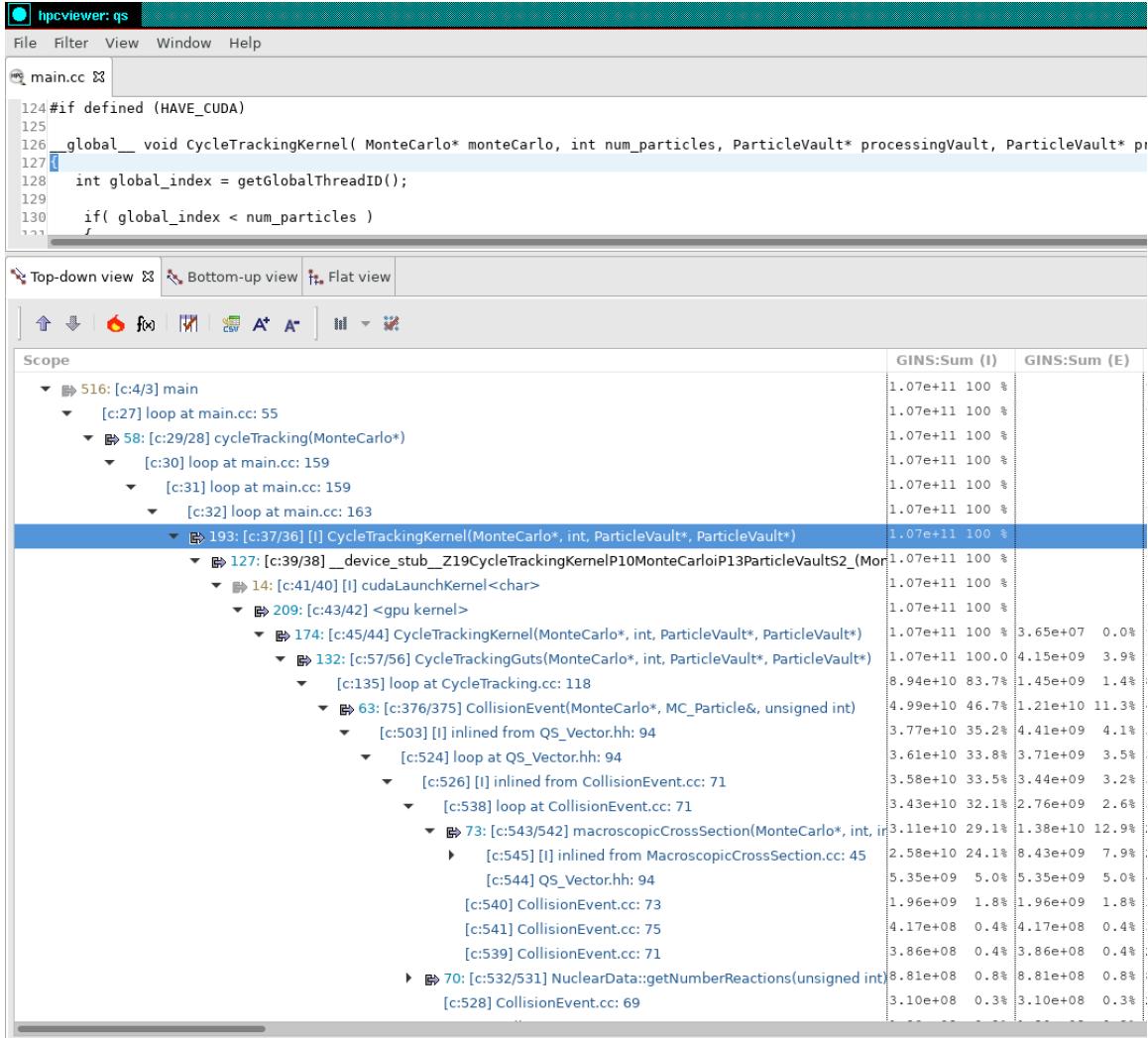


Figure 6.4: A screenshot of `hpcviewer` for the GPU-accelerated Quicksilver proxy app with GPU CCT reconstruction.

BIN's machine instructions) and PC sample counts for some or all GPU instructions in the CUBIN. Figure 6.4 shows an `hpcviewer` screenshot for the GPU-accelerated Quicksilver proxy app following reconstruction of GPU calling contexts using the algorithm described in this section. Notice that after the reconstruction, one can see that `CycleTrackingKernel` calls `CycleTrackingGuts`, which calls `CollisionEvent`, which eventually calls `macroscopicCrossSection` and `NuclearData::getNumberOfReactions`. The rich approximate GPU calling context tree reconstructed by `hpcprof` also shows loop nests and inlined code.²

²The control flow graph used to produce this reconstruction for Quicksilver was computed with a pre-release version of `nvdismasm`. You will not be able to reproduce these results with CUDA 10's `nvdismasm` because the released version of `nvdismasm` fails to analyze some of the key functions in Quicksilver.

Argument to <code>hpcrun</code>	What is monitored
<code>-e gpu=amd</code>	profiling of AMD GPU operations
<code>-e gpu=amd -t</code>	profiling and tracing of AMD GPU operations

Table 6.10: Monitoring performance on AMD GPUs.

6.2 AMD GPUs

HPCToolkit provides *alpha*-level support for coarse-grain profiling of GPU-accelerated applications that offload computation onto AMD GPUs using AMD’s HIP programming model. Table 6.10 shows arguments to `hpcrun` to monitor the performance of GPU operations on AMD GPUs. With this coarse-grain profiling support, HPCToolkit can collect GPU operation timings (Table 6.2) and a subset of standard metrics for GPU operations such as memory allocation and deallocation (Table 6.3), memory set (Table 6.4), explicit memory copies (Table 6.5), and synchronization (Table 6.6).

At present, HPCToolkit’s support for monitoring on AMD GPUs is incomplete. The HPCToolkit team is awaiting enhancements to AMD’s ROC-tracer monitoring library that are needed to complete support for adding GPU kernels to call paths collected at the point where GPU operations are invoked. Currently, the leaves of call paths for GPU kernel invocations end in the placeholder `<gpu kernel>` and lack the name of the kernel invoked. This more greatly affects trace lines for GPU streams in Trace view where one can only see GPU activity marked by `<gpu kernel>` without being able to readily identify the specific kernels being invoked. The HPCToolkit project team is awaiting completion of additional support in AMD’s ROC-tracer monitoring library to address this deficiency.

At present, ROC-tracer lacks support for fine-grain performance measurement of GPU computations on AMD GPUs like the PC sampling supported by NVIDIA’s GPUs.

6.3 Intel GPUs

The released version of HPCToolkit does not yet support measurement and analysis of performance on Intel GPUs. A development branch supports measurement of GPU computations launched on Intel GPUs using OpenCL or Intel’s new Level 0 runtime. A release of these capabilities is forthcoming.

Chapter 7

Measurement and Analysis of OpenMP Multithreading

Note: This release contains beta support for monitoring OpenMP computations using the OpenMP 5.0 OpenMP Tools API support known as OMPT.

HPCToolkit includes an implementation of the OpenMP 5.0 Tools API known as OMPT. The OMPT interface enables HPCToolkit to extract enough information to reconstruct user-level calling contexts from implementation-level measurements.

Support for OpenMP 5.0 and OMPT is emerging in OpenMP runtimes. Recent versions of LLVM's OpenMP runtime, IBM's LOMP (Lightweight OpenMP Runtime) and Intel's OpenMP runtime provide emerging support for OMPT. At present, support in these implementations is known to be incomplete, especially with respect to offloading computation onto TARGET devices.

If an interaction between HPCToolkit's support for the OMPT interface and an OpenMP runtime causes problems, OMPT support may be disabled when using HPCToolkit by setting OpenMP environment variable `OMP_TOOL` to `disabled`.

Chapter 8

The `hpcviewer` User Interface

HPCTOOLKIT provides the `hpcviewer` [2, 19] performance presentation tool for interactive examination of performance databases. `hpcviewer` presents a heterogenous calling context tree that spans both CPU and GPU contexts, annotated with measured or derived metrics to help users assess code performance and identify bottlenecks.

The database generated by `hpcprof` consists of 4 dimensions: *profile*, *time*, *context*, and *metric*. We employ the term *profile* to include any logical threads (such as OpenMP, pthread and C++ threads), and also MPI processes and GPU streams. The *time* dimension represents the timeline of the program’s execution, *context* depicts the path in calling-context tree, and *metric* constitutes program measurements performed by `hpcreun` such as cycles, number of instructions, stall percentages and ratio of idleness. The *time* dimension is available if the application is profiled with `hpcreun -t` option.

To simplify performance data visualization, `hpcviewer` restricts display two dimensions at a time: the *profile viewer* displays pairs of \langle context, metric \rangle or \langle profile, metric \rangle dimensions; and the *trace viewer* visualizes the behavior of threads or streams over time. To enable to display a pair of \langle profile, metric \rangle dimensions, the database has to be generated with `hpcprof --metric-db yes` option.

8.1 Launching

Requirements to launch `hpcviewer`:

- On all platforms: Java 8 or newer.
- On Linux: GTK.

`hpcviewer` can either be launched from a command line (Linux platforms) or by clicking the `hpcviewer` icon (for Windows, Mac OS X and Linux platforms). The command line syntax is as follows:

```
hpcviewer [options] [<hpctoolkit-database>]
```

Here, `<hpctoolkit-database>` is an optional argument to load a database automatically. Without this argument, `hpcviewer` will prompt for the location of a database. Possible

options for `hpcviewer` are shown in the table below:

<code>-h, --help</code>	Print a help message.
<code>-jh, --java-heap size</code>	Set the JVM maximum heap size for this execution of <code>hpcviewer</code> . The value of <i>size</i> must be in megabytes (M) or gigabytes (G). For example, one can specify a <i>size</i> of 3 gigabytes as either 3076M or 3G.

On Linux, when `hpcviewer` is installed using its `install` script, the `install` script chooses a default maximum size for the Java heap on the current platform. When analyzing measurements for large and complex applications, it may be necessary to use the `--java-heap` option to specify a larger heap size for `hpcviewer` to accommodate many metrics for many contexts.

8.2 Profile View

This view is the default view and displays the pair of \langle context, metric \rangle dimensions. It interactively presents context-sensitive performance metrics correlated to program structure and mapped to a program’s source code, if available. It can present an arbitrary collection of performance metrics gathered during one or more runs or compute derived metrics.

Figure 8.1 shows an annotated screenshot of `hpcviewer`’s user interface presenting a call path profile. The annotations highlight `hpcviewer`’s principal window panes and key controls. The browser window is divided into three panes. The Source pane (top) displays program source code. The Navigation and Metric panes (bottom) associate a table of performance metrics with static or dynamic program structure. These panes are discussed in more detail in Section 8.3.

`hpcviewer` displays calling-context-sensitive performance data in three different views: a top-down *Top-down View*, a bottom-up *Bottom-up View*, and a *Flat View*. One selects the desired view by clicking on the corresponding view control tab. We briefly describe the three views and their corresponding purposes.

- **Top-down View.** This top-down view shows the dynamic calling contexts (call paths) in which costs were incurred. Using this view, one can explore performance measurements of an application in a top-down fashion to understand the costs incurred by calls to a procedure in a particular calling context. We use the term *cost* rather than simply *time* since `hpcviewer` can present a multiplicity of metrics such as cycles, or cache misses) or derived metrics (e.g. cache miss rates or bandwidth consumed) that are other indicators of execution cost.

A calling context for a procedure *f* consists of the stack of procedure frames active when the call was made to *f*. Using this view, one can readily see how much of the application’s cost was incurred by *f* when called from a particular calling context. If finer detail is of interest, one can explore how the costs incurred by a call to *f* in a particular context are divided between *f* itself and the procedures it calls. HPCTOOLKIT’s call path profiler `hpcrun` and the `hpcviewer` user interface distinguish calling context

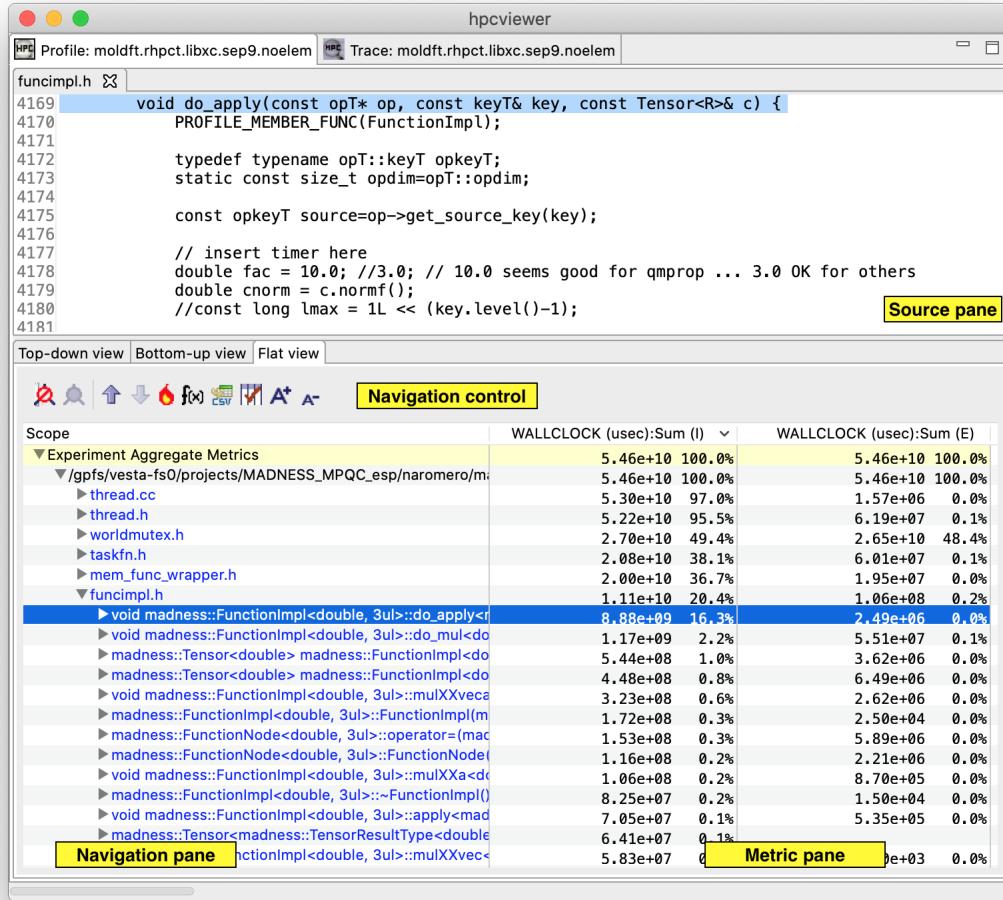


Figure 8.1: An annotated screenshot of hpcviewer's interface.

precisely by individual call sites; this means that if a procedure *g* contains calls to procedure *f* in different places, these represent separate calling contexts.

- **Bottom-up View.** This bottom-up view enables one to look upward along call paths. The view apportions a procedure's costs to its callers and, more generally, its calling contexts. This view is particularly useful for understanding the performance of software components or procedures that are used in more than one context. For instance, a message-passing program may call `MPI_Wait` in many different calling contexts. The cost of any particular call will depend upon the structure of the parallelization in which the call is made. Serialization or load imbalance may cause long waits in some calling contexts while other parts of the program may have short waits because computation is balanced and communication is overlapped with computation.

