



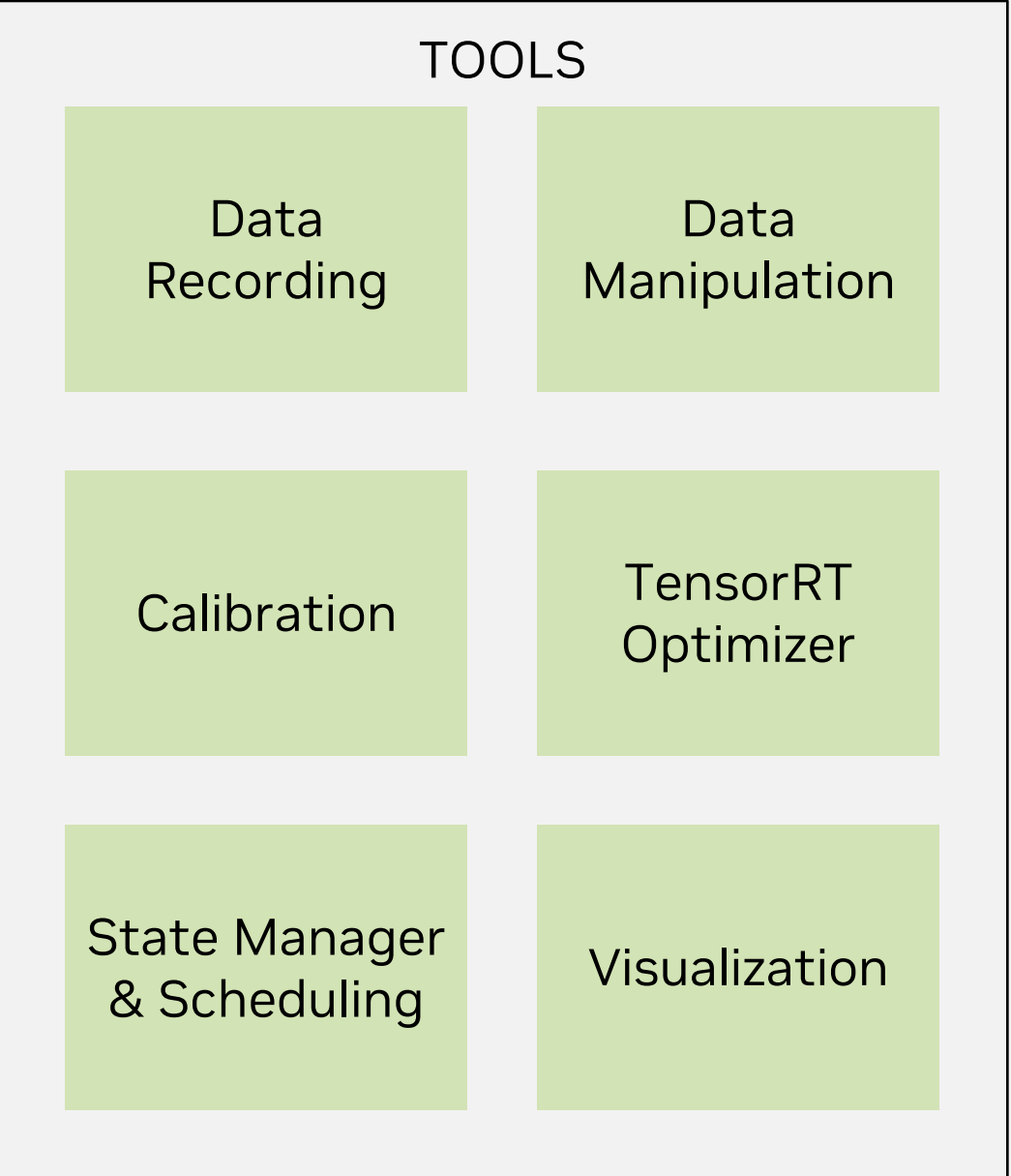
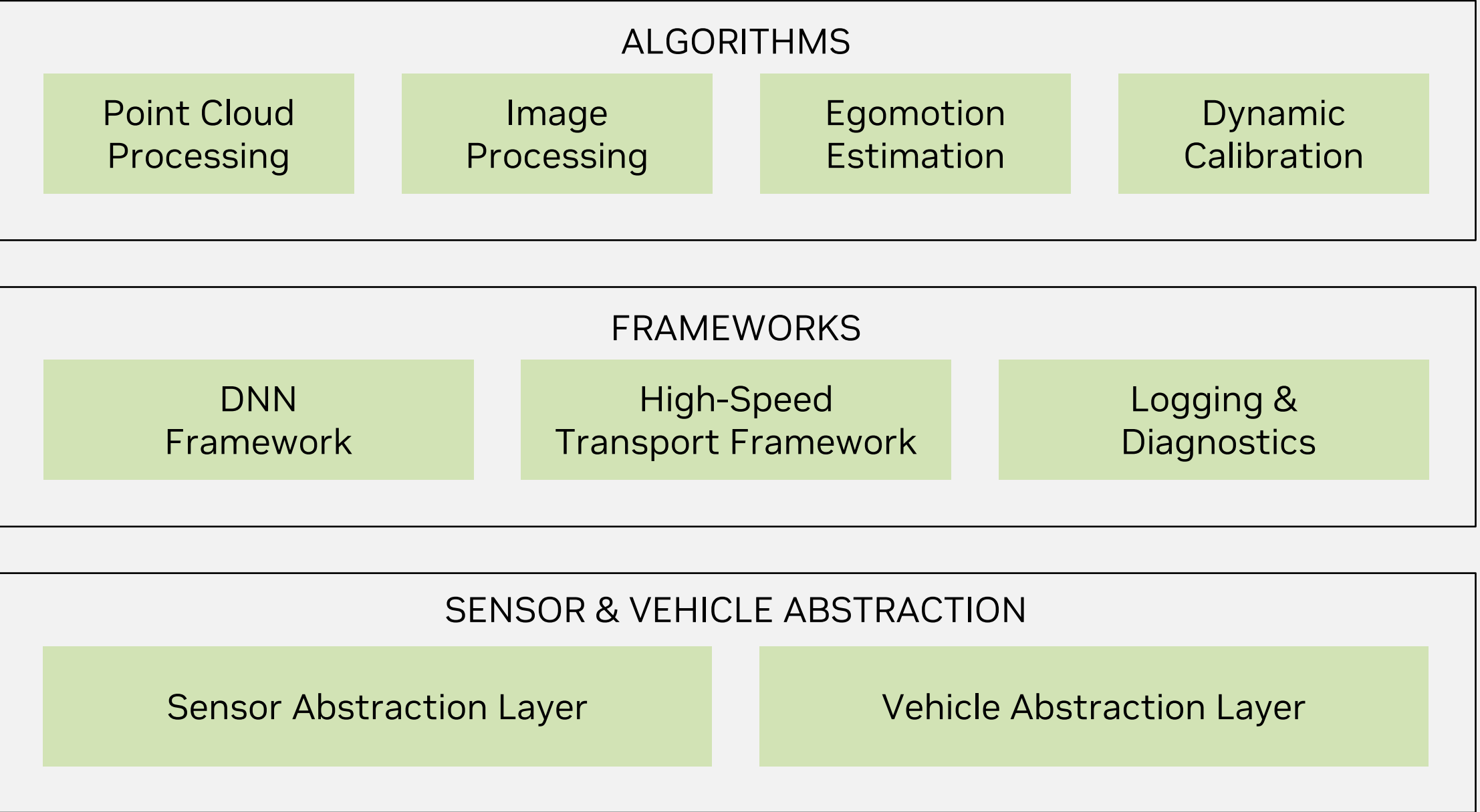
Performance Oriented Scheduling with System Task Manager

