

October 10-12, 2023



ALCF Hands-on HPC Workshop

THAPI and iprof Hands On

Aurelio A. Vivas Meza; Solomon Bekele; Thomas Applencourt; **Brice Videau**

October 11, 2023

Table of Contents

Installation

Basic Usage: iprof

Advanced Usage: THAPI

Questions?

Installation

Install via Spack

If THAPI is not already installed on your system, You can easily install it via Spack:

- Clone Spack:

```
1 git clone -c feature.manyFiles=true -b releases/latest https://github.com/spack/spack.git
```

- Load Spack into your environment:

```
1 . spack/share/spack/setup-env.sh
```

- Clone THAPI-spack:

```
1 git clone https://github.com/argonne-lcf/THAPI-spack
```

- Add THAPI-spack to Spack and install:

```
1 spack repo add ./THAPI-spack/
2 spack install thapi
```

Install specific versions via Spack

Specific versions can be specified on the spack install line.

- Example using 0.0.10:

```
1 spack install thapi@0.0.10
```

- Can also use git hash or branch name

For more detailed instruction see:

https://spack-tutorial.readthedocs.io/en/latest/tutorial_basics.html

Load THAPI via Spack

To use a package installed by Spack.

- Load Spack into your environment (if it is not already loaded):

```
1 . spack/share/spack/setup-env.sh
```

- Load the default THAPI version:

```
1 spack load thapi
```

- Or load a specific version, 0.0.10 for example:

```
1 spack load thapi@0.0.10
```

Building THAPI from Sources

Sometimes, you may want to build THAPI from source. It is easy to achieve using spack.

- Load Spack into your environment (if it is not already loaded):

```
1 . spack/share/spack/setup-env.sh
```

- Clone THAPI-spack:

```
1 git clone https://github.com/argonne-lcf/THAPI
```

- Create a shell with the THAPI build environment:

```
1 spack build-env thapi bash
```

- Build as a regular autotools project:

```
1 cd THAPI
2 ./autogen.sh
3 mkdir build
4 cd build/
5 ../configure [--prefix=...]
6 make install
```

Basic Usage: iprof

Loading THAPI for the Workshop

As the spack install can be long, installation on Polaris has been done already.

- Load Spack into your environment on Polaris:

```
1 cd /eagle/projects/fallwkshp23/THAPI/
2 . spack/share/spack/setup-env.sh
```

- Load THAPI:

```
1 spack load thapi
```

- Test that iprof is working:

```
1 > iprof --help
2 iprof: a tracer / summarizer of OpenCL, L0, CUDA, HIP, and OMPT driver calls
3 Usage:
4   iprof -h | --help
5   iprof [option]... <application> <application-arguments>
6   iprof [option]... -r [<trace>]...
7 [...]
```

Profiling a Simple CUDA Applications with iprof: Polaris

- Several sample applications are located inside the `/eagle/projects/fallwkshp23/THAPI` directory.
- Each application is located in 3 sub-directories:
 - *CUDA*: as simple addition CUDA application
 - *SYCL*: as simple hello world SYCL application that runs on the GPU
 - *OMP*: as simple benchmark OPENMP offload application taht runs on the GPU
- Some directories may contain a *README.md* file with important information regarding the environment required to run the application.

Profiling a Simple CUDA Applications with iprof

Example: running the simple CUDA application:

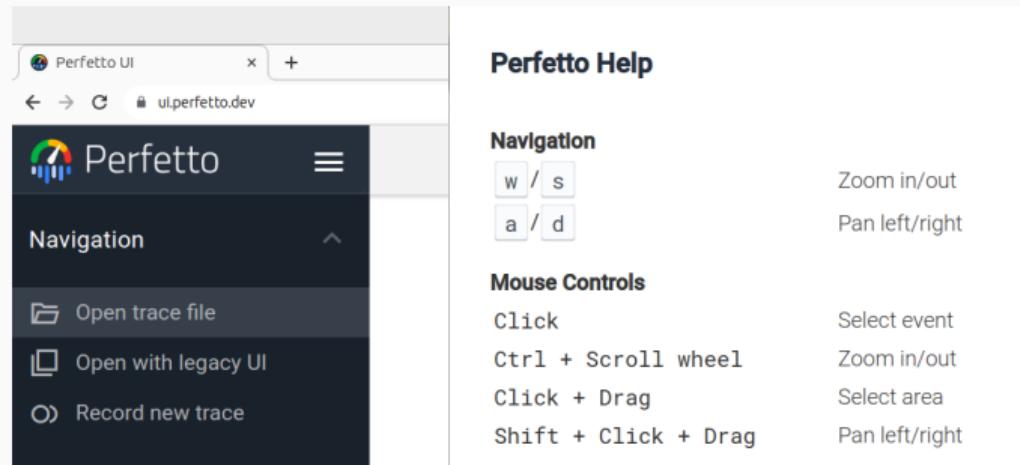
```
1 > cd /eagle/projects/fallwkshp23/THAPI/CUDA
2 > ./cuda_hello
3 Max error: 0
4 > iprof ./cuda_hello
5 Trace location: /home/videau/lttng-traces/iprof-20231010-025849
6
7 BACKEND_CUDA | 1 Hostnames | 1 Processes | 1 Threads |
8
9          Name |      Time | Time(%) | Calls |   Average |      Min |      Max | Error |
10         cuInit | 154.38ms | 43.41% |    1 | 154.38ms | 154.38ms | 154.38ms | 0 |
11        cuCtxSynchronize | 122.67ms | 34.49% |    1 | 122.67ms | 122.67ms | 122.67ms | 0 |
12 cuDevicePrimaryCtxRetain | 71.31ms | 20.05% |    1 | 71.31ms | 71.31ms | 71.31ms | 0 |
13 cuDeviceTotalMem_v2 | 3.01ms | 0.85% |    4 | 753.07us | 731.66us | 772.24us | 0 |
14 cuDeviceGetAttribute | 1.81ms | 0.51% | 392 | 4.61us | 838ns | 172.09us | 0 |
15 cuGetProcAddress | 1.64ms | 0.46% | 373 | 4.39us | 908ns | 531.57us | 0 |
16 [...]
17        cuCtxGetDevice | 2.44us | 0.00% |    1 | 2.44us | 2.44us | 2.44us | 0 |
18 cuDeviceGetCount | 1.47us | 0.00% |    1 | 1.47us | 1.47us | 1.47us | 0 |
19 cuDevicePrimaryCtxRelease | | | 1 | | | | | 1 |
20        cuModuleUnload | | | 1 | | | | | 1 |
21          Total | 355.67ms | 100.00% | 803 | | | | 2 |
22
23 Device profiling | 1 Hostnames | 1 Processes | 1 Threads | 1 Devices | 1 Subdevices |
24
25          Name |      Time | Time(%) | Calls |   Average |      Min |      Max |
26         add | 122.68ms | 100.00% |    1 | 122.68ms | 122.68ms | 122.68ms |
27       Total | 122.68ms | 100.00% |    1 |
```

- By default traces are stored on the disc in the `~/lttng-traces/` directory.
- Last trace can be replayed using the `-r` option: `iprof -r`
- A specific trace can be replayed using the `-r` option followed by the trace location. If using the trace from the previous slide:

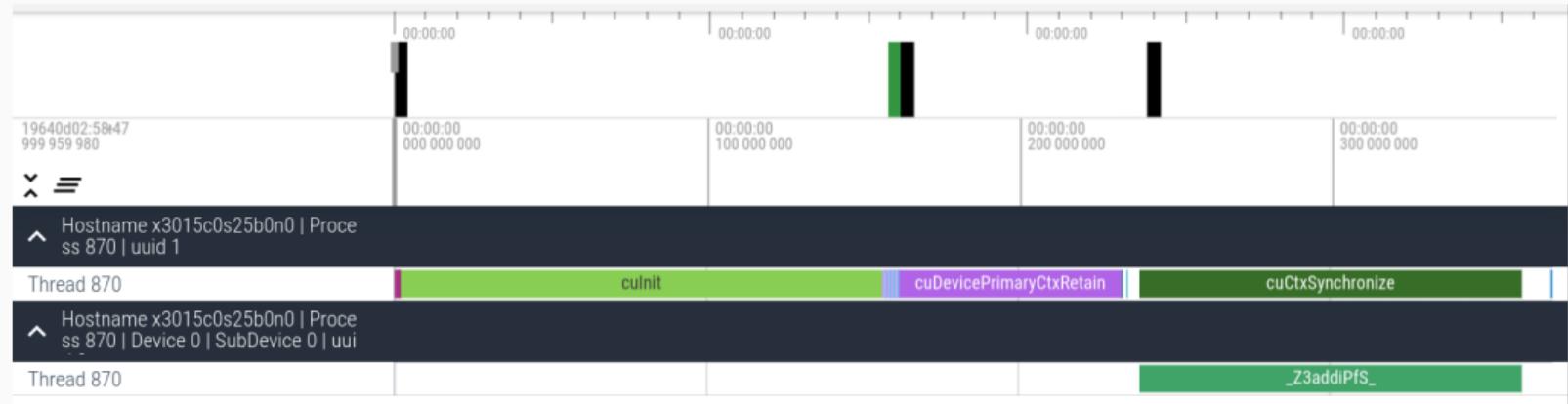
```
iprof -r /home/videau/lttng-traces/iprof-20231010-025849
```

Generating a Timeline

- A timeline can be generated using the `-l` option.
- Often used while replaying a trace, like this:
`iprof -l -r /home/videau/lttng-traces/iprof-20231010-025849`
- Generates a `out.pftrace` file that can be loaded in Perfetto UI:
`https://ui.perfetto.dev/`

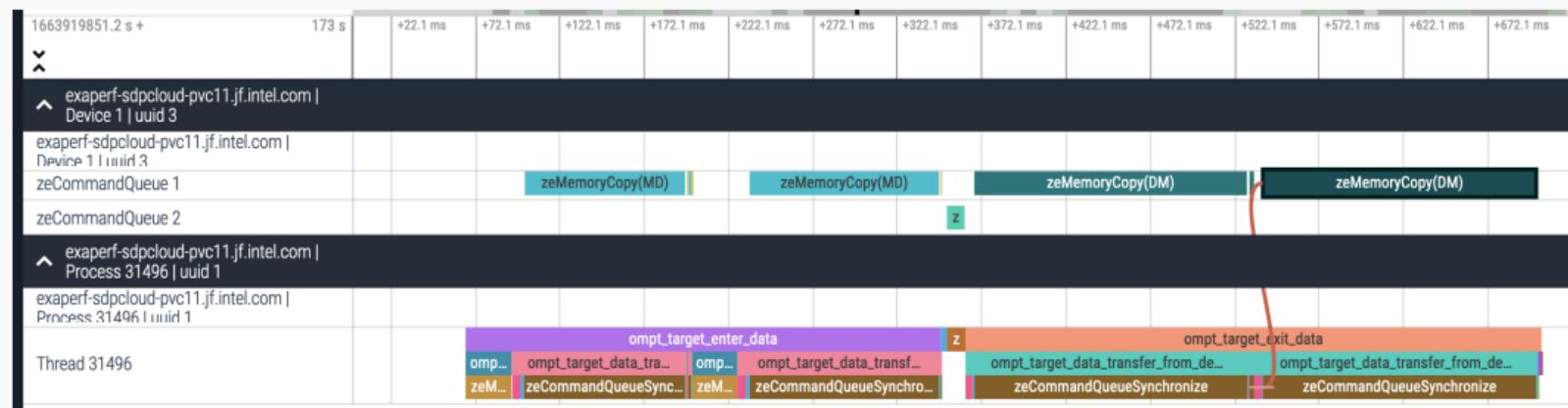


Perfetto Timeline



Perfetto Timeline: OpenMP on Level Zero

The OpenMP compiler on Polaris does not support OMPT. But the Intel compiler on JLSE does:



Displaying a Trace

- A trace can be generated using the `-t` option.
- Often used while replaying a trace, like this:

```
iprof -l -t /home/videau/lttng-traces/iprof-20231010-025849
```

- Generates a trace containing all the API calls (and more)

Simple CUDA Trace

```
1  [...]
2  02:58:48.234853435 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuModuleGetFunction_entry: { hfunc: 0x00007ffc1ebf8048,
3   ↪ hmod: 0x00000000012e43d0, name: 0x000000000046dcbd, name_val: "_Z3addiPfS_" }
4  02:58:48.234858185 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuModuleGetFunction_exit: { cuResult: CUDA_SUCCESS,
5   ↪ hfunc_val: 0x00000000012e1b60 }
6  02:58:48.234867264 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemAllocManaged_entry: { dptr: 0x00007ffc1ebf8380,
7   ↪ bytesize: 4194304, flags: 1 }
8  02:58:48.234900648 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemAllocManaged_exit: { cuResult: CUDA_SUCCESS,
9   ↪ dptr_val: 0x00001471e4000000 }
10 02:58:48.234902604 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemAllocManaged_entry: { dptr: 0x00007ffc1ebf8378,
11   ↪ bytesize: 4194304, flags: 1 }
12 02:58:48.234911613 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemAllocManaged_exit: { cuResult: CUDA_SUCCESS,
13   ↪ dptr_val: 0x00001471e4400000 }
14 02:58:48.239007688 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuLaunchKernel_entry: { f: 0x000000000012e1b60, gridDimX:
15   ↪ 1, gridDimY: 1, gridDimZ: 1, blockDimX: 1, blockDimY: 1, blockDimZ: 1, sharedMemBytes: 0, hStream: 0x0000000000000000,
16   ↪ kernelParams: 0x00007ffc1ebf8300, extra: 0x0000000000000000, extra_vals: [ ] }
17 02:58:48.239034018 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda_profiling:event_profiling: { hStart: 0x000000000012e7810,
18   ↪ hStop: 0x000000000012e1780 }
19 02:58:48.239035415 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuLaunchKernel_exit: { cuResult: CUDA_SUCCESS }
20 02:58:48.239038627 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuCtxSynchronize_entry: { }
21 02:58:48.361712445 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuCtxSynchronize_exit: { cuResult: CUDA_SUCCESS }
22 02:58:48.370911152 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemFree_v2_entry: { dptr: 0x00001471e4000000 }
23 02:58:48.371248487 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemFree_v2_exit: { cuResult: CUDA_SUCCESS }
24 02:58:48.371249954 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemFree_v2_entry: { dptr: 0x00001471e4400000 }
25 02:58:48.371527993 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemFree_v2_exit: { cuResult: CUDA_SUCCESS }
26 02:58:48.371538120 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda_profiling:event_profiling_results: { hStart:
27   ↪ 0x000000000012e7810, hStop: 0x000000000012e1780, startStatus: CUDA_SUCCESS, stopStatus: CUDA_SUCCESS, status: CUDA_SUCCESS,
28   ↪ milliseconds: 122.68134307861328 }
29 [...]
```

Cleaning up

- When *iprof* seems to be acting weird, it is usually following a *Ctrl-C* or a crash.
- Cleaning up leftover daemons usually solves issues: *iprof --clean-up*

Advanced Usage: THAPI

THAPI Consists of Two Big Components

Open source at: <https://github.com/argonne-lcf/THAPI>

- The tracing of events
 - Use low level tracing: Linux Tracing Toolkit Next Generation (LTTng):
 - Tracepoints are generated from APIs' headers
- The parsing of the trace
 - Use Babeltrace2 library and tools
 - Pretty Printer, Tally, Timeline/Flamegraph, ...