When several levels of the Bottom-up View are expanded, saying that the Bottom-up View apportions metrics of a callee on behalf of its callers can be confusing. More precisely, the Bottom-up View apportions the metrics of a procedure on behalf of the various *calling contexts* that reach it.

- **Flat View.** This view organizes performance measurement data according to the static structure of an application. All costs incurred in any calling context by a procedure are aggregated together in the Flat View. This complements the Top-down View, in which the costs incurred by a particular procedure are represented separately for each call to the procedure from a different calling context.

8.3 Panes

`hpcviewer`'s browser window is divided into three panes: the *Navigation pane*, *Source pane*, and the *Metrics pane*. We briefly describe the role of each pane.

8.3.1 Source Pane

The source pane displays the source code associated with the current entity selected in the navigation pane. When a performance database is first opened with `hpcviewer`, the source pane is initially blank because no entity has been selected in the navigation pane. Selecting any entity in the navigation pane will cause the source pane to load the corresponding file, scroll to and highlight the line corresponding to the selection. Switching the source pane to view to a different source file is accomplished by making another selection in the navigation pane.

8.3.2 Navigation Pane

The navigation pane presents a hierarchical tree-based structure that is used to organize the presentation of an application's performance data. Entities that occur in the navigation pane's tree include load modules, files, procedures, procedure activations, inlined code, loops, and source lines. Selecting any of these entities will cause its corresponding source code (if any) to be displayed in the source pane. One can reveal or conceal children in this hierarchy by 'opening' or 'closing' any non-leaf (i.e., individual source line) entry in this view.

The nature of the entities in the navigation pane's tree structure depends upon whether one is exploring the Top-down View, the Bottom-up View, or the Flat View of the performance data.

- In the **Top-down View**, entities in the navigation tree represent procedure activations, inlined code, loops, and source lines. While most entities link to a single location in source code, procedure activations link to two: the call site from which a procedure was called and the procedure itself.
- In the **Bottom-up View**, entities in the navigation tree are procedure activations. Unlike procedure activations in the top-down view in which call sites are paired with the called procedure, in the bottom-up view, call sites are paired with the calling

procedure to facilitate attribution of costs for a called procedure to multiple different call sites and callers.

- In the **Flat View**, entities in the navigation tree correspond to source files, procedure call sites (which are rendered the same way as procedure activations), loops, and source lines.

Navigation Control

The header above the navigation pane contains some controls for the navigation and metric view. In Figure 8.1, they are labeled as “navigation/metric control.”

- **Flatten  / Unflatten ** (available for the Flat View):

Enabling to flatten and unflatten the navigation hierarchy. Clicking on the flatten button (the icon that shows a tree node with a slash through it) will replace each top-level scope shown with its children. If a scope has no children (i.e., it is a leaf), the node will remain in the view. This flattening operation is useful for relaxing the strict hierarchical view so that peers at the same level in the tree can be viewed and ranked together. For instance, this can be used to hide procedures in the Flat View so that outer loops can be ranked and compared to one another. The inverse of the flatten operation is the unflatten operation, which causes an elided node in the tree to be made visible once again.

- **Zoom-in  / Zoom-out **:

Depressing the up arrow button will zoom in to show only information for the selected line and its descendants. One can zoom out (reversing a prior zoom operation) by depressing the down arrow button.

- **Hot call path **:

This button is used to automatically reveal and traverse the hot call path rooted at the selected node in the navigation pane with respect to the selected metric column. Let n be the node initially selected in the navigation pane. A hot path from n is traversed by comparing the values of the selected metric for n and its children. If one child accounts for $T\%$ or more (where T is the threshold value for a hot call path) of the cost at n , then that child becomes n and the process repeats recursively. The default value for T is 50. One can change T by using the menu File—Preferences—hpcviewer preferences.

- **Add derived metric **:

Create a new metric by specifying a mathematical formula. See Section 8.7 for more details.

- **Hide/show metrics **:

Show or hide metric columns. A dialog box will appear and the user can select which metric columns should be shown. See Section 8.10.2 section for more details.

- **Export into a CSV format file**  :

Export the current metric table into a comma separated value (CSV) format file. This feature only exports all metrics that are currently shown. Metrics that are not shown in the view (whose scopes are not expanded) will not be exported (we assume these metrics are not significant).

- **Increase font size  / Decrease font size ** :

Increase or decrease the size of the navigation and metric panes.

- **Show a graph of metric values**  :

Show a graph (a plot, a sorted plot or a histogram) of metric values associated with the selected node in CCT for all processes or threads (Section 8.8.1). This menu is only available if the database is generated by `hpcprof-mpi` instead of `hpcprof`.

- **Show the metrics of a set of threads**  :

Show the CCT and the metrics of a selected threads (Section 8.8.2). This menu is only available if the database is generated by `hpcprof-mpi` instead of `hpcprof`.

Context menus

Navigation control also provides several context menus by clicking the right-button of the mouse.

- **Copy to clipboard:** Copy into clipboard the selected line in navigation pane which includes the name of the node in the tree, and the values of visible metrics in metric pane (Section 8.3.3). The values of hidden metrics will not be copied.

8.3.3 Metric Pane

The metric pane displays one or more performance metrics associated with entities to the left in the navigation pane. Entities in the tree view of the navigation pane are sorted at each level of the hierarchy by the metric in the selected column. When `hpcviewer` is launched, the leftmost metric column is the default selection and the navigation pane is sorted according to the values of that metric in descending order. One can change the selected metric by clicking on a column header. Clicking on the header of the selected column toggles the sort order between descending and ascending.

During analysis, one often wants to consider the relationship between two metrics. This is easier when the metrics of interest are in adjacent columns of the metric pane. One can change the order of columns in the metric pane by selecting the column header for a metric and then dragging it left or right to its desired position. The metric pane also includes scroll bars for horizontal scrolling (to reveal other metrics) and vertical scrolling (to reveal other scopes). Vertical scrolling of the metric and navigation panes is synchronized.

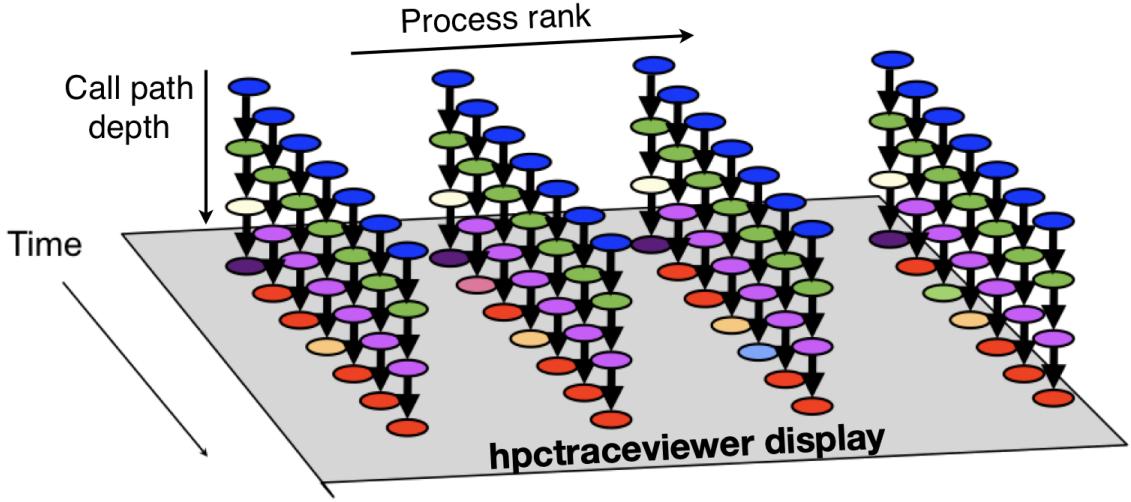


Figure 8.2: Logical view of trace call path samples on three dimensions: time, process rank and call path depth.

8.4 Trace view

Trace view [19] is a time-centric user interface for interactive examination of a sample-based time series (hereafter referred to as a trace) view of a program execution. Trace view can interactively present a large-scale execution trace without concern for the scale of parallelism it represents.

To collect a trace for a program execution, one must instruct HPCTOOLKIT’s measurement system to collect a trace. When launching a dynamically-linked executable with `hpcrun`, add the `-t` flag to enable tracing. When launching a statically-linked executable, set the environment variable `HPCRUN_TRACE=1` to enable tracing. When collecting a trace, one must also specify a metric to measure. The best way to collect a useful trace is to asynchronously sample the execution with a time-based metric such as `REALTIME`, `CYCLES`, or `WALLCLOCK` (Blue Gene/Q only).

As shown in Figure 8.2, call path traces consist of data in three dimensions: *process/thread rank*, *time*, and *call path depth*. A *crosshair* in Trace view is defined by a triplet (p, t, d) where p is the selected process/thread rank, t is the selected time, and d is the selected call path depth.

Trace view’s *Trace View* (Section 8.4.1) renders a view of processes and threads over time. Trace view’s *Depth View* (Section 8.4.2) shows the call path depth over time for the thread selected by the cursor. Trace view’s *Call Path View* (Section 8.4.4) shows the call path associated with the thread and time pair specified by the cursor. Each of these views plays a role for understanding an application’s performance.

In Trace view, each procedure is assigned specific color. Figure 8.2 shows that at depth 1 each call path has the same color: blue. This node represents the main program that serves as the root of the call chain in all process at all times. At depth 2, all processes have a green node, which indicates another procedure. At depth 3, in the first time step all

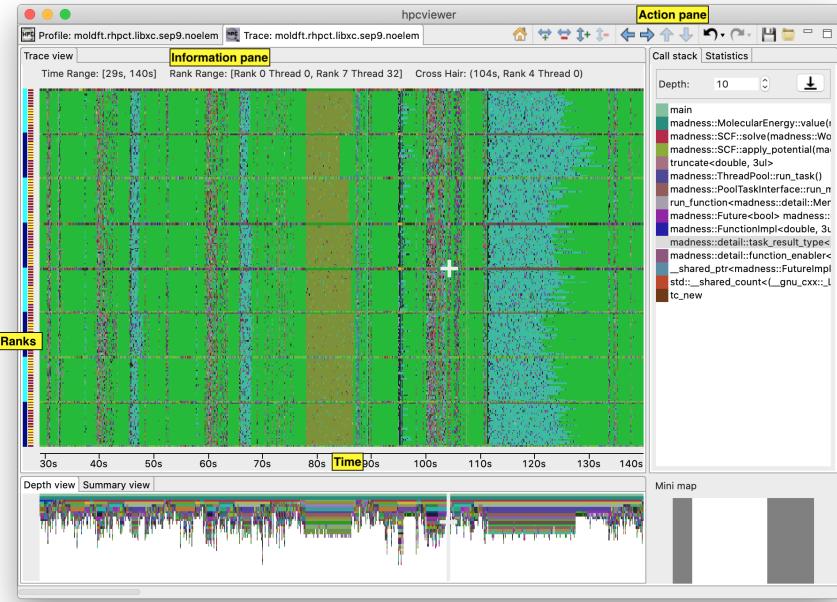


Figure 8.3: A screenshot of Trace view’s interface.

processes have a yellow node; in subsequent time steps they have purple nodes. This might indicate that the processes first are observed in an initialization procedure (represented by yellow) and later observed in a solve procedure (represented by purple). The pattern of colors that appears in a particular depth slice of the Trace View enables a user to visually identify inefficiencies such as load imbalance and serialization.

Figures 8.3 and 8.4 show screenshots of Trace view’s user interface presenting a call path profile. Figure 8.3 highlights Trace view’s four principal window panes: Trace View, Depth View, Call Path View and Mini Map View, while Figure 8.4 shows additional two window panes: Summary View and Statistics View.

- **Trace View** (top, left pane): This is Trace view’s primary view. This view, which is similar to a conventional process/time (or space/time) view, shows time on the horizontal axis and process (or thread) rank on the vertical axis; time moves from left to right. Compared to typical process/time views, there is one key difference. To show call path hierarchy, the view is actually a user-controllable slice of the process/time/call-path space. Given a call path depth, the view shows the color of the currently active procedure at a given time and process rank. (If the requested depth is deeper than a particular call path, then Trace view simply displays the deepest procedure frame and, space permitting, overlays an annotation indicating the fact that this frame represents a shallower depth.)

Trace view assigns colors to procedures based on (static) source code procedures. Although the color assignment is currently random, it is consistent across the different

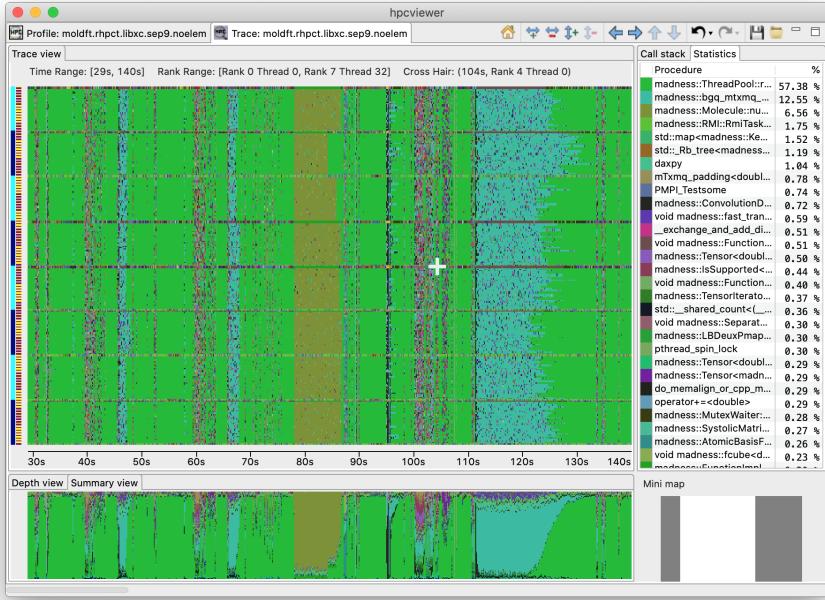


Figure 8.4: A screenshot of Trace view’s interface showing the Summary View and Statistics View.

views. Thus, the same color within the Trace and Depth Views refers to the same procedure.

The Trace View has a white crosshair that represents a selected point in time and process space. For this selected point, the Call Path View shows the corresponding call path. The Depth View shows the selected process.

- **Depth View** (tab in bottom, left pane): This is a call-path/time view for the process rank selected by the Trace View’s crosshair. Given a process rank, the view shows for each virtual time along the horizontal axis a stylized call path along the vertical axis, where ‘main’ is at the top and leaves (samples) are at the bottom. In other words, this view shows for the whole time range, in qualitative fashion, what the Call Path View shows for a selected point. The horizontal time axis is exactly aligned with the Trace View’s time axis; and the colors are consistent across both views. This view has its own crosshair that corresponds to the currently selected time and call path depth.
- **Summary View** (tab in bottom, left pane): The view shows for the whole time range displayed, the proportion of each subroutine in a certain time. Similar to Depth view, the time range in Summary reflects to the time range in the Trace view.
- **Call Path View** (tab in top, right pane): This view shows two things: (1) the current call path depth that defines the hierarchical slice shown in the Trace View; and (2) the actual call path for the point selected by the Trace View’s crosshair. (To easily coordinate the call path depth value with the call path, the Call Path View currently

suppresses details such as loop structure and call sites; we may use indentation or other techniques to display this in the future.)

