# THAPI/Iprof: Design and Implementation of a Tracer for Heterogeneous API

Thomas Applencourt, Solomon Bekele, Brice Videau and many others
September 25, 2025

- "The audience is very technical", Mike
- So I made the slides in Beamer...

# Table of Contents

What?

Why?

Pretty Pictures

How?

Conclusion / Open Questions

One Last note about MPI deployment

Argonne
NATIONAL LABORATORY

# What?

- `mpirun -n 60000 -- iprof ./a.out`
- THAPI (Tracing Heterogeneous API) is... a tracer for heterogeneous API.
  - We support OpeCL, L0, Cuda (Driver and Runtime) HIP, OpenMP, MPI
  - We Dump All Arguments. Traces should contain enough information to reconstruct the programming model state.
  - We do hardware counter sampling (frequency, power, traffic...)
- Iprof take the trace, and post-process it (high-level summary, timeline, ...)
- Scalable (tested on 10k GPUs), Low/Reasonable overhead (about 0.2us overhead ns per tracepoint)

Argonne ◆
NATIONAL LABORATORY

- It's not a full-blow performance analysis framework (Vtune, NSigh, HPC Toolkit, Tau)
- It's not a line-level profiler [1]

---

[1]We give you Kernel Time. And we know have sampling support of HW counter, but we stop here

Argonne
NATIONAL LABORATORY

# Why?

We work on runtime. So we Need to understand what is going on to solve bug.

- Why no data-transfer H2D, D2H overlap?!
- Why OpenMP Mapping take 10 min?!
- Why my SYCL queue in-order have so much submission overhead?!

Application:

- Does I'm GPU bound? MPI Bound? Data-transfer bound?
- What is my memory footprint?
- How does my scaling affect my offload
- And give me some timeline to see if I see some "bubble"

Argonne
NATIONAL LABORATORY

- When starting working on GPU-Aurora, their was not tracer for Level Zero[2].
- We wanted just a tracer. Nothing else.
- Want to use "industry standard" binary trace format so we can develop tools on top of it.
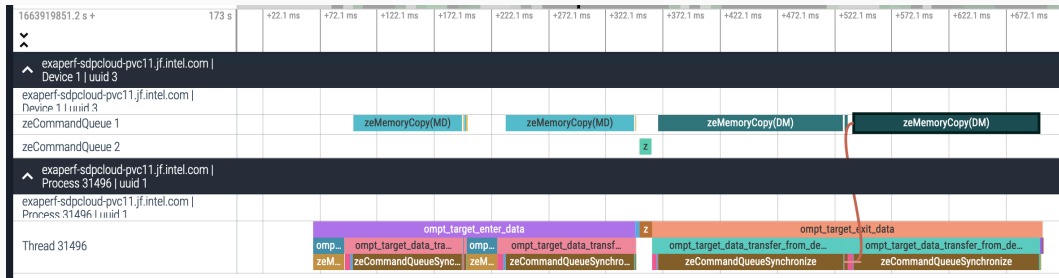- Easy to maintain