DRIVE SDK Supercharges AV Development

One architecture
DRIVE OS & DriveWorks
DRIVE AGX Orin

DriveWorks
Middleware

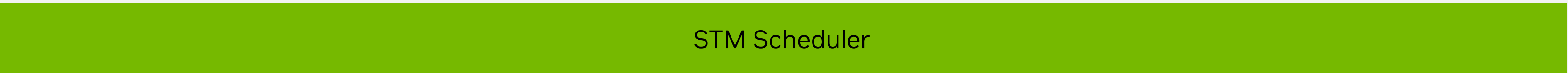
Application Building Blocks



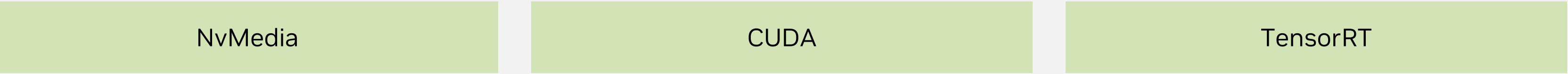
Application Framework



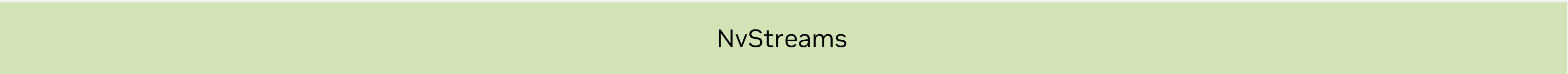
Deterministic Scheduler



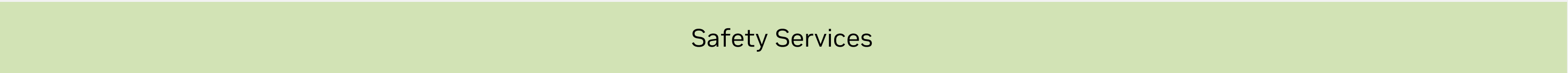
Application APIs



Hardware Synchronization



Safety Framework
(QNX Only)

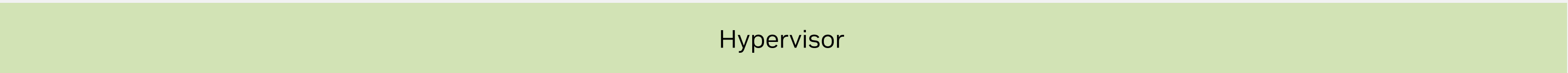


DRIVE OS
Linux / QNX