- **Statistics View** (tab in top, right pane): This view shows the list of procedures active in the space-time region shown in the Trace View at the current Call Path Depth. Each procedure's percentage in the Statistics View indicates the percentage of pixels in the Trace View pane that are filled with this procedure's color at the current Call Path Depth. When the Trace View is navigated to show a new time-space interval or the Call Path Depth is changed, the statistics view will update its list of procedures and the percentage of execution time to reflect the new space-time interval or depth selection.
- **Mini Map View** (right, bottom): The Mini Map shows, relative to the process/time dimensions, the portion of the execution shown by the Trace View. The Mini Map enables one to zoom and to move from one close-up to another quickly.

8.4.1 Trace View

Trace View is divided into two parts: the top part which contains *action pane* and the *information pane*, and the main view which displays the traces.

The buttons in the action pane are the following:

- **Home**  : Resetting the view configuration into the original view, i.e., viewing traces for all times and processes.
- **Horizontal zoom in**  / **out**  : Zooming in/out the time dimension of the traces.
- **Vertical zoom in**  / **out**  : Zooming in/out the process dimension of the traces.
- **Navigation buttons** , , ,  : Navigating the trace view to the left, right, up and bottom, respectively. It is also possible to navigate with the arrow keys in the keyboard. Since Trace View does not support scroll bars, the only way to navigate is through navigation buttons (or arrow keys).
- **Undo**  : Canceling the action of zoom or navigation and returning back to the previous view configuration.
- **Redo**  : Redoing of previously undo change of view configuration.
- **Save**  / **Open**  **a view configuration** : Saving/loading a saved view configuration. A view configuration file contains the information about the process/thread and time ranges shown, the selected depth, and the position of the crosshair. It is recommended to store the view configuration file in the same directory as the database to ensure that the view configuration file matches the database since a configuration does not store its associated database. Although it is possible to open a view configuration file associated with a different database, it is not recommended since each database has different time/process dimensions and depth.

At the top of an execution's Trace View pane is some information about the data shown in the pane.

- **Time Range.** The time interval shown along the horizontal dimension.
- **Rank Range.** The range of process/thread ranks shown along the vertical dimension.
- **Cross Hair.** The crosshair indicates the current cursor position in the time and rank dimensions.

8.4.2 Depth View

Depth View shows all the call path for a certain time range $[t_1, t_2] = \{t | t_1 \leq t \leq t_2\}$ in a specified process rank p . The content of Depth View is always consistent with the position of the crosshair in Trace View. For instance once the user clicks in process p and time t , while the current depth of call path is d , then the Depth View's content is updated to display all the call path of process p and shows its crosshair on the time t and the call path depth d .

On the other hand, any user action such as crosshair and time range selection in Depth View will update the content within Trace View. Similarly, the selection of new call path depth in Call Path View invokes a new position in Depth View.

In Depth View a user can specify a new crosshair time and a new time range.

Specifying a new crosshair time. Selecting a new crosshair time t can be performed by clicking a pixel within Depth View. This will update the crosshair in Trace View and the call path in Call Path View.

Selecting a new time range. Selecting a new time range $[t_m, t_n] = \{t | t_m \leq t \leq t_n\}$ is performed by first clicking the position of t_m and drag the cursor to the position of t_n . A new content in Depth View and Trace View is then updated. Note that this action will not update the call path in Call Path View since it does not change the position of the crosshair.

8.4.3 Summary View

Summary View presents the proportion of number of calls of time t across the current displayed rank of process p . Similar to Depth View, the time range in Summary View is always consistent with the time range in Trace View.

8.4.4 Call Path View

This view lists the call path of process p and time t specified in Trace View and Depth View. Figure 8.5 shows a call path from depth 0 to depth 14, and the current depth is 9 as shown in the depth editor (located on the top part of the view).

In this view, the user can select the depth dimension of Trace View by either typing the depth in the depth editor or selecting a procedure in the table of call path.

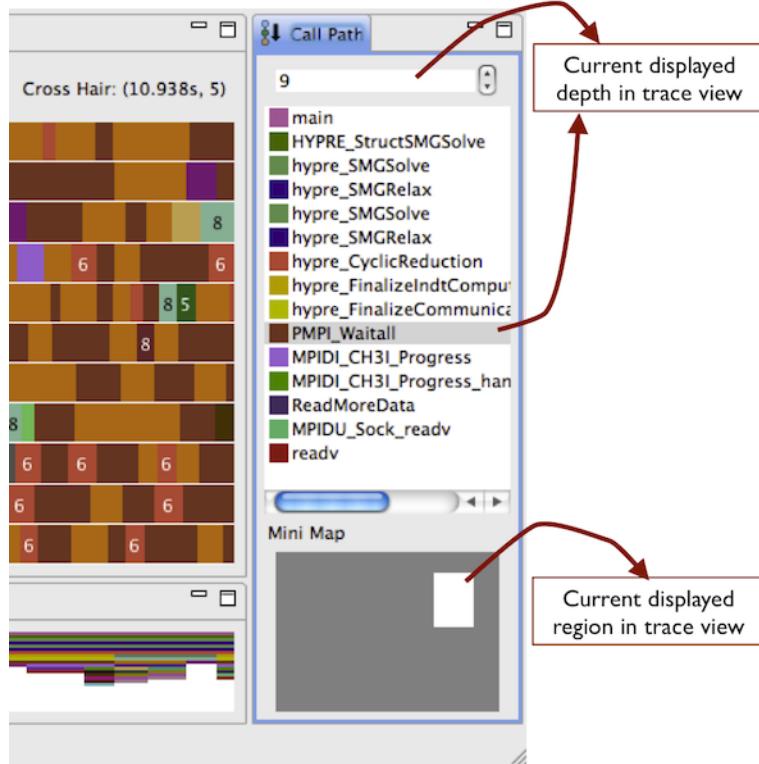


Figure 8.5: An annotated screenshot of Trace view’s Call Path View.

8.4.5 Mini Map View

The Mini Map View shows, relative to the process/time dimensions, the portion of the execution shown by the Trace View. In Mini Map View, the user can select a new process/time $(p_a, t_a), (p_b, t_b)$ dimensions by clicking the first process/time position (p_a, t_a) and then drag the cursor to the second position (p_b, t_b) . The user can also move the current selected region to another region by clicking the white rectangle and drag it to the new place.

Trace view also provides a context menu to save the current image of the view. This context menu is available in three views: trace view, depth view and summary view.

8.5 GPU Threads

Trace view displays activity over time in GPU streams as well as CPU threads. Presently, GPU streams and CPU threads are indexed within a process or MPI rank using a single integer coordinate. To distinguish between GPU threads and CPU threads in this prototype, CPU thread numbers begin at 0 and GPU stream numbers begin at 500.¹

¹A future version of HPCToolkit will clearly distinguish GPU streams from CPU threads. This will require changes to measurement, file formats, post-processing, and the user interfaces. Treating GPU streams like threads with indices beginning at 500 is a workaround until the overhaul is complete.

Unlike CPU threads, which are sampled, each GPU activity (e.g., a kernel execution or a memory copy) on a GPU stream or queue is traced. To avoid giving a misleading view of GPU activity on a stream or queue, intervals between GPU activities where the GPU is idle are marked with the placeholder `<no activity>`.

8.6 Understanding Metrics

`hpcviewer` can present an arbitrary collection of performance metrics gathered during one or more runs, or compute derived metrics expressed as formulae. A derived metric may be specified with a formula that typically uses one or more existing metrics as terms in an expression.

For any given scope in `hpcviewer`'s three views, `hpcviewer` computes both *inclusive* and *exclusive* metric values. First, consider the Top-down View. Inclusive metrics reflect costs for the entire subtree rooted at that scope. Exclusive metrics are of two flavors, depending on the scope. For a procedure, exclusive metrics reflect all costs within that procedure but excluding callees. In other words, for a procedure, costs are exclusive with respect to dynamic call chains. For all other scopes, exclusive metrics reflect costs for the scope itself; i.e., costs are exclusive with respect to static structure. The Bottom-up and Flat Views contain inclusive and exclusive metric values that are relative to the Top-down View. This means, e.g., that inclusive metrics for a particular scope in the Bottom-up or Flat View are with respect to that scope's subtree in the Top-down View.

8.6.1 How Metrics are Computed

Call path profile measurements collected by `hpcrun` correspond directly to the Top-down View. `hpcviewer` derives all other views from exclusive metric costs in the Top-down View. For the Bottom-up View, `hpcviewer` collects the cost of all samples in each function and attribute that to a top-level entry in the Bottom-up View. Under each top-level function, `hpcviewer` can look up the call chain at all of the context in which the function is called. For each function, `hpcviewer` apportions its costs among each of the calling contexts in which they were incurred. `hpcviewer` computes the Flat View by traversing the calling context tree and attributing all costs for a scope to the scope within its static source code structure. The Flat View presents a hierarchy of nested scopes for load modules, files, procedures, loops, inlined code and statements.

8.6.2 Example

Figure 8.6 shows an example of a recursive program separated into two files, `file1.c` and `file2.c`. In this figure, we use numerical subscripts to distinguish between different instances of the same procedure. In the other parts of this figure, we use alphabetic subscripts. We use different labels because there is no natural one-to-one correspondence between the instances in the different views.

Routine `g` can behave as a recursive function depending on the value of the condition branch (lines 3–4). Figure 8.7 shows an example of the call chain execution of the program annotated with both inclusive and exclusive costs. Computation of inclusive costs from

file1.c	file2.c
<pre>f () { g (); } // m is the main routine m () { f (); g (); }</pre>	<pre>// g can be a recursive function g () { if (. .) g (); if (. .) h (); } h () {</pre>

Figure 8.6: A sample program divided into two source files.

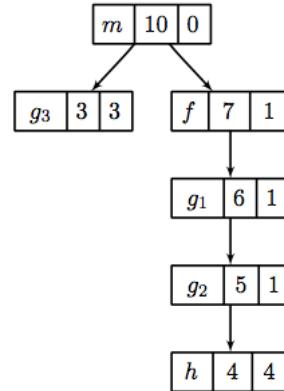


Figure 8.7: Top-down View. Each node of the tree has three boxes: the left-most is the name of the node (or in this case the name of the routine, the center is the inclusive value, and on the right is the exclusive value.

exclusive costs in the Top-down View involves simply summing up all of the costs in the subtree below.

In this figure, we can see that on the right path of the routine `m`, routine `g` (instantiated in the diagram as `g1`) performed a recursive call (`g2`) before calling routine `h`. Although `g1`, `g2` and `g3` are all instances from the same routine (i.e., `g`), we attribute a different cost for each instance. This separation of cost can be critical to identify which instance has a performance problem.

Figure 8.8 shows the corresponding scope structure for the Bottom-up View and the costs we compute for this recursive program. The procedure `g` noted as `ga` (which is a root node in the diagram), has different cost to `g` as a callsite as noted as `gb`, `gc` and `gd`. For instance, on the first tree of this figure, the inclusive cost of `ga` is 9, which is the sum of the highest cost for each path in the calling context tree shown in Figure 8.7 that includes `g`: the inclusive cost of `g3` (which is 3) and `g1` (which is 6). We do not attribute the cost of `g2` here since it is a descendant of `g1` (in other term, the cost of `g2` is included in `g1`).

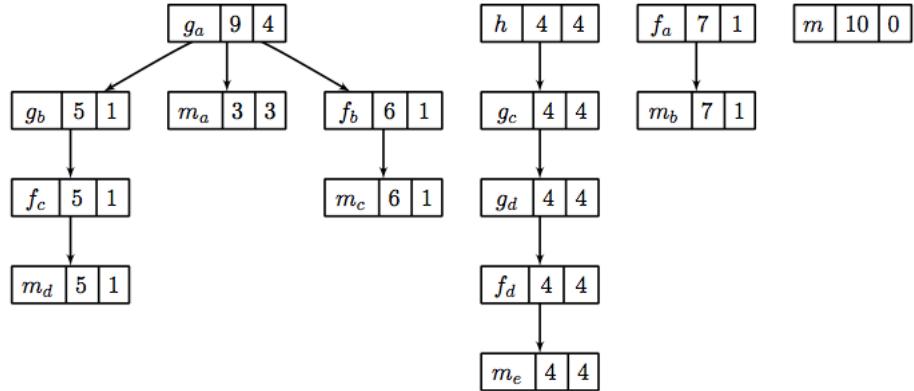


Figure 8.8: Bottom-up View

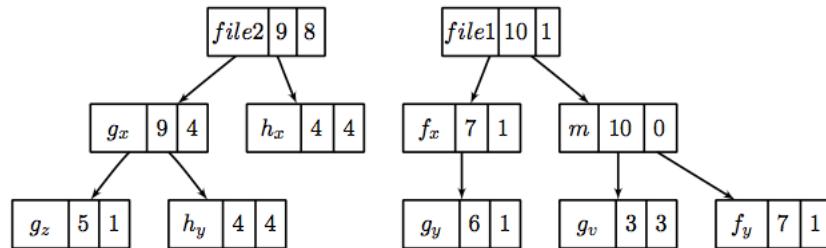


Figure 8.9: Flat View

Inclusive costs need to be computed similarly in the Flat View. The inclusive cost of a recursive routine is the sum of the highest cost for each branch in calling context tree. For instance, in Figure 8.9, The inclusive cost of g_x , defined as the total cost of all instances of g , is 9, and this is consistently the same as the cost in the bottom-up tree. The advantage of attributing different costs for each instance of g is that it enables a user to identify which instance of the call to g is responsible for performance losses.

8.7 Derived Metrics

Frequently, the data become useful only when combined with other information such as the number of instructions executed or the total number of cache accesses. While users don't mind a bit of mental arithmetic and frequently compare values in different columns to see how they relate for a scope, doing this for many scopes is exhausting. To address this problem, `hpcviewer` provides a mechanism for defining metrics. A user-defined metric is called a "derived metric." A derived metric is defined by specifying a spreadsheet-like mathematical formula that refers to data in other columns in the metric table by using `$n` to refer to the value in the n^{th} column.

8.7.1 Formulae

The formula syntax supported by `hpcviewer` is inspired by spreadsheet-like in-fix mathematical formulae. Operators have standard algebraic precedence.

8.7.2 Examples

Suppose the database contains information from five executions, where the same two metrics were recorded for each:

1. Metric 0, 2, 4, 6 and 8: total number of cycles
2. Metric 1, 3, 5, 7 and 9: total number of floating point operations

To compute the average number of cycles per floating point operation across all of the executions, we can define a formula as follows:

```
avg($0, $2, $4. $6. $8) / avg($1, $3, $5, $7, $9)
```

8.7.3 Creating Derived Metrics

A derived metric can be created by clicking the **Derived metric** tool item in the navigation/control pane. A derived metric window will then appear as shown in Figure 8.10.