---

[2]L0 is the cuda-driver of Intel

Argonne ▲
NATIONAL LABORATORY

# Pretty Pictures

```
 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 |
```

# Perfetto Timeline: OpenMP on top of Level Zero

Argonne
NATIONAL LABORATORY

# Simple CUDA Trace

```
 1   [...]
 2   02:58:48.234853435 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuModuleGetFunction_entry: { hfunc: 0x00007ffc1ebf8048,
     ↪ hmod: 0x0000000012e43d0, name: 0x000000000046dcbd, name_val: "_Z3addiPfS_" }
 3   02:58:48.234858185 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuModuleGetFunction_exit: { cuResult: CUDA_SUCCESS,
     ↪ hfunc_val: 0x0000000012e1b60 }
 4   02:58:48.234867264 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemAllocManaged_entry: { dptr: 0x00007ffc1ebf8380,
     ↪ bytesize: 4194304, flags: 1 }
 5   02:58:48.234900648 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemAllocManaged_exit: { cuResult: CUDA_SUCCESS,
     ↪ dptr_val: 0x00001471e4000000 }
 6   02:58:48.234902604 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemAllocManaged_entry: { dptr: 0x00007ffc1ebf8378,
     ↪ bytesize: 4194304, flags: 1 }
 7   02:58:48.234911613 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemAllocManaged_exit: { cuResult: CUDA_SUCCESS,
     ↪ dptr_val: 0x00001471e4400000 }
 8   02:58:48.239007688 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuLaunchKernel_entry: { f: 0x0000000012e1b60, gridDimX:
     ↪ 1, gridDimY: 1, gridDimZ: 1, blockDimX: 1, blockDimY: 1, blockDimZ: 1, sharedMemBytes: 0, hStream: 0x0000000000000000,
     ↪ kernelParams: 0x00007ffc1ebf8300, extra: 0x0000000000000000, extra_vals: [   ] }
 9   02:58:48.239034018 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda_profiling:event_profiling: { hStart: 0x0000000012e7810,
     ↪ hStop: 0x0000000012e1780 }
10   02:58:48.239035415 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuLaunchKernel_exit: { cuResult: CUDA_SUCCESS }
11   02:58:48.239038627 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuCtxSynchronize_entry: {   }
12   02:58:48.361712445 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuCtxSynchronize_exit: { cuResult: CUDA_SUCCESS }
13   02:58:48.370911152 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemFree_v2_entry: { dptr: 0x00001471e4000000 }
14   02:58:48.371248487 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemFree_v2_exit: { cuResult: CUDA_SUCCESS }
15   02:58:48.371249954 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemFree_v2_entry: { dptr: 0x00001471e4400000 }
16   02:58:48.371527993 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda:cuMemFree_v2_exit: { cuResult: CUDA_SUCCESS }
17   02:58:48.371538120 - x3015c0s25b0n0 - vpid: 870, vtid: 870 - lttng_ust_cuda_profiling:event_profiling_results: { hStart:
     ↪ 0x0000000012e7810, hStop: 0x0000000012e1780, startStatus: CUDA_SUCCESS, stopStatus: CUDA_SUCCESS, status: CUDA_SUCCESS,
     ↪ milliseconds: 122.68134307861328 }
18   [...]
```

Argonne
NATIONAL LABORATORY

How?

- With lot of love
- … and metaprograming.

Two Big Components:

- The tracing of events
  - Using: Linux Tracing Toolkit Next Generation (LTTng) to generate CTF Trace.
  - Tracepoints are generated from APIs' headers
- The parsing of the trace
  - Using: Babeltrace2 [3]
  - Babaltrace2 plugin infrastructure generated by MetaBabel
  - Pretty Printer, Tally, Timeline/Flamegraph, …

iprof is an orchestrator around THAPI, lTTng, and Babeltrace2.

LTTng, and Babeltrace2 are developed by EfficiOS (`https://www.efficios.com/`)

---

[3]Similar to ffmpeg pipeline approach

Argonne
NATIONAL LABORATORY

# LTTng

State of the art tracing infrastructure for kernel and user-space.

- Well maintained and established (used in industry leading data-centers)
- Binary format (CTF: Common Trace Format) open standard
- About 0.2us overhead per Tracepoint (in our case: blocking mode)
  - can be relaxed if use case tolerate event losses
- LTTng relay daemons can be setup to stream over the network in complex topologies, or block for "on-the-fly" analysis
  - ideal to deploy at scale
- Fine granularity, you can enable/disable individual tracepoint

Argonne
NATIONAL LABORATORY

To Intersect Symbol:

- we either use `LD_PRELOAD="fake_lib.so"` + у `dlopen("real_lib.so", RTLD_LAZY | RTLD_LOCAL)`

- or we use Native Support for Interception (OpenCL Layers, OpenMP [4])

Implementation of Tracing functions:

```
1  CUresult cuDeviceGetCount(int *count) {
2    tracepoint(lttng_ust_cuda, cuDeviceGetCount_entry, count);
3    CUresult _retval = CU_DEVICE_GET_COUNT_PTR(count);
4    tracepoint(lttng_ust_cuda, cuDeviceGetCount_exit, count, _retval);
5    return _retval;
6  }
```

```
1  21:03:53.070592532 - x3006c0s25b0n0 - vpid: 36056, vtid: 36056
2  - lttng_ust_cuda:cuDeviceGetCount_entry: { count: 0x00007ffe93bec390 }
3  21:03:53.070593929 - x3006c0s25b0n0 - vpid: 36056, vtid: 36056
4  - lttng_ust_cuda:cuDeviceGetCount_exit: { cuResult: CUDA_SUCCESS, count_val: 6 }
```

[4]But not PMPI, we want to trace people who use PMPI…

Argonne ▲
NATIONAL LABORATORY

- We trace all APIs entry points (multiple thousand tracepoint...).
- Tedious, error prone, and hard to maintain by hand
- Automatic generation from headers or API description (OpenCL)
  - C99 parser => YAML intermediary representation
  - YAML + user provided meta information + user provided tracepoints => wrapper functions + Trace Model
  - Trace Model => tracepoints