Supported APIs

- OpenCL, Level Zero, Cuda Runtime/Driver, HIP
- OMPT

iprof is an orchestrator around THAPI, lttng, and babeltrace2.

iprof is a rather complex script, but simpler tools dedicated to each API exist. Example:
tracer_cuda.sh

- Will capture a CUDA trace
- Provides additional options, useful for debugging applications
- Collected trace can be replayed with *iprof*
- Can be used as a base to design user tailored tracers

How Tracers Work - 1

- Use LTTng to setup a capture daemon

```
1 lttnng-sessiond --daemonize --quiet
```

- Create a LTTng session

```
1 lttnng create thapi-cuda-session
```

- Create a blocking channel and add context information. Context can also be perf counters.

```
1 lttnng enable-channel --userspace --blocking-timeout=inf blocking-channel  
2 lttnng add-context --userspace --channel=blocking-channel -t vpid -t vtid
```

- Enable desired events, can be individual events or regular expressions:

```
1 lttnng enable-event --channel=blocking-channel --userspace lttnng_ust_cuda:*
```

How Tracers Work - 2

- Setup interposition libraries:

```
1 # Point tracer to CUDA driver library
2 export LTTNG_UST_CUDA_LIBCUDA=/usr/lib64/libcuda.so
3 # Put our interception version of the CUDA library in the path
4 export LD_LIBRARY_PATH=$pkglibdir/cuda:$LD_LIBRARY_PATH
5 # Preload it as well
6 export LD_PRELOAD=$libdir/libTracerCUDA.so:$LD_PRELOAD
7 # Enable LTTng client side blocking mode
8 export LTTNG_UST_ALLOW_BLOCKING=1
```

- Start LTTng

```
1 lttnng start
```

- Launch the application
- Stop LTTng

```
1 lttnng stop
```

Reading a Trace with Babeltrace2 - 1

Processing the trace can be done through Babeltrace 2 bindings. Here is a ruby example:

```
1 require 'babeltrace2'
2 # Find babeltrace 2 plugins, the CTF reader, and the stream muxer
3 ctf_fs = BT2::BTPlugin.find("ctf").get_source_component_class_by_name("fs")
4 utils_muxer = BT2::BTPlugin.find("utils").get_filter_component_class_by_name("muxer")
5 # trace locations have the metadata file generated by lttng at runtime
6 trace_locations = ARGV
7 # Create a processing graph:
8 graph = BT2::BTGraph.new
9 # Add each trace to the graph:
10 sources = trace_locations.each_with_index.collect do |trace_location, i|
11   graph.add_component(ctf_fs, "trace_#{i}", params: {"inputs" => [ trace_location ] })
12 end
13 # Mux the traces together
14 muxer = graph.add_component(utils_muxer, "mux")
15 sources.map(&:output_ports).flatten.each_with_index do |port, i|
16   graph.connect_ports(port, muxer.input_port(i))
17 end
```

Reading a Trace with Babeltrace2 - 2

```
1 # Simple pretty printer
2 sink = graph.add_simple_sink("babeltrace_thapi", lambda do |it, _|
3     it.next_messages.select { |m| m.type == :BT_MESSAGE_TYPE_EVENT }.each do |message|
4         event = message.event
5         str = message.get_default_clock_snapshot.ns_from_origin.to_s
6         str << " - #{event.stream.trace.get_environment}"
7         str << " - #{event.common_context_field.value}" if event.get_common_context_field
8         str << " - #{event.name}: #{event.payload_field}"
9         puts str
10    end
11 end)
12 # Connect the muxer output
13 graph.connect_ports(muxer.output_port(0), sink.input_port(0))
14 # Execute the Babeltrace 2 graph
15 graph.run
```

Similar Python bindings exist, or plugins can be written in C.

Questions?