The window has two main parts:

- **Derived metric definition**, which consists of:
 - *New name for the derived metric.* Supply a string that will be used as the column header for the derived metric. If you don't supply one, the metric will have no name.
 - *Formula definition field.* In this field the user can define a formula with spreadsheet-like mathematical formula. This field must be filled. A user can type a formula into this field, or use the buttons in the Assistance pane below to help insert metric terms or function templates.
 - *Metrics.* This is used to find the *ID* of a metric. For instance, in this snapshot, the metric `WALLCLOCK` has the ID 2. By clicking the button **Insert metric**, the metric ID will be inserted in formula definition field. A metric may refer to the value at an individual node in the calling context tree (point-wise) or the value at the root of the calling context tree (aggregate).
 - *Functions.* This is to guide the user who wants to insert functions in the formula definition field. Some functions require only one metric as the argument, but some can have two or more arguments. For instance, the function `avg()` which computes the average of some metrics, needs at least two arguments.
- **Advanced options:**

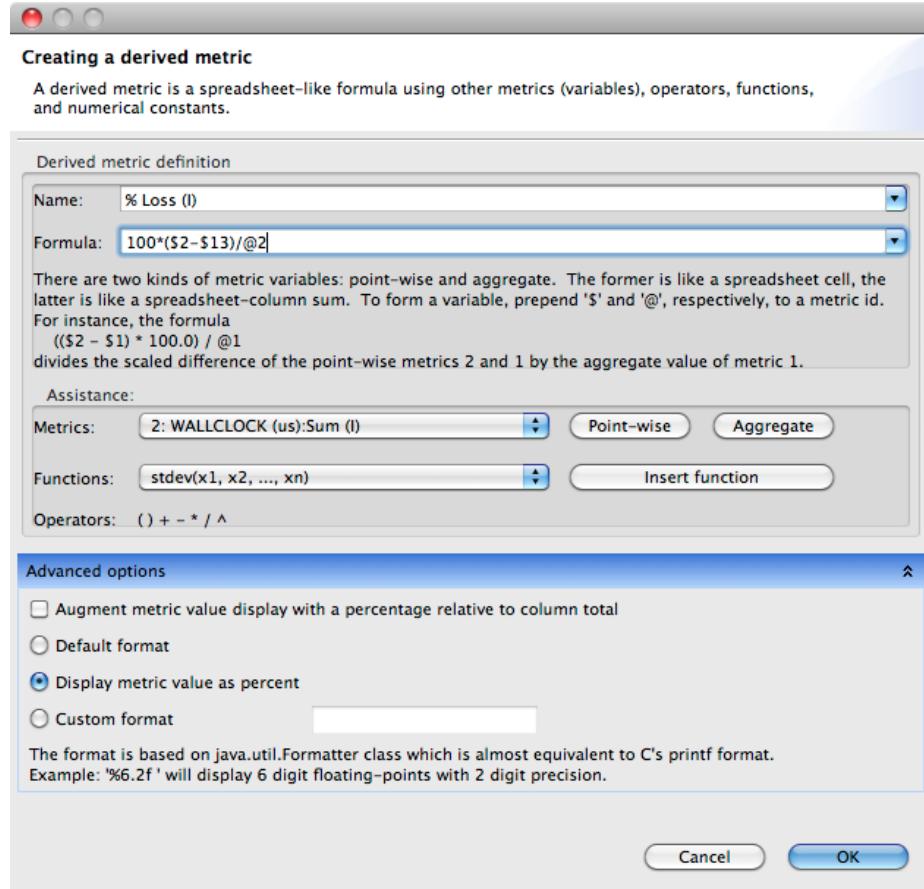


Figure 8.10: Derived metric dialog box

- *Augment metric value display with a percentage relative to column total.* When this box is checked, each scope's derived metric value will be augmented with a percentage value, which for scope s is computed as $100 * (s\text{'s derived metric value}) / (\text{the derived metric value computed by applying the metric formula to the aggregate values of the input metrics for the entire execution})$. Such a computation can lead to nonsensical results for some derived metric formulae. For instance, if the derived metric is computed as a ratio of two other metrics, the aforementioned computation that compares the scope's ratio with the ratio for the entire program won't yield a meaningful result. To avoid a confusing metric display, think before you use this button to annotate a metric with its percent of total.
- *Default format.* This option will display the metric value using scientific notation with three digits of precision, which is the default format.
- *Display metric value as percent.* This option will display the metric value formatted as a percent with two decimal digits. For instance, if the metric has a value 12.3415678, with this option, it will be displayed as 12.34%.

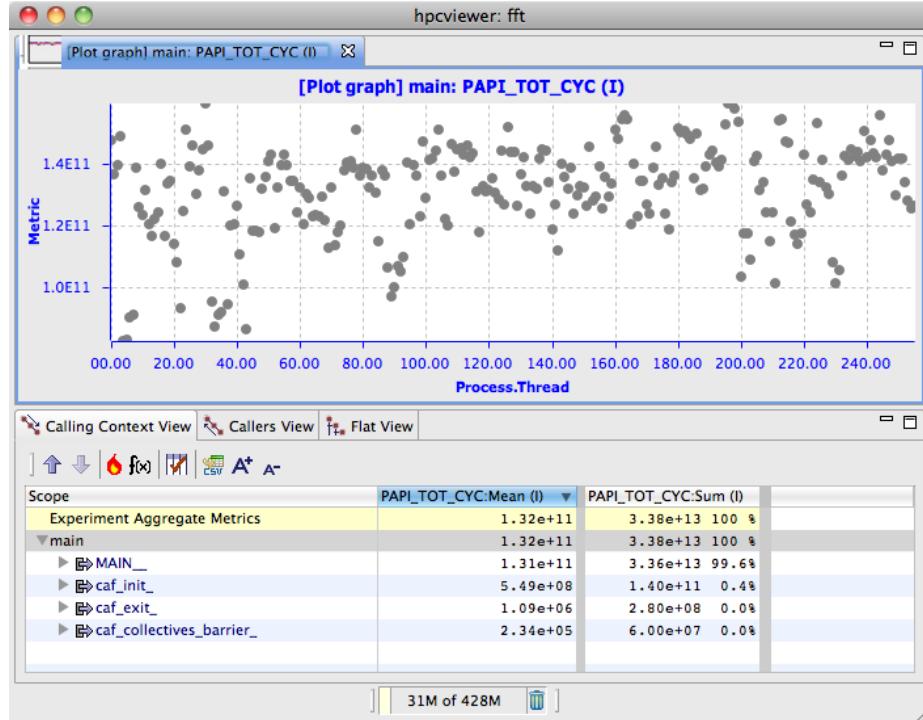


Figure 8.11: Plot graph view of main procedure in a Coarray Fortran application.

- *Custom format.* This option will present the metric value with your customized format. The format is equivalent to Java’s Formatter class, or similar to C’s printf format. For example, the format ”%6.2f” will display six digit floating-points with two digits to the right of the decimal point.

Note that the entered formula and the metric name will be stored automatically. One can then review again the formula (or metric name) by clicking the small triangle of the combo box (marked with a red circle).

8.8 Thread-level Metric Values

8.8.1 Plotting Graphs

HPCTOOLKIT Experiment databases that have been generated by `hpcprof-mpi` (in contrast to `hpcprof`) can be used by `hpcviewer` to plot graphs of thread-level metric values. This is particularly useful for quickly assessing load imbalance *in context* across the several threads or processes of an execution. Figure 8.11 shows `hpcviewer` rendering such a plot. The horizontal axis shows application processes, ordered by MPI rank. The vertical axis shows metric values for each process. Because `hpcviewer` can generate scatter plots for any node in the Top-down View, these graphs are calling-context sensitive.

To create a graph, first select a scope in the Top-down View; in the Figure, the top-level procedure `main` is selected. Then, click the graph button to show the associated sub-menus. At the bottom of the sub-menu is a list of metrics that `hpcviewer` can graph.

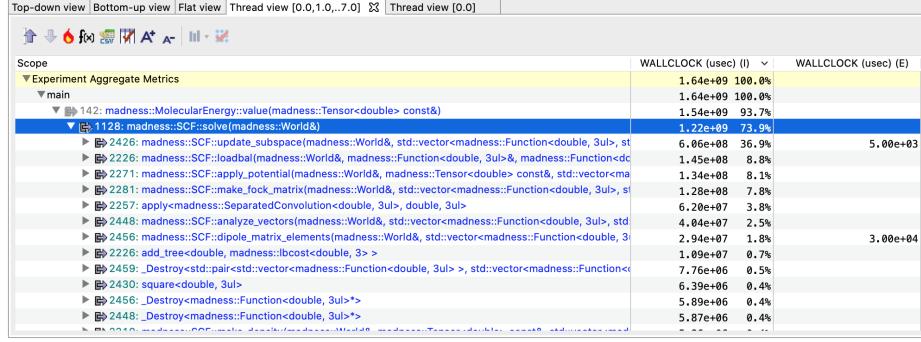


Figure 8.12: Example of a Thread View which display thread-level metrics of a set of threads. The first column is a CCT equivalent to the CCT in the Top-down View, the second and third columns represent the metrics of the selected threads (in this case they are the sum of metrics from threads 0.0, to 7.0)

Each metric contains a sub-menu that lists the three different types of graphs `hpcviewer` can plot:

- **Plot graph.** This standard graph plots metric values by their MPI rank (if available) and thread id (where ids are assigned by thread creation).
- **Sorted plot graph.** This graph plots metric values in ascending order.
- **Histogram graph.** This graph is a histogram of metric values. It divides the range of metric values into a small number of sub-ranges. The graph plots the frequency that a metric value falls into a particular sub-range.

Note that the viewers have the following notation for the ranks:

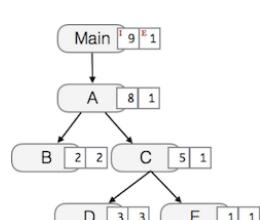
`<process_id> . <thread_id>`

Hence, if the ranks are 0.0, 0.1, ... 31.0, 31.1 it means MPI process 0 has two threads: thread 0 and thread 1 (similarly with MPI process 31).

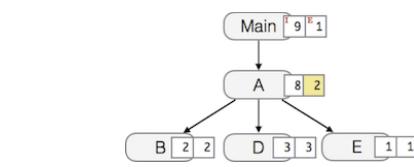
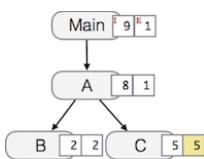
Currently, it is only possible to generate scatter plots for metrics directly collected by `hpcrun`, which excludes derived metrics created within `hpcviewer`.

8.8.2 Thread View

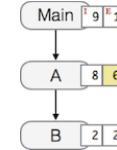
`hpcviewer` also provides a feature to view the metrics of a certain threads (or processes) named Thread View. To select a thread or group of threads, you need to use the thread selection window by clicking button from the calling-context view. On the thread selection window, you need to select the checkbox of the threads of interest. To narrow the list, you can specify the thread name on the filter part of the window. Recall that the format of the thread is “`process_id.thread_id`” (see Section 8.8). Hence, to specify just a main thread (thread zero), you can type ‘`.0`’ on the filter, and the view only list threads 0 (such as 1.0, 2.0, 3.0 ...).



(a) The original CCT tree.



(b) The result of applying *self only* filter on node C. Node C is elided and its children (nodes D and E) are augmented to the parent of node C. The exclusive cost of node C is also augmented to node A.



(d) The result of applying *self and descendants* filter on node C. Nodes C and its descendants are elided, and their exclusive cost is augmented to node A which is the parent of node C.

Figure 8.13: Different results of filtering on node C from Figure 8.13a (the original CCT). Figure 8.13b shows the result of *self only* filter, Figure 8.13c shows the result of *descendants only* filter, and Figure 8.13d shows the result of *self and descendants* filter. Each node is attributed with two boxes on its right. The left box represents the node's inclusive cost, while the right box represents the exclusive cost.

Once threads have been selected, you can click **OK**, and the Thread view (Figure 8.12) will be activated. The tree of the view is the same as the tree from calling context view, with the metrics only from the selected threads. If there are more than one selected threads, the metrics are the sum of the values of the selected threads.

8.9 Filtering Tree Nodes

Occasionally, It is useful to omit uninterested nodes of the tree to enable to focus on important parts. For instance, you may want to hide all nodes associated with OpenMP runtime and just show all nodes and metrics from the application. For this purpose, `hpcviewer` provides *filtering* to elide nodes that match a filter pattern. `hpcviewer` allows users to define multiple filters, and each filter is associated with a glob pattern² and a type. There are three types of filter: “*self only*” to omit matched nodes, “*descendants only*” to exclude only the subtree of the matched nodes, and “*self and descendants*” to remove matched nodes and its descendants.

Self only : This filter is useful to hide intermediary runtime functions such as pthread or OpenMP runtime functions. All nodes that match filter patterns will be removed, and

²A glob pattern specifies which name to be removed by using wildcard characters such as *, ? and +

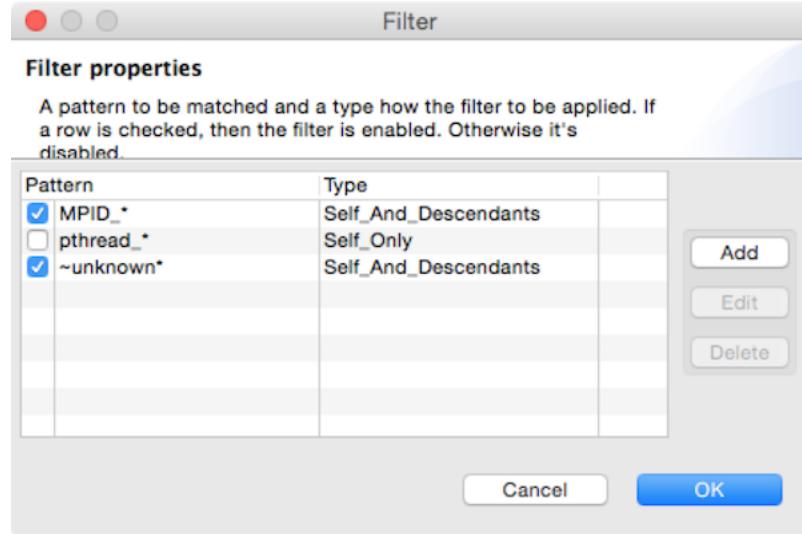


Figure 8.14: The window of filter property.

their children will be augmented to the parent of the elided nodes. The exclusive cost of the elided nodes will be also augmented into the exclusive cost of the parent of the elided nodes. Figure 8.13b shows the result of filtering node C of the CCT from Figure 8.13a. After filtering, node C is elided and its exclusive cost is augmented into the exclusive cost of its parent (node A). The children of node C (nodes D and E) are now the children of node A.

Descendants only : This filter elides only the subtree of the matched node, while the matched node itself is not removed. A common usage of this filter is to exclude any call chains after MPI functions. As shown in Figure 8.13c, filtering node C incurs nodes D and E to be elided and their exclusive cost is augmented to node C.

Self and descendants : This filter elides both the matched node and its subtree. This type is useful to exclude any unnecessary details such as glibc or malloc functions. Figure 8.13d shows that filtering node C will elide the node and its children (nodes D and E). The total of the exclusive cost of the elided nodes is augmented to the exclusive cost of node A.