Argonne
NATIONAL LABORATORY

Cuda Header:
```
CUresult cuDeviceGetCount(int* count);
```

Metadata:
```
1  cuDeviceGetCount:
2  - [OutScalar, count]
```

```
1   - name: cuDeviceGetCount
2     type:
3       kind: custom_type
4       name: CUresult
5     params:
6     - name: count
7       type:
8         kind: pointer
9         type:
10          kind: int
11          name: int
```

The Final tracepoints:

```
1   - :name: cuDeviceGetCount_entry
2     :payload:
3     - :name: count
4       :cast_type: int *
5       :class: unsigned
6       :class_properties:
7        :field_value_range: 64
8        :preferred_display_base: 16
```

```
1   - :name: cuDeviceGetCount_exit
2     :payload:
3     - :name: cuResult
4       :cast_type: CUresult
5       :class: signed
6       :class_properties:
7        :field_value_range: 32
8       :be_class: CUDA::CUResult
9     - :name: count_val
10      :cast_type: int
11      :class: signed
12      :class_properties:
13       :field_value_range: 32
```

Argonne
NATIONAL LABORATORY

- Thanks you for asking!

Argonne
NATIONAL LABORATORY

Reference parser implementation of Common Trace Format

- Modular plugin infrastructure
- Compose Babeltrace 2 components into trace processing graphs:
  - Sources, Filters, Sources

```
1  babeltrace2 --plugin-path=$libdir \
2              --component=filter.zeinterval.interval \
3              --component=filter.ompinterval.interval \
4              --component=sink.xprof.tally
```

THAPI use Pipeline of plugins

- Source are the generated traces
- Filters which aggregate messages
- Sinks which create outputs: Tally, Pretty Print Timeline

Automatic Plugins generation for Babeltrace 2 from the Trace Model

Argonne
NATIONAL LABORATORY

- Problem: Writing Babeltrace 2 plugin by hand is tedious, error prone and hard to maintain.
  - Using Python bindings is too slow -> Use C or C++
- Main Idea: Attaching User-Callbacks to Trace Events
- Metababel generates Babeltrace 2 calls to read, write and dispatch events to User-Callbacks
  - Generate State Machine to handle Babeltrace 2 messages queues
- Open Source: `https://github.com/TApplencourt/metababel`

Argonne ▲
NATIONAL LABORATORY

# Code Generated for the Cu Exit

```
1   CUresult cuResult;
2   int count_val;
3   const bt_field *payload_field = bt_event_borrow_payload_field_const(bt_evt);
4   {
5       const bt_field *_field = NULL;
6       _field = bt_field_structure_borrow_member_field_by_index_const(payload_field, 0);
7       cuResult = (CUresult)bt_field_integer_signed_get_value(_field);
8   }
9   {
10      const bt_field *_field = NULL;
11      _field = bt_field_structure_borrow_member_field_by_index_const(payload_field, 1);
12      count_val = (int)bt_field_integer_signed_get_value(_field);
13  }
14  [...]
```

User code:

```
1   #include <metababel/metababel.h>
2   void cuDeviceGetCount_exit_callback(void *btx_handle, CUresult cuResult, int count_val) {
3       std::cout << "cuResult: " << cuResult << ", count_val: "  << count_val << std::endl;
4   }
5   void btx_register_usr_callbacks(void *btx_handle) {
6       btx_register_cuDeviceGetCount_exit(btx_handle, &cuDeviceGetCount_exit_callback);
7   }
```

Argonne ◢◣
NATIONAL LABORATORY

# Conclusion / Open Questions

- Scalable Tracer with low-overhead
- Based on industry standard technology (LTTng, Babeltrace2)
- Lot of meta-programming (so maintainable but a small team)
- Should we converge on a meta-data description of API format?

# One Last note about MPI deployment

# How to Launch the LTTng Daemon and Post-Processing?

- Problem 1: We need to spawn one daemon per node.
- We need first to spawn the daemon and then launch user-application
- Problem 2: `iprof mpirun ./a.out` or `mpirun iprof ./a.out`
- I was afraid of doing 1. MPI launcher scare me[5].

So new Problem… Problem 3:

- Application will call `MPI_init`, but we need a barrier before app run. And We cannot call 2 `MPI_init`
- `MPI_Session_init` "aka" initialize a mpi communicator without call mpi init.

---

[5]vni, weird argument parsing, …

Argonne
NATIONAL LABORATORY