# Proton

A ***Pro***filer for Tri***ton***

Keren Zhou
kzhou6@gmu.edu

# Goals

- Provide a quick, intuitive, and simple way to check kernel performance

    - Open source

    - Multiple vendor GPUs

    - Flexible metrics collection

        - Hardware metrics

        - Software metrics

    - Call path profiling

# Call Path Profiling
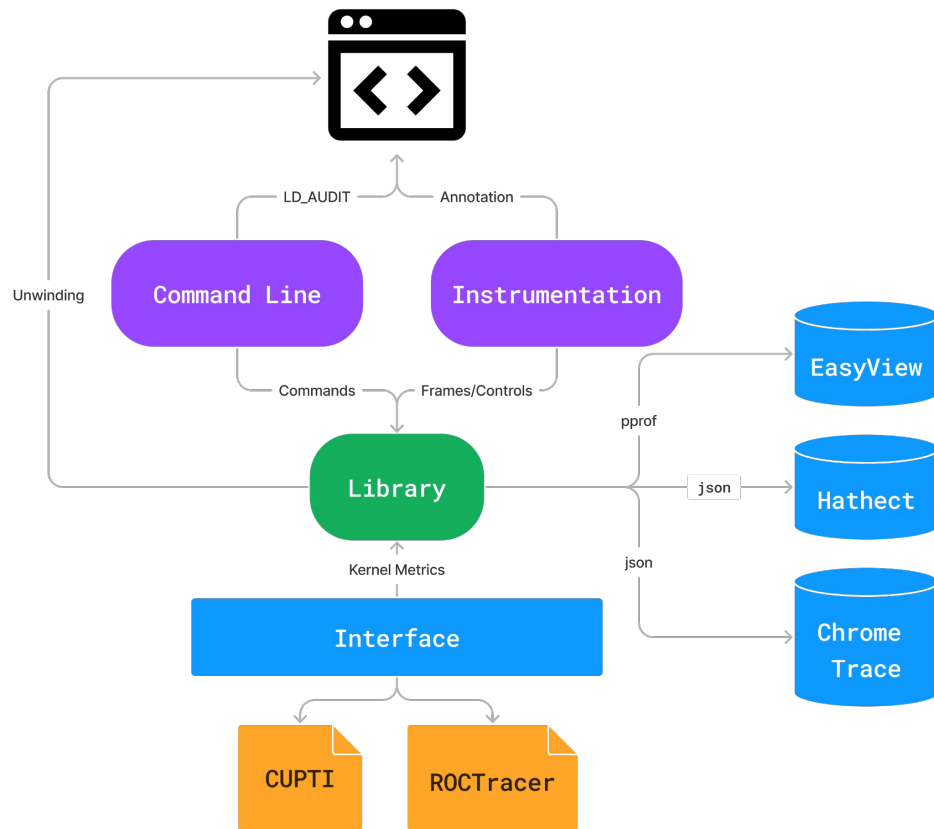
- Profile kernel running time



Python Context



Custom Context

# Design

# Inside the Library



Aggregate timing into kernels with the same "group name"

min/max/mean/stddev

# User Interface

- Lightweight instrumentation

    - Profile start/stop/finalize (*torchinductor compatible*)

    - Scopes

    - Metrics

    - Renaming

    - Hooks

- Command line (*torchinductor compatible*)

# Profile Start/Stop/Finalize

- Profile only interesting regions
    - **proton.start**(name: str, *, backend: str = "cupti", context: str = "shadow", data: str = "tree", hook: Optional[str|callable] = None) -> *session_id*: int
    - **proton.finalize**()
- Skip some regions, but accumulate to the same profile
    - **proton.start(...)**
    - **proton.deactive**(*session_id*)
    - ... # region skipped
    - **proton.activate**(*session_id*)
- Profile with multiple concurrent sessions
    - Different views (e.g., *tree*, *trace*, …)

# Scopes

- Only collect the "Master Thread" scope

  - In PyTorch, the thread that train and test models

```
with proton.scope("test0"):
    with proton.scope("test1"):
        foo[1,](x, y)
with proton.scope("test2"):
    foo[1,](x, y)
```

```
2368.000 ROOT
├── 1344.000 test0
│   └── 1344.000 test1
│       └── 1344.000 foo_0d1d
└── 1024.000 test2
    └── 1024.000 foo_0d1d

Legend (Metric: Time (ns) (inc) Min: 1024.00 Max: 2368.00)
2233.60 — 2368.00
1964.80 — 2233.60
1696.00 — 1964.80
1427.20 — 1696.00
1158.40 — 1427.20
1024.00 — 1158.40
```

# Metrics

- Asynchronous metrics

  - Come from profilers

- Synchronous metrics

  - Come from users

    - Theoretical flops, bytes

    - Loss

    - Counts

  - Dict[str, Union[int, float]]

*with* **proton.scope**("test0", {"foo_metric": 1.0}):
   foo[1,](x, y)

"test0" scope ends with multiple metrics.
*Two* metrics can be display the same time.

```
(triton) kzhou6@x-d-e5309-n05223:~/Code/proton/test$ proton-viewer -l ./test.hatchet
Available metrics:
Count
Time (ns)
foo_metric
```

```
1313.000 1.000 ROOT
└─ 1313.000 1.000 test0
   └─ 1313.000 nan foo_0d1d

Legend (Metric: Time (ns) (inc) Min: 1313.00 Max: 1313.00)
  1313.00 - 1313.00
  1313.00 - 1313.00
  1313.00 - 1313.00
  1313.00 - 1313.00
  1313.00 - 1313.00
  1313.00 - 1313.00
```

# Renaming

- Rename the triton function with a custom name

    - Append launch configurations

    - Append runtime dynamic

    - Append constants

    - e.g., foo_<num_warps:4>_<fast_math:4>_<branch_0:1>

- Can be used together with flexible metrics

```
with proton.Rename(foo_rename_fn):
    with proton.Metrics(foo_metric_fn):
        foo[1,](x, y, num_warps=4)
```

# Hooks

- Decorators that are less intrusive

```
@proton.metrics(metrics_fn)
@proton.rename(rename_fn)
@triton.jit
def triton_fn():
    ...
```

```
def metric_fn(grid_x, grid_y, grid_z,
    num_warps, num_ctas,
    cluster_x, cluster_y, cluster_z,
    shared_memory, stream, function, metadata,
    *args) -> Dict[str, Union[int, float]]
```

```
def rename_fn(grid_x, grid_y, grid_z,
    num_warps, num_ctas,
    cluster_x, cluster_y, cluster_z,
    shared_memory, stream, function, metadata,
    *args) -> str
```

# Plan

- Integrate into triton
    - third_party/proton
- AMD GPUs
    - ROCTracer
- Command line interface
- Fine-grained metrics
    - Instruction samples
    - Binary instrumentation-based metrics
    - …
- VSCode Integration