The filter feature can be accessed by clicking the menu “Filter” and then submenu “Show filter property”, which will then show a Filter property window (Figure 8.14). The window consists of a table of filters, and a group of action buttons: *add* to create a new filter; *edit* to modify a selected filter; and *delete* to remove a set of selected filters.. The table comprises of two columns: the left column is to display a filter’s switch whether the filter is enabled or disabled, and a glob-like filter pattern; and the second column is to show the type of pattern (self only, children only or self and children). If a checkbox is checked, it signifies the filter is enabled; otherwise the filter is disabled.

Cautious is needed when using filter feature since it can change the shape of the tree, thus affects different interpretation of performance analysis. Furthermore, if the filtered

nodes are children of a “fake” procedures (such as `<program root>` and `<thread root>`), the exclusive metrics in Bottom-up view and flat view can be misleading. This occurs since these views do not show “fake” procedures.

8.10 Convenience Features

In this section we describe some features of `hpcviewer` that help improve productivity.

8.10.1 Editor Pane

The editor pane is used to display *a copy* of your program’s source code or HPCTOOLKIT’s performance data in XML format; for this reason, it does not support editing of the pane’s contents. To edit your program, you should use your favorite editor to edit *your* original copy of the source, not the one stored in HPCTOOLKIT’s performance database. Thanks to built-in capabilities in Eclipse, `hpcviewer` supports some useful shortcuts and customization:

- **Find.** To search for a string in the current source pane, `<ctrl>-f` (Linux and Windows) or `<command>-f` (Mac) will bring up a find dialog that enables you to enter the target string.

8.10.2 Metric Pane

For the metric pane, `hpcviewer` has some convenient features:

- **Sorting the metric pane contents by a column’s values.** First, select the column on which you wish to sort. If no triangle appears next to the metric, click again. A downward pointing triangle means that the rows in the metric pane are sorted in descending order according to the column’s value. Additional clicks on the header of the selected column will toggle back and forth between ascending and descending.
- **Changing column width.** To increase or decrease the width of a column, first put the cursor over the right or left border of the column’s header field. The cursor will change into a vertical bar between a left and right arrow. Depress the mouse and drag the column border to the desired position.
- **Changing column order.** If it would be more convenient to have columns displayed in a different order, they can be permuted as you wish. Depress and hold the mouse button over the header of column that you wish to move and drag the column right or left to its new position.
- **Copying selected metrics into clipboard.** In order to copy selected lines of scopes/metrics, one can right click on the metric pane or navigation pane then select the menu **Copy**. The copied metrics can then be pasted into any text editor.
- **Hiding or showing metric columns.** Sometimes, it may be more convenient to suppress the display of metrics that are not of current interest. When there are too many metrics to fit on the screen at once, it is often useful to suppress the display of

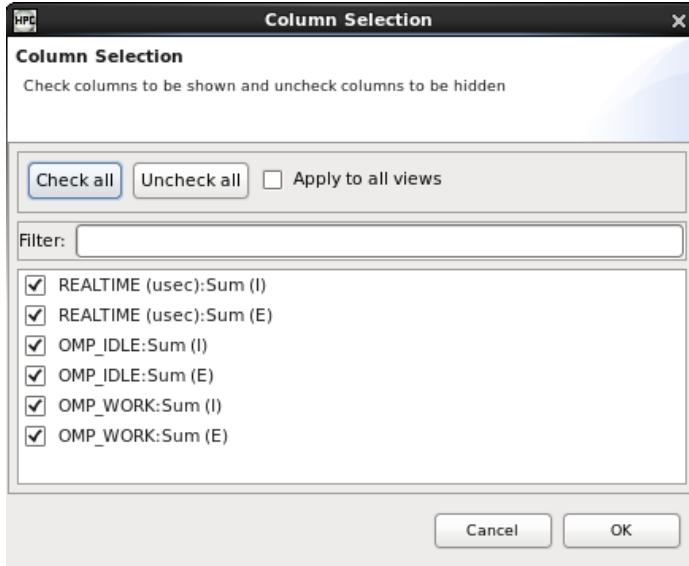


Figure 8.15: Hide/Show columns dialog box

some. The icon above the metric pane will bring up the column selection dialog shown in Figure 8.15.

The dialog box contains a list of metric columns sorted according to their order in HPCTOOLKIT’s performance database for the application. Each metric column is prefixed by a check box to indicate if the metric should be *displayed* (if checked) or *hidden* (unchecked). To display all metric columns, one can click the **Check all** button. A click to **Uncheck all** will hide all the metric columns.

Finally, an option **Apply to all views** will set the configuration into all views when checked. Otherwise, the configuration will be applied only on the current view.

8.11 Menus

`hpcviewer` provides four main menus:

8.11.1 File

This menu includes several menu items for controlling basic viewer operations.

- **New window** Open a new `hpcviewer` window that is independent from the existing one.
- **Open database** Load a performance database into the current `hpcviewer` window replacing the existing opened databases (if any).
- **Add database** Open another database without replacing the existing one. This menu can be used to compare two databases. Currently `hpcviewer` restricts maximum of two database open at a time.

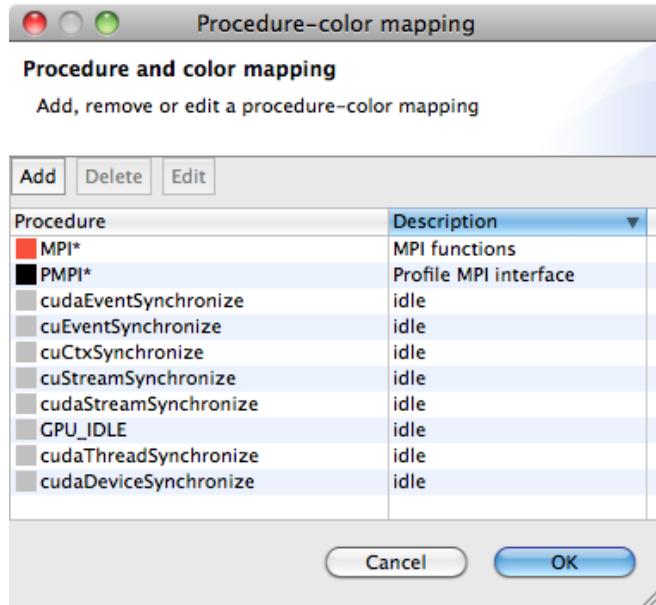


Figure 8.16: Procedure-color mapping dialog box. This window shows that any procedure names that match with "MPI*" pattern are assigned with red, while procedures that match with "PMPI*" pattern are assigned with color black.

- **Close database** Unloading one of more open performance database.
- **Merge databases** Merging two database that are currently in the viewer. At the moment `hpcviewer` doesn't support storing a merged database into a file.
 - **Top-down tree** Merging the top-down tree of the databases.
 - **Flat tree** Merging the flat (static) tree of the databases.
- **Preferences** Display the settings dialog box.
- **Exit** Quit the `hpcviewer` application.

8.11.2 Filter

This menu only contains one submenu:

- **Filter CCT nodes** (*Profile view only*) Open a filter property window which lists a set of filters and its properties (Section 8.9).
- **Filter ranks** (*Trace view only*) Open a window for selecting which nodes will be hidden in the tree. Currently filtering CCT nodes only affect the Profile view, and doesn't affect the Trace view.

8.11.3 View

This menu is only visible if at least one database is loaded. All actions in this menu are intended primarily for tool developer use. By default, the menu is hidden. Once you open a database, the menu is then shown.

- **Show metrics** (*Profile view only*) Display a list of metrics in a window. From this window, you can modify the name of the metric. For derived metrics, this also allows to modify the formula as well as the format.
- **Show procedure-color mapping** (*Trace view only*) Open a window which shows customized mapping between a procedure pattern and a color (Figure 8.16). Trace view allows users to customize assignment of a pattern of procedure names with a specific color.
- **Debug** (*if the debug mode is enabled*)
 - **Show database raw's XML** Enable one to request display of HPCTOOLKIT's raw XML representation for performance data.

8.11.4 Help

This menu displays information about the viewer. The menu contains two items:

- **About.** Displays brief information about the viewer, including used plug-ins and error log.

8.12 Limitations

Some important `hpcviewer` limitations are listed below:

- **No support for ARM platforms.** `hpcviewer` is not currently available on ARM platforms because the Eclipse Rich Client Platform does not yet export binaries for ARM.
- **Limited number of metric columns.** With a large number of metric columns, `hpcviewer`'s response time may become sluggish as this requires a large amount of memory.

Chapter 9

Monitoring MPI Applications

HPCTOOLKIT's measurement subsystem can measure each process and thread in an execution of an MPI program. HPCTOOLKIT can be used with pure MPI programs as well as hybrid programs that use multithreading, e.g. OpenMP or Pthreads, within MPI processes.

HPCTOOLKIT supports C, C++ and Fortran MPI programs. It has been successfully tested with MPICH, MVAPICH and OpenMPI and should work with almost all MPI implementations.

9.1 Running and Analyzing MPI Programs

Q: How do I launch an MPI program with `hpcrun`?

A: For a dynamically linked application binary `app`, use a command line similar to the following example:

```
<mpi-launcher> hpcrun -e <event>:<period> ... app [app-arguments]
```

Observe that the MPI launcher (`mpirun`, `mpiexec`, etc.) is used to launch `hpcrun`, which is then used to launch the application program.

Q: How do I compile and run a statically linked MPI program?

A: On systems such IBM's Blue Gene/Q microkernel that are designed to run statically linked binaries, use `hpclink` to build a statically linked version of your application that includes HPCTOOLKIT's monitoring library. For example, to link your application binary `app`:

```
hpclink <linker> -o app <linker-arguments>
```

Then, set the `HPCRUN_EVENT_LIST` environment variable in the launch script before running the application:

```
export HPCRUN_EVENT_LIST="CYCLES@f200"
<mpi-launcher> app [app-arguments]
```

See the Chapter 10 for more information.

Q: What files does hpcrun produce for an MPI program?

A: In this example, `s3d_f90.x` is the Fortran S3D program compiled with OpenMPI and run with the command line

```
mpiexec -n 4 hpcrun -e PAPI_TOT_CYC:2500000 ./s3d_f90.x
```

This produced 12 files in the following abbreviated `ls` listing:

```
krentel 1889240 Feb 18 s3d_f90.x-000000-000-72815673-21063.hpcrun
krentel    9848 Feb 18 s3d_f90.x-000000-001-72815673-21063.hpcrun
krentel 1914680 Feb 18 s3d_f90.x-000001-000-72815673-21064.hpcrun
krentel    9848 Feb 18 s3d_f90.x-000001-001-72815673-21064.hpcrun
krentel 1908030 Feb 18 s3d_f90.x-000002-000-72815673-21065.hpcrun
krentel    7974 Feb 18 s3d_f90.x-000002-001-72815673-21065.hpcrun
krentel 1912220 Feb 18 s3d_f90.x-000003-000-72815673-21066.hpcrun
krentel    9848 Feb 18 s3d_f90.x-000003-001-72815673-21066.hpcrun
krentel 147635 Feb 18 s3d_f90.x-72815673-21063.log
krentel 142777 Feb 18 s3d_f90.x-72815673-21064.log
krentel 161266 Feb 18 s3d_f90.x-72815673-21065.log
krentel 143335 Feb 18 s3d_f90.x-72815673-21066.log
```

Here, there are four processes and two threads per process. Looking at the file names, `s3d_f90.x` is the name of the program binary, 000000-000 through 000003-001 are the MPI rank and thread numbers, and 21063 through 21066 are the process IDs.

We see from the file sizes that OpenMPI is spawning one helper thread per process. Technically, the smaller `.hpcrun` files imply only a smaller calling-context tree (CCT), not necessarily fewer samples. But in this case, the helper threads are not doing much work.

Q: Do I need to include anything special in the source code?

A: Just one thing. Early in the program, preferably right after `MPI_Init()`, the program should call `MPI_Comm_rank()` with communicator `MPI_COMM_WORLD`. Nearly all MPI programs already do this, so this is rarely a problem. For example, in C, the program might begin with:

```
int main(int argc, char **argv)
{
    int size, rank;

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    ...
}
```

Note: The first call to `MPI_Comm_rank()` should use `MPI_COMM_WORLD`. This sets the process's MPI rank in the eyes of `hpcrun`. Other communicators are allowed, but the first call should use `MPI_COMM_WORLD`.

Also, the call to `MPI_Comm_rank()` should be unconditional, that is all processes should make this call. Actually, the call to `MPI_Comm_size()` is not necessary (for `hpcrun`), although most MPI programs normally call both `MPI_Comm_size()` and `MPI_Comm_rank()`.

Q: What MPI implementations are supported?

A: Although the matrix of all possible MPI variants, versions, compilers, architectures and systems is very large, HPCTOOLKIT has been tested successfully with MPICH, MVAPICH and OpenMPI and should work with most MPI implementations.

Q: What languages are supported?

A: C, C++ and Fortran are supported.

9.2 Building and Installing HPCToolkit

Q: Do I need to compile HPCToolkit with any special options for MPI support?

A: No, HPCTOOLKIT is designed to work with multiple MPI implementations at the same time. That is, you don't need to provide an `mpi.h` include path, and you don't need to compile multiple versions of HPCTOOLKIT, one for each MPI implementation.

The technically-minded reader will note that each MPI implementation uses a different value for `MPI_COMM_WORLD` and may wonder how this is possible. `hpcrun` (actually `libmonitor`) waits for the application to call `MPI_Comm_rank()` and uses the same communicator value that the application uses. This is why we need the application to call `MPI_Comm_rank()` with communicator `MPI_COMM_WORLD`.

Chapter 10

Monitoring Statically Linked Applications

On modern Linux systems, dynamically linked executables are the default. With dynamically linked executables, HPCTOOLKIT’s `hpcrun` script uses library preloading to inject HPCTOOLKIT’s monitoring code into an application’s address space. However, in some cases, statically-linked executables are necessary or desirable.

- One might prefer a statically linked executable because they are generally faster if the executable spends a significant amount of time calling functions in libraries.
- On Blue Gene and Cray supercomputers, statically-linked executables are preferred.

For statically linked executables, preloading HPCTOOLKIT’s monitoring code into an application’s address space at program launch is not an option. Instead, monitoring code must be added at link time; HPCTOOLKIT’s `hpclink` script is used for this purpose.

10.1 Linking with `hpclink`

Adding HPCTOOLKIT’s monitoring code into a statically linked application is easy. This does not require any source-code modifications, but it does involve a small change to your build procedure. You continue to compile all of your object (`.o`) files exactly as before, but you will need to modify your final link step to use `hpclink` to add HPCTOOLKIT’s monitoring code to your executable.

In your build scripts, locate the last step in the build, namely, the command that produces the final statically linked binary. Edit that command line to add the `hpclink` command at the front.

For example, suppose that the name of your application binary is `app` and the last step in your `Makefile` links various object files and libraries as follows into a statically linked executable:

```
mpicc -o app -static file.o ... -l<lib> ...
```

To build a version of your executable with HPCTOOLKIT’s monitoring code linked in, you would use the following command line:

```
hpmlink mpicc -o app -static file.o ... -l<lib> ...
```

In practice, you may want to edit your `Makefile` to always build two versions of your program, perhaps naming them `app` and `app.hpc`.

10.2 Running a Statically Linked Binary

For dynamically linked executables, the `hpcrun` script sets environment variables to pass information to the HPCTOOLKIT monitoring library. On standard Linux systems, statically linked `hpmlink`-ed executables can still be launched with `hpcrun`.

On Cray and Blue Gene systems, the `hpcrun` script is not applicable because of differences in application launch procedures. On these systems, you will need to use the `HPCRUN_EVENT_LIST` environment variable to pass a list of events to HPCTOOLKIT's monitoring code, which was linked into your executable using `hpmlink`. Typically, you would set `HPCRUN_EVENT_LIST` in your launch script.

The `HPCRUN_EVENT_LIST` environment variable should be set to a space-separated list of `EVENT@COUNT` pairs. For example, in a PBS script for a Cray system, you might write the following in Bourne shell or bash syntax:

```
#!/bin/sh
#PBS -l size=64
#PBS -l walltime=01:00:00
cd $PBS_O_WORKDIR
export HPCRUN_EVENT_LIST="CYCLES@f200 PERF_COUNT_HW_CACHE_MISSES@f200"
aprun -n 64 ./app arg ...
```

Using the Cobalt job launcher on Argonne National Laboratory's Blue Gene/Q system, you would use the `--env` option to pass environment variables. For example, you might submit a job with:

```
qsub -t 60 -n 64 --env HPCRUN_EVENT_LIST="WALLCLOCK@1000" \
      /path/to/app <app arguments> ...
```

To collect sample traces of an execution of a statically linked binary (for visualization with Trace view), one needs to set the environment variable `HPCRUN_TRACE=1` in the execution environment.

10.3 Troubleshooting

With some compilers you need to disable interprocedural optimization to use `hpmlink`. To instrument your statically linked executable at link time, `hpmlink` uses the `ld` option `--wrap` (see the `ld(1)` man page) to interpose monitoring code between your application and various process, thread, and signal control operations, e.g., `fork`, `pthread_create`, and `sigprocmask` to name a few. For some compilers, e.g., IBM's XL compilers and Pathscale's compilers, interprocedural optimization interferes with the `--wrap` option and prevents `hpmlink` from working properly. If this is the case, `hpmlink` will emit error messages and

fail. If you want to use `hpmlink` with such compilers, sadly, you must turn off interprocedural optimization.

Note that interprocedural optimization may not be explicitly enabled during your compiles; it might be implicitly enabled when using a compiler optimization option such as `-fast`. In cases such as this, you can often specify `-fast` along with an option such as `-no-ipa`; this option combination will provide the benefit of all of `-fast`'s optimizations *except* interprocedural optimization.

Chapter 11

Known Issues

11.1 A confusing label for GPU theoretical occupancy

Affected architectures: NVIDIA GPUs

Description: When analyzing a GPU-accelerated application that employs NVIDIA GPUs, HPCToolkit estimates percent GPU theoretical occupancy as the ratio of active GPU threads divided by the maximum number of GPU threads available. In multi-threaded or multi-rank programs, HPCToolkit reports GPU theoretical occupancy with the label

Sum over rank/thread of exclusive 'GPU kernel: theoretical occupancy
(FGP_ACT / FGP_MAX)'

rather than its correct label

GPU kernel: theoretical occupancy (FGP_ACT / FGP_MAX)

The metric is computed correctly by summing the fine-grain parallelism used in each kernel launch across all threads and ranks and dividing it by the sum of the maximum fine-grain parallelism available to each kernel launch across all threads and ranks, and presenting the value as a percent.

Explanation: This metric is unlike others computed by HPCToolkit. Rather than being computed by `hpcprof`, it is computed by having `hpcviewer` interpret a formula.

Workaround: Pay attention to the metric value, which is computed correctly and ignore its awkward label.

Development Plan: Add additional support to `hpcrun` and `hpcprof` to understand how derived metrics are computed and avoid spoiling their labels.

11.2 Deadlock when using Darshan

Affected architectures: x86_64 and ARM

Description: Darshan is a library for monitoring POSIX I/O. When using asynchronous sampling on the CPU to monitor a program that is being monitored with Darshan, your program may deadlock.

Explanation: Darshan hijacks calls to `open`. HPCToolkit uses the `libunwind` library. Under certain circumstances, `libunwind` uses `open` to inspect an application's executable or one of the shared libraries it uses to look for unwinding information recorded by the compiler. The following sequence of actions leads to a problem:

1. A user application calls `malloc` and acquires a mutex lock on an allocator data structure.
2. HPCToolkit's signal handler is invoked to record an asynchronous sample.
3. `libunwind` is invoked to obtain the calling context for the sample.
4. `libunwind` calls `open` to look for compiler-based unwind information.
5. A Darshan wrapper for `open` executes in HPCToolkit's signal handler.
6. The Darshan wrapper for `open` may try to allocate data to record statistics for the application's calls to `open`, deadlocking because a non-reentrant allocator lock is already held by this thread.

Workaround: Unload the Darshan module before compiling a statically-linked application or running a dynamically-linked application.

Development Plan: Ensure that `libunwind`'s calls to `open` are never intercepted by Darshan.

11.3 Deadlock when using UCX

Affected architectures: x86_64 and ARM

Description: OpenUCX (<https://www.openucx.org>) is a communication layer developed at ORNL. When using asynchronous sampling on the CPU to monitor a program that uses a communication layer based on UCX, your program may deadlock. We have observed this behavior with a program that used a version of OpenMPI that employed UCX as a communication substrate.

Explanation: UCX hijacks calls to `mmap` and `munmap`. HPCToolkit uses the `libunwind` library. Under certain circumstances, `libunwind` uses `mmap` and `munmap` to map a section of your application's executable or one of the shared libraries it uses to obtain unwinding information recorded by the compiler. The following sequence of actions leads to a problem:

1. UCX is interrupted with an asynchronous sample while holding a lock in its implementation.

2. HPCToolkit's signal handler is invoked to record the sample.
3. `libunwind` is invoked to obtain the calling context for the sample.
4. `libunwind` calls `mmap` and `munmap` when inspecting compiler-based unwind information.
5. UCX wrappers for `mmap` and `munmap` execute in HPCToolkit's signal handler.
6. A UCX wrapper for `mmap` or `munmap` may try to acquire a non-reentrant lock that is already held by the UCX code that was interrupted by the sample trigger.

Workaround: Use a communication library that doesn't employ UCX.

Development Plan: Ensure that `libunwind`'s calls to `mmap` and `munmap` are never intercepted by UCX.

Chapter 12

FAQ and Troubleshooting

To measure an application’s performance with HPCTOOLKIT, one must add HPCTOOLKIT’s measurement subsystem to an application’s address space.

- For a statically-linked binary, one adds HPCTOOLKIT’s measurement subsystem directly into the binary by prefixing your link command with HPCTOOLKIT’s `hpclink` command.
- For a dynamically-linked binary, launching your application with HPCTOOLKIT’s `hpcrun` command pre-loads HPCTOOLKIT’s measurement subsystem into your application’s address space before the application begins to execute.

In this Chapter, for convenience, we refer to HPCToolkit’s measurement system simply as `hpcrun` since the measurement subsystem is most commonly used with dynamically-linked binaries. From the context, it should be clear enough whether we are talking about HPCTOOLKIT’s measurement subsystem or the `hpcrun` command itself.

12.1 Why do I see partial unwinds?

Under certain circumstances, HPCToolkit can’t fully unwind the call stack to determine the full calling context where a sample event occurred. Most often, this occurs when `hpcrun` tries to unwind through functions in a shared library or executable that has not been compiled with `-g` as one of its options. The `-g` compiler flag can be used in addition to optimization flags. On Power and `x86_64` processors, `hpcrun` can often compensate the lack of the availability of unwind recipes in a binary by using binary analysis to compute recipes itself. However, since `hpcrun` lacks binary analysis capabilities for ARM processors, there is a higher likelihood that the lack of a `-g` compiler option for an executable or a shared library will lead to partial unwinds.

There are two ways to deal with this problem:

- Analyze the execution using information from partial unwinds. Often knowing several levels of calling context is enough for analysis without full calling context for sample events.
- Recompile the binary or shared library causing the problem and add `-g` to the list of its compiler options.

12.2 How do I choose `hpcrun` sampling periods?

When using sample sources for hardware counter and software counter events provided by Linux `perf_events`, we recommend that you use frequency-based sampling. The default frequency is 300 samples/second.

Statisticians use samples sizes of approximately 3500 to make accurate projections about the voting preferences of millions of people. In an analogous way, rather than collect unnecessary large amounts of performance information, sampling-based performance measurement collects “just enough” representative performance data. You can control `hpcrun`’s sampling periods to collect “just enough” representative data even for very long executions and, to a lesser degree, for very short executions.

For reasonable accuracy ($\pm 5\%$), there should be at least 20 samples in each context that is important with respect to performance. Since unimportant contexts are irrelevant to performance, as long as this condition is met (and as long as samples are not correlated, etc.), HPCTOOLKIT’s performance data should be accurate.

We typically recommend targeting a frequency of hundreds of samples per second. For very short runs, you may need to collect thousands of samples per second to record an adequate number of samples. For long runs, tens of samples per second may suffice for performance diagnosis.

Choosing sampling periods for some events, such as Linux timers, cycles and instructions, is easy given a target sampling frequency. Choosing sampling periods for other events such as cache misses is harder. In principle, an architectural expert can easily derive reasonable sampling periods by working backwards from (a) a maximum target sampling frequency and (b) hardware resource saturation points. In practice, this may require some experimentation.

See also the `hpcrun` man page.

12.3 `hpcrun` incurs high overhead! Why?

For reasonable sampling periods, we expect `hpcrun`’s overhead percentage to be in the low single digits, e.g., less than 5%. The most common causes for unusually high overhead are the following:

- Your sampling frequency is too high. Recall that the goal is to obtain a representative set of performance data. For this, we typically recommend targeting a frequency of hundreds of samples per second. For very short runs, you may need to try thousands of samples per second. For very long runs, tens of samples per second can be quite reasonable. See also Section 12.2.
- `hpcrun` has a problem unwinding. This causes overhead in two forms. First, `hpcrun` will resort to more expensive unwind heuristics and possibly have to recover from self-generated segmentation faults. Second, when these exceptional behaviors occur, `hpcrun` writes some information to a log file. In the context of a parallel application and overloaded parallel file system, this can perturb the execution significantly. To diagnose this, execute the following command and look for “Errant Samples”:

```
hpcsummary --all <hpctoolkit-measurements>
```

Note: The `hpcsummary` script is no longer included in the `bin` directory of an HPCTOOLKIT installation; it is a developer script that can be found in the `libexec/hpctoolkit` directory. Let us know if you encounter significant problems with bad unwinds.

- You have very long call paths where long is in the hundreds or thousands. On x86-based architectures, try additionally using `hpcrun`'s RETCNT event. This has two effects: It causes `hpcrun` to collect function return counts and to memoize common unwind prefixes between samples.
- Currently, on very large runs the process of writing profile data can take a long time. However, because this occurs after the application has finished executing, it is relatively benign overhead. (We plan to address this issue in a future release.)

12.4 Fail to run `hpcviewer`: executable launcher was unable to locate its companion shared library

Although this error mostly incurs on Windows platform, but it can happen in other environment. The cause of this issue is that the permission of one of Eclipse launcher library (`org.eclipse.equinox.launcher.*`) is too restricted. To fix this, set the permission of the library to 0755, and launch again the viewer.

12.5 Mac only: `hpcviewer` runs on Java X instead of “Java 8”

If your system has multiple versions of Java and Java 8 is not the newer version, you need to set Java 8 as the default JVM:

1. Leave all JDKs at their default location (usually under `/Library/Java/JavaVirtualMachines`). The system will pick the highest version by default.
2. To exclude a JDK from being picked by default, rename `Contents/Info.plist` file to other name like `Info.plist.disabled`. That JDK can still be used when `$JAVA_HOME` points to it, or explicitly referenced in a script or configuration. It will simply be ignored by system's `java` command.

12.6 When executing `hpcviewer`, it complains cannot create “Java Virtual Machine”

If you encounter this problem, we recommend that you edit the `hpcviewer.ini` file which is located in HPCToolkit installation directory to reduce the Java heap size. By default, the content of the file on `ppc64le` is as follows:

```
-consoleLog  
-clearPersistedState  
-vmargs  
-Dosgi.locking=none  
-Dosgi.requiredJavaVersion=1.8  
-Xmx2048m
```

You can decrease the maximum size of the Java heap from 2048MB to 1GB by changing the `Xmx` specification in the `hpcviewer.ini` file as follows:

```
-Xmx1024m
```

12.7 hpcviewer fails to launch due to java.lang.NoSuchMethodError exception.

The root cause of the error is due to a mix of old new hpcviewer binaries. To solve this problem, you need to remove your hpcviewer workspace (usually in `$HOME/.hpctoolkit/hpcviewer` directory, and run `hpcviewer` again.

12.8 hpcviewer fails due to java.lang.OutOfMemoryError exception.

If you see this error, the memory footprint that `hpcviewer` needs to store and the metrics for your measured program execution exceeds the maximum size for the Java heap specified at program launch. On Linux, `hpcviewer` accepts a command-line option `--java-heap` that enables you to specify a larger non-default value for the maximum size of the Java heap. Run `hpcviewer --help` for the details of how to use this option.

12.9 hpcviewer writes a long list of Java error messages to the terminal!

The Eclipse Java framework that serves as the foundation for `hpcviewer` can be somewhat temperamental. If the persistent state maintained by Eclipse for `hpcviewer` gets corrupted, `hpcviewer` may spew a list of errors deep within call chains of the Eclipse framework. Below are a few suggestions that may fix the problem:

On Linux, try removing your `hpcviewer` Eclipse workspace with default location:
`$HOME/.hpctoolkit/hpcviewer`
and run `hpcviewer` again.

On MacOS, persistent state is currently stored within Mac app. If the Eclipse persistent state gets corrupted, one can't simply clear the workspace because some initial persistent state is needed for Eclipse to function properly. For MacOS, the thing to try is downloading a fresh copy of `hpcviewer` and running the freshly downloaded copy.

If one of the aforementioned suggestions doesn't fix the problem, report a bug.

12.10 hpcviewer attributes performance information only to functions and not to source code loops and lines! Why?

Most likely, your application’s binary either lacks debugging information or is stripped. A binary’s (optional) debugging information includes a line map that is used by profilers and debuggers to map object code to source code. HPCTOOLKIT can profile binaries without debugging information, but without such debugging information it can only map performance information (at best) to functions instead of source code loops and lines.

For this reason, we recommend that you always compile your production applications with optimization *and* with debugging information. The options for doing this vary by compiler. We suggest the following options:

- GNU compilers (`gcc`, `g++`, `gfortran`): `-g`
- Intel compilers (`icc`, `icpc`, `ifort`): `-g -debug inline_debug_info`
- Pathscale compilers (`pathcc`, `pathCC`, `pathf95`): `-g1`
- PGI compilers (`pgcc`, `pgCC`, `pgf95`): `-gopt`.

We generally recommend adding optimization options *after* debugging options — e.g., ‘`-g -O2`’ — to minimize any potential effects of adding debugging information.¹ Also, be careful not to strip the binary as that would remove the debugging information. (Adding debugging information to a binary does not make a program run slower; likewise, stripping a binary does not make a program run faster.)

Please note that at high optimization levels, a compiler may make significant program transformations that do not cleanly map to line numbers in the original source code. Even so, the performance attribution is usually very informative.

12.11 hpcviewer hangs trying to open a large database! Why?

The most likely problem is that the Java virtual machine is low on memory and thrashing. The memory footprint that `hpcviewer` needs to store and the metrics for your measured program execution is likely near the maximum size for the Java heap specified at program launch.

On Linux, `hpcviewer` accepts a command-line option `--java-heap` that enables you to specify a larger non-default value for the maximum size of the Java heap. Run `hpcviewer --help` for the details of how to use this option.

¹In general, debugging information is compatible with compiler optimization. However, in a few cases, compiling with debugging information will disable some optimization. We recommend placing optimization options *after* debugging options because compilers usually resolve option incompatibilities in favor of the last option.

12.12 hpcviewer runs glacially slowly! Why?

There are three likely reasons why `hpcviewer` might run slowly. First, you may be running `hpcviewer` on a remote system with low bandwidth, high latency or an otherwise unsatisfactory network connection to your desktop. If any of these conditions are true, `hpcviewer`'s otherwise snappy GUI can become sluggish if not downright unresponsive. The solution is to install `hpcviewer` on your local system, copy the database onto your local system, and run `hpcviewer` locally. We almost always run `hpcviewer` on our local desktops or laptops for this reason.

Second, HPCTOOLKIT's database may contain too many metrics. A common reason why this can occur is if rather than analyzing an HPCTOOLKIT measurements directory, you chose to directly specify a large collection of `.hpcrun` files directly on the command line. When you do this, it treats each `.hpcrun` file as a separate experiment and gives you separate metrics for each.

You can check the number of columns in your database by running

```
grep -e "<Metric" experiment.xml | wc -l
```

If that command yields a number greater than 50 or so, `hpcviewer` is likely slow because you are working with so many columns of metrics. In this case, you might consider using `hpcprof-mpi` or run `hpcprof` to analyze the measurements directory as a whole, so you get one set of metrics, or you can build a database based on fewer separate `.hpcrun` files on the command line to reduce the number of metrics overall.

Third, HPCTOOLKIT's database may be very large, which can cause the Java virtual machine to run short on memory and thrash. The memory footprint that `hpcviewer` needs to store and the metrics for your measured program execution is likely near the maximum size for the Java heap specified at program launch.

On Linux, `hpcviewer` accepts a command-line option `--java-heap` that enables you to specify a larger non-default value for the maximum size of the Java heap. Run `hpcviewer --help` for the details of how to use this option.

12.13 hpcviewer does not show my source code! Why?

Assuming you compiled your application with debugging information (see Issue 12.10), the most common reason that `hpcviewer` does not show source code is that `hpcprof/mpi` could not find it and therefore could not copy it into the HPCTOOLKIT performance database.

12.13.1 Follow ‘Best Practices’

When running `hpcprof/mpi`, we recommend using an `-I/--include` option to specify a search directory for each distinct top-level source directory (or build directory, if it is separate from the source directory). Assume the paths to your top-level source directories are `<dir1>` through `<dirN>`. Then, pass the the following options to `hpcprof/mpi`:

```
-I <dir1>/+ -I <dir2>/+ ... -I <dirN>/+
```

These options instruct `hpcprof/mpi` to search for source files that live within any of the source directories `<dir1>` through `<dirN>`. Each directory argument can be either absolute or relative to the current working directory.

It will be instructive to unpack the rationale behind this recommendation. `hpcprof/mpi` obtains source file names from your application binary's debugging information. These source file paths may be either absolute or relative. Without any `-I/--include` options, `hpcprof/mpi` can find source files that either (1) have absolute paths (and that still exist on the file system) or (2) are relative to the current working directory. However, because the nature of these paths depends on your compiler and the way you built your application, it is not wise to depend on either of these default path resolution techniques. For this reason, we always recommend supplying at least one `-I/--include` option.

There are two basic forms in which the search directory can be specified: non-recursive and recursive. In most cases, the most useful form is the recursive search directory, which means that the directory should be searched *along with all of its descendants*. A non-recursive search directory `dir` is simply specified as `dir`. A recursive search directory `dir` is specified as the base search directory followed by the special suffix '`/+`': `dir/+`. The paths above use the recursive form.

12.13.2 Additional Background

`hpcprof/mpi` obtains source file names from your application binary's debugging information. If debugging information is unavailable, such as is often the case for system or math libraries, then source files are unknown. Two things immediately follow from this. First, in most normal situations, there will always be some functions for which source code cannot be found, such as those within system libraries. Second, to ensure that `hpcprof/mpi` has file names for which to search, make sure as much of your application as possible (including libraries) contains debugging information.

If debugging information is available, source files can come in two forms: absolute and relative. `hpcprof/mpi` can find source files under the following conditions:

- If a source file path is absolute and the source file can be found on the file system, then `hpcprof/mpi` will find it.
- If a source file path is relative, `hpcprof/mpi` can only find it if the source file can be found from the current working directory or within a search directory (specified with the `-I/--include` option).
- Finally, if a source file path is absolute and cannot be found by its absolute path, `hpcprof/mpi` uses a special search mode. Let the source file path be p/f . If the path's base file name f is found within a search directory, then that is considered a match. This special search mode accommodates common complexities such as: (1) source file paths that are relative not to your source code tree but to the directory where the source was compiled; (2) source file paths to source code that is later moved; and (3) source file paths that are relative to file system that is no longer mounted.

Note that given a source file path p/f (where p may be relative or absolute), it may be the case that there are multiple instances of a file's base name f within one search directory,

e.g., p_1/f through p_n/f , where p_i refers to the i^{th} path to f . Similarly, with multiple search-directory arguments, f may exist within more than one search directory. If this is the case, the source file p/f is resolved to the first instance p'/f such that p' best corresponds to p , where instances are ordered by the order of search directories on the command line.

For any functions whose source code is not found (such as functions within system libraries), `hpcviewer` will generate a synopsis that shows the presence of the function and its line extents (if known).

12.14 `hpcviewer`'s reported line numbers do not exactly correspond to what I see in my source code! Why?

To use a cliché, “garbage in, garbage out”. HPCTOOLKIT depends on information recorded in the symbol table by the compiler. Line numbers for procedures and loops are inferred by looking at the symbol table information recorded for machine instructions identified as being inside the procedure or loop.

For procedures, often no machine instructions are associated with a procedure’s declarations. Thus, the first line in the procedure that has an associated machine instruction is the first line of executable code.

Inlined functions may occasionally lead to confusing data for a procedure. Machine instructions mapped to source lines from the inlined function appear in the context of other functions. While `hpcprof`’s methods for handling inlined functions are good, some codes can confuse the system.

For loops, the process of identifying what source lines are in a loop is similar to the procedure process: what source lines map to machine instructions inside a loop defined by a backward branch to a loop head. Sometimes compilers do not properly record the line number mapping.

When the compiler line mapping information is wrong, there is little you can do about it other than to ignore its imperfections, or hand-edit the XML program structure file produced by `hpcstruct`. This technique is used only when truly desperate.

12.15 `hpcviewer` claims that there are several calls to a function within a particular source code scope, but my source code only has one! Why?

In the course of code optimization, compilers often replicate code blocks. For instance, as it generates code, a compiler may peel iterations from a loop or split the iteration space of a loop into two or more loops. In such cases, one call in the source code may be transformed into multiple distinct calls that reside at different code addresses in the executable.

When analyzing applications at the binary level, it is difficult to determine whether two distinct calls to the same function that appear in the machine code were derived from the same call in the source code. Even if both calls map to the same source line, it may be wrong to coalesce them; the source code might contain multiple calls to the same function on the same line. By design, HPCTOOLKIT does not attempt to coalesce distinct calls to the same function because it might be incorrect to do so; instead, it independently reports each

call site that appears in the machine code. If the compiler duplicated calls as it replicated code during optimization, multiple call sites may be reported by `hpcviewer` when only one appeared in the source code.

12.16 Trace view shows lots of white space on the left. Why?

At startup, Trace view renders traces for the time interval between the minimum and maximum times recorded for any process or thread in the execution. The minimum time for each process or thread is recorded when its trace file is opened as HPCToolkit’s monitoring facilities are initialized at the beginning of its execution. The maximum time for a process or thread is recorded when the process or thread is finalized and its trace file is closed. When an application uses the `hpctoolkit_start` and `hpctoolkit_stop` primitives, the minimum and maximum time recorded for a process/thread are at the beginning and end of its execution, which may be distant from the start/stop interval. This can cause significant white space to appear in Trace view’s display to the left and right of the region (or regions) of interest demarcated in an execution by start/stop calls.

12.17 I get a message about “Unable to find HPCTOOLKIT root directory”

On some systems, you might see a message like this:

```
/path/to/copy/of/hpcrun: Unable to find HPCTOOLKIT root directory.  
Please set HPCTOOLKIT to the install prefix, either in this script,  
or in your environment, and try again.
```

The problem is that the system job launcher copies the `hpcrun` script from its install directory to a launch directory and runs it from there. When the system launcher moves `hpcrun` to a different directory, this breaks `hpcrun`’s method for finding its own install directory. The solution is to add `HPCTOOLKIT` to your environment so that `hpcrun` can find its install directory. See section 5.7 for general notes on environment variables for `hpcrun`. Also, see section 5.8, as this problem occurs on Cray XE and XK systems.

Note: Your system may have a module installed for `hpctoolkit` with the correct settings for `PATH`, `HPCTOOLKIT`, etc. In that case, the easiest solution is to load the `hpctoolkit` module. If there is such a module, Try “`module show hpctoolkit`” to see if it sets `HPCTOOLKIT`.

12.18 Some of my syscalls return EINTR when run under hpcrun

When profiling a threaded program, there are times when it is necessary for `hpcrun` to signal another thread to take some action. When this happens, if the thread receiving the signal is blocked in a syscall, the kernel may return `EINTR` from the syscall. This would happen only in a threaded program and mainly with “slow” syscalls such as `select()`, `poll()` or `sem_wait()`.

12.19 How do I debug HPCToolkit's measurement?

Assume you want to debug HPCTOOLKIT's measurement subsystem when collecting measurements for an application named `app`.

12.19.1 Tracing libmonitor

HPCTOOLKIT's measurement subsystem uses `libmonitor` for process/thread control. To collect a debug trace of `libmonitor`, use either `monitor-run` or `monitor-link`, which are located within:

```
<externals-install>/libmonitor/bin
```

Launch your application as follows:

- Dynamically linked applications:

```
[<mpi-launcher>] monitor-run --debug app [app-arguments]
```

- Statically linked applications:

Link `libmonitor` into `app`:

```
monitor-link <linker> -o app <linker-arguments>
```

Then execute `app` under special environment variables:

```
export MONITOR_DEBUG=1  
[<mpi-launcher>] app [app-arguments]
```

12.19.2 Tracing HPCToolkit's Measurement Subsystem

Broadly speaking, there are two levels at which a user can test `hpcrun`. The first level is tracing `hpcrun`'s application control, that is, running `hpcrun` without an asynchronous sample source. The second level is tracing `hpcrun` with a sample source. The key difference between the two is that the former uses the `--event NONE` or `HPCRUN_EVENT_LIST="NONE"` option (shown below) whereas the latter does not (which enables the default CPUTIME sample source). With this in mind, to collect a debug trace for either of these levels, use commands similar to the following:

- Dynamically linked applications:

```
[<mpi-launcher>] \  
hpcrun --monitor-debug --dynamic-debug ALL --event NONE \  
app [app-arguments]
```

- Statically linked applications:

Link `hpcrun` into `app` (see Section 3.1.2). Then execute `app` under special environment variables:

```

export MONITOR_DEBUG=1
export HPCRUN_EVENT_LIST="NONE"
export HPCRUN_DEBUG_FLAGS="ALL"
[<mpi-launcher>] app [app-arguments]

```

Note that the `*debug*` flags are optional. The `--monitor-debug/MONITOR_DEBUG` flag enables `libmonitor` tracing. The `--dynamic-debug/HPCRUN_DEBUG_FLAGS` flag enables `hpcrun` tracing.

12.19.3 Using a debugger to inspect an execution being monitored by HPCToolkit

If HPCTOOLKIT has trouble monitoring an application, you may find it useful to execute an application being monitored by HPCTOOLKIT under the control of a debugger to observe how HPCTOOLKIT’s measurement subsystem interacts with the application.

HPCTOOLKIT’s measurement subsystem is easiest to debug if you configure and build HPCTOOLKIT by adding the `--enable-develop` option as an argument to `configure` when preparing to build HPCTOOLKIT. (It is not necessary to rebuild HPCTOOLKIT’s `hpctoolkit-externals`.)

One can debug a statically-linked or a dynamically-linked applications being measured by HPCTOOLKIT’s measurement subsystem.

- Dynamically-linked applications. When launching an application with `hpcrun`, add the `--debug` option to `hpcrun`.
- Statically-linked applications. To debug a statically-linked application that has HPCTOOLKIT’s measurement subsystem linked into it, set `HPCRUN_WAIT` in the environment before launching the application, e.g.

```

export HPCRUN_WAIT=1
export HPCRUN_EVENT_LIST="... the metric(s) you want to measure ..."
app [app-arguments]

```

There are two ways to use launch an application with a debugger when using To attach a debugger when monitoring an application using `hpcrun`, add `hpcrun`’s `--debug` option o debug `hpcrun` with a debugger use the following approach.

1. Launch your application. To debug `hpcrun` without controlling sampling signals, launch normally. To debug `hpcrun` with controlled sampling signals, launch as follows:

```
hpcrun --debug --event REALTIME@0 app [app-arguments]
```

or

```

export HPCRUN_WAIT=1
export HPCRUN_EVENT_LIST="REALTIME@0"
app [app-arguments]

```

2. Attach a debugger. The debugger should be spinning in a loop whose exit is conditioned by the `HPCRUN_DEBUGGER_WAIT` variable.
3. Set any desired breakpoints. To send a sampling signal at a particular point, make sure to stop at that point with a *one-time* or *temporary* breakpoint (`tbreak` in GDB).
4. Call `hpcrun_continue()` or set the `HPCRUN_DEBUGGER_WAIT` variable to 0 and continue.
5. To raise a controlled sampling signal, raise a `SIGPROF`, e.g., using GDB's command `signal SIGPROF`.

12.19.4 Using `hpmlink` with `cmake`

When creating a statically-linked executable with `cmake`, it is not obvious how to add `hpmlink` as a prefix to a link command. Unless it is overridden somewhere along the way, the following rule found in `Modules/CMakeCXXInformation.cmake` is used to create the link command line for a C++ executable:

```
if(NOT CMAKE_CXX_LINK_EXECUTABLE)
  set(CMAKE_CXX_LINK_EXECUTABLE
      "<CMAKE_CXX_COMPILER> <FLAGS> <CMAKE_CXX_LINK_FLAGS> <LINK_FLAGS>
      <OBJECTS> -o <TARGET> <LINK_LIBRARIES>")
endif()
```

As the rule shows, by default, the C++ compiler is used to link C++ executables. One way to change this is to override the definition for `CMAKE_CXX_LINK_EXECUTABLE` on the `cmake` command line so that it includes the necessary `hpmlink` prefix, as shown below:

```
cmake srcdir ... \
-DCMAKE_CXX_LINK_EXECUTABLE="hpmlink <CMAKE_CXX_COMPILER> \
<FLAGS> <CMAKE_CXX_LINK_FLAGS> <LINK_FLAGS> <OBJECTS> -o <TARGET> \
<LINK_LIBRARIES>" ...
```

If your project has executables linked with a C or Fortran compiler, you will need analogous redefinitions for `CMAKE_C_LINK_EXECUTABLE` or `CMAKE_Fortran_LINK_EXECUTABLE` as well.

Rather than adding the redefinitions of these linker rules to the `cmake` command line, you may find it more convenient to add definitions of these rules to your `CMakeLists.cmake` file.

Bibliography

- [1] L. Adhianto, S. Banerjee, M. Fagan, M. Krentel, G. Marin, J. Mellor-Crummey, and N. R. Tallent. HPCToolkit: Tools for performance analysis of optimized parallel programs. *Concurrency and Computation: Practice and Experience*, 22(6):685–701, 2010.
- [2] L. Adhianto, J. Mellor-Crummey, and N. R. Tallent. Effectively presenting call path profiles of application performance. In *PSTI 2010: Workshop on Parallel Software Tools and Tool Infrastructures, in conjunction with the 2010 International Conference on Parallel Processing*, 2010.
- [3] Advanced Micro Devices. ROCm Tracer Callback/Activity Library for Performance tracing AMD GPU’s. [Accessed February 27, 2020]. <https://github.com/ROCm-Developer-Tools/roctracer>.
- [4] C. Coarfa, J. Mellor-Crummey, N. Froyd, and Y. Dotsenko. Scalability analysis of SPMD codes using expectations. In *ICS ’07: Proc. of the 21st International Conference on Supercomputing*, pages 13–22, New York, NY, USA, 2007. ACM.
- [5] N. Corporation. Pc sampling, 2019. [Accessed January 26, 2019].
- [6] N. Froyd, J. Mellor-Crummey, and R. Fowler. Low-overhead call path profiling of unmodified, optimized code. In *Proc. of the 19th International Conference on Supercomputing*, pages 81–90, New York, NY, USA, 2005. ACM.
- [7] Lawrence Livermore National Laboratory. Laghos: High-order Lagrangian Hydrodynamics Miniapp. [Accessed February 27, 2020]. <https://computing.llnl.gov/projects/co-design/laghos>.
- [8] Lawrence Livermore National Laboratory. Quicksilver: A Proxy App for the Monte Carlo Transport Code, Mercury. [Accessed February 27, 2020]. <https://github.com/LLNL/Quicksilver>.
- [9] Libpfm4. Libpfm4: a helper library for performance tools using hardware counters. <http://perfmon2.sf.net/>, 2008.
- [10] P. E. McKenney. Differential profiling. *Software: Practice and Experience*, 29(3):219–234, 1999.
- [11] T. Mytkowicz, A. Diwan, M. Hauswirth, and P. F. Sweeney. Producing wrong data without doing anything obviously wrong! *SIGARCH Comput. Archit. News*, 37(1):265–276, Mar. 2009.

- [12] NVIDIA Corporation. *CUPTI User’s Guide DA-05679-001_v10.1*, 2019. https://docs.nvidia.com/cuda/pdf/CUPTI_Library.pdf.
- [13] Rice University. HPCToolkit performance tools. <http://hpctoolkit.org>.
- [14] N. Tallent, J. Mellor-Crummey, L. Adhianto, M. Fagan, and M. Krentel. HPCToolkit: Performance tools for scientific computing. *Journal of Physics: Conference Series*, 125:012088 (5pp), 2008.
- [15] N. R. Tallent, L. Adhianto, and J. M. Mellor-Crummey. Scalable identification of load imbalance in parallel executions using call path profiles. In *SC ’10: Proc. of the 2010 ACM/IEEE Conference on Supercomputing*, pages 1–11, Washington, DC, USA, 2010. IEEE Computer Society.
- [16] N. R. Tallent and J. Mellor-Crummey. Effective performance measurement and analysis of multithreaded applications. In *PPoPP ’09: Proc. of the 14th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 229–240, New York, NY, USA, 2009. ACM.
- [17] N. R. Tallent, J. Mellor-Crummey, and M. W. Fagan. Binary analysis for measurement and attribution of program performance. In *PLDI ’09: Proc. of the 2009 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 441–452, New York, NY, USA, 2009. ACM. **Distinguished Paper**.
- [18] N. R. Tallent, J. M. Mellor-Crummey, L. Adhianto, M. W. Fagan, and M. Krentel. Diagnosing performance bottlenecks in emerging petascale applications. In *SC ’09: Proc. of the 2009 ACM/IEEE Conference on Supercomputing*, pages 1–11, New York, NY, USA, 2009. ACM.
- [19] N. R. Tallent, J. M. Mellor-Crummey, M. Franco, R. Landrum, and L. Adhianto. Scalable fine-grained call path tracing. In *ICS ’11: Proc. of the 25th International Conference on Supercomputing*, pages 63–74, New York, NY, USA, 2011. ACM.
- [20] N. R. Tallent, J. M. Mellor-Crummey, and A. Porterfield. Analyzing lock contention in multithreaded applications. In *PPoPP ’10: Proc. of the 15th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 269–280, New York, NY, USA, 2010. ACM.

Appendix A

Environment Variables

HPCTOOLKIT’s measurement subsystem decides what and how to measure using information it obtains from environment variables. This chapter describes all of the environment variables that control HPCTOOLKIT’s measurement subsystem.

When using HPCTOOLKIT’s `hpcrun` script to measure the performance of dynamically-linked executables, `hpcrun` takes information passed to it in command-line arguments and communicates it to HPCTOOLKIT’s measurement subsystem by appropriately setting environment variables. To measure statically-linked executables, one first adds HPCTOOLKIT’s measurement subsystem to a binary as it is linked by using HPCTOOLKIT’s `hpclink` script. Prior to launching a statically-linked binary that includes HPCTOOLKIT’s measurement subsystem, a user must manually set environment variables.

Section A.1 describes environment variables of interest to users. Section A.2 describes environment variables designed for use by HPCTOOLKIT developers. In some cases, HPCTOOLKIT’s developers will ask a user to set some of the environment variables described in Section A.2 to generate a detailed error report when problems arise.

A.1 Environment Variables for Users

HPCTOOLKIT. Under normal circumstances, there is no need to use this environment variable. However, there are two situations, however, `hpcrun` *must* consult the HPCTOOLKIT environment variable to determine the location of HPCTOOLKIT’s top-level installation directory:

- On some systems, parallel job launchers (e.g., Cray’s `aprun`) *copy* the `hpcrun` script to a different location. In this case, for `hpcrun` to find libraries and utilities it needs at runtime, you must set the `HPCTOOLKIT` environment variable to HPCTOOLKIT’s top-level installation directory.
- If you launch the `hpcrun` script via a file system link, you must set `HPCTOOLKIT` for the same reason.

HPCRUN_EVENT_LIST. This environment variable is used provide a set of (event, period) pairs that will be used to configure HPCTOOLKIT’s measurement subsystem to

perform asynchronous sampling. The HPCRUN_EVENT_LIST environment variable must be set otherwise HPCToolkit's measurement subsystem will terminate execution. If an application should run with sampling disabled, HPCRUN_EVENT_LIST should be set to NONE. Otherwise, HPCToolkit's measurement subsystem expects an event list of the form shown below.

$$event1[@period1]; \dots; eventN[@periodN]$$

As denoted by the square brackets, periods are optional. The default period is 1 million.

Flags to add an event with `hpcrun`: `-e/--event event1[@period1]`

Multiple events may be specified using multiple instances of `-e/--event` options.

HPCRUN_TRACE. If this environment variable is set, HPCToolkit's measurement subsystem will collect a trace of sample events as part of a measurement database in addition to a profile. HPCToolkit's `hpctraceviewer` utility can be used to view the trace after the measurement database are processed with either HPCToolkit's `hpcprof` or `hpcprofmpi` utilities.

Flags to enable tracing with `hpcrun`: `-t/--trace`

HPCRUN_OUT_PATH If this environment variable is set, HPCToolkit's measurement subsystem will use the value specified as the name of the directory where output data will be recorded. The default directory for a command *command* running under control of a job launcher with as job ID *jobid* is `hptoolkit-command-measurements[-jobid]`. (If no job ID is available, the portion of the directory name in square brackets will be omitted. Warning: Without a *jobid* or an output option, multiple profiles of the same *command* will be placed in the same output directory.)

Flags to set output path with `hpcrun`: `-o/--output directoryName`

HPCRUN_PROCESS_FRACTION If this environment variable is set, HPCTOOLKIT's measurement subsystem will measure only a fraction of an execution's processes. The value of HPCRUN_PROCESS_FRACTION may be written as a floating point number or as a fraction. So, '0.10' and '1/10' are equivalent. If HPCRUN_PROCESS_FRACTION is set to a value with an unrecognized format, HPCTOOLKIT's measurement subsystem will use the default probability of 0.1. For each process, HPCTOOLKIT's measurement subsystem will generate a pseudo-random value in the range [0.0, 1.0). If the generated random number is less than the value of HPCRUN_PROCESS_FRACTION, then HPCTOOLKIT will collect performance measurements for that process.

Flags to set process fraction with `hpcrun`: `-f/-fp/--process-fraction frac`

HPCRUN_MEMLEAK_PROB If this environment variable is set, HPCTOOLKIT's measurement subsystem will measure only a fraction of an execution's memory allocations, e.g., calls to `malloc`, `calloc`, `realloc`, `posix_memalign`, `memalign`, and `valloc`. All allocations monitored will have their corresponding calls to free monitored as well. The value of HPCRUN_MEMLEAK_PROB may be written as a floating point number or as a fraction. So, '0.10' and '1/10' are equivalent. If HPCRUN_MEMLEAK_PROB is set to a value with

an unrecognized format, HPCTOOLKIT’s measurement subsystem will use the default probability of 0.1. For each memory allocation, HPCTOOLKIT’s measurement subsystem will generate a pseudo-random value in the range [0.0, 1.0). If the generated random number is less than the value of HPCRUN_MEMLEAK_PROB, then HPCTOOLKIT will monitor that allocation.

Flags to set process fraction with `hpcrun`: `-mp/--memleak-prob prob`

HPCRUN_DELAY_SAMPLING If this environment variable is set, HPCToolkit’s measurement subsystem will initialize itself but not begin measurement using sampling until the program turns on sampling by calling `hpctoolkit_sampling_start()`. To measure only a part of a program, one can bracket that with `hpctoolkit_sampling_start()` and `hpctoolkit_sampling_stop()`. Sampling may be turned on and off multiple times during an execution, if desired.

Flags to delay sampling with `hpcrun`: `-ds/--delay-sampling`

HPCRUN_RETAIN_RECURSION Unless this environment variable is set, by default HPCToolkit’s measurement subsystem will summarize call chains from recursive calls at a depth of two. Typically, application developers have no need to see performance attribution at all recursion depths when an application calls recursive procedures such as quicksort. Setting this environment variable may dramatically increase the size of calling context trees for applications that employ bushy subtrees of recursive calls.

Flags to retain recursion with `hpcrun`: `-r/--retain-recursion`

HPCRUN_MEMSIZE If this environment variable is set, HPCToolkit’s measurement subsystem will allocate memory for measurement data in segments using the value specified for HPCRUN_MEMSIZE (rounded up to the nearest enclosing multiple of system page size) as the segment size. The default segment size is 4M.

Flags to set memsize with `hpcrun`: `-ms/--memsize bytes`

HPCRUN_LOW_MEMSIZE If this environment variable is set, HPCToolkit’s measurement subsystem will allocate another segment of measurement data when the amount of free space available in the current segment is less than the value specified by HPCRUN_LOW_MEMSIZE. The default for low memory size is 80K.

Flags to set low memsize with `hpcrun`: `-lm/--low-memsize bytes`

A.2 Environment Variables for Developers

HPCRUN_WAIT If this environment variable is set, HPCToolkit’s measurement subsystem will spin wait for a user to attach a debugger. After attaching a debugger, a user can set breakpoints or watchpoints in the user program or HPCToolkit’s measurement subsystem before continuing execution. To continue after attaching a debugger, use the debugger to set the program variable DEBUGGER_WAIT=0 and then continue. Note: Setting HPCRUN_WAIT can only be cleared by a debugger if HPCTOOLKIT has been

built with debugging symbols. Building HPCTOOLKIT with debugging symbols requires configuring HPCTOOLKIT with `--enable-develop`.

HPCRUN_DEBUG_FLAGS HPCTOOLKIT supports a multitude of debugging flags that enable a developer to log information about HPCToolkit’s measurement subsystem as it records sample events. If `HPCRUN_DEBUG_FLAGS` is set, this environment variable is expected to contain a list of tokens separated by a space, comma, or semicolon. If a token is the name of a debugging flag, the flag will be enabled, it will cause HPCToolkit’s measurement subsystem to log messages guarded with that flag as an application executes. The complete list of dynamic debugging flags can be found in HPCToolkit’s source code in the file `src/tool/hpcrun/messages/messages.flag-defns`. A special flag value “`ALL`” enables all flags. Note: not all debugging flags are meaningful on all architectures.

Caution: turning on debugging flags will typically result in voluminous log messages, which will typically will dramatically slow measurement of the execution under study.

Flags to set debug flags with `hpcrun`: `-dd/--dynamic-debug flag`

HPCRUN_ABORT_TIMEOUT If an execution hangs when profiled with HPCToolkit’s measurement subsystem, the environment variable `HPCRUN_ABORT_TIMEOUT` can be used to specify the number of seconds that an application should be allowed to execute. After executing for the number of seconds specified in `HPCRUN_ABORT_TIMEOUT`, HPCToolkit’s measurement subsystem will forcibly terminate the execution and record a core dump (assuming that core dumps are enabled) to aid in debugging.

Caution: for a large-scale parallel execution, this might cause a core dump for each process, depending upon the settings for your system. Be careful!

HPCRUN_FNBOUNDS_CMD For dynamically-linked executables, this environment variable must be set to the full path of a copy of HPCToolkit’s `hpcfnbounds` utility. There are presently two versions of this utility. One, known as `hpcfnbounds`, analyzes program load modules (the executable and shared libraries) using Dyninst to recover a table of addresses that represent the beginning of each function. A second version of the tool, known as `hpcfnbounds2`, was designed to compute the same set of addresses for a load module using only a lightweight inspection of the load module’s symbol table and DWARF information. `hpcfnbounds2` is over a factor of ten faster and uses over a factor of 10 less memory than the original. `hpcfnbounds2` is the default. If `hpcfnbounds2` delivers an unsatisfactory result, a user can employ `hpcfnbounds` instead by setting this environment variable using the `--fnbounds` command line argument to `hpcrun`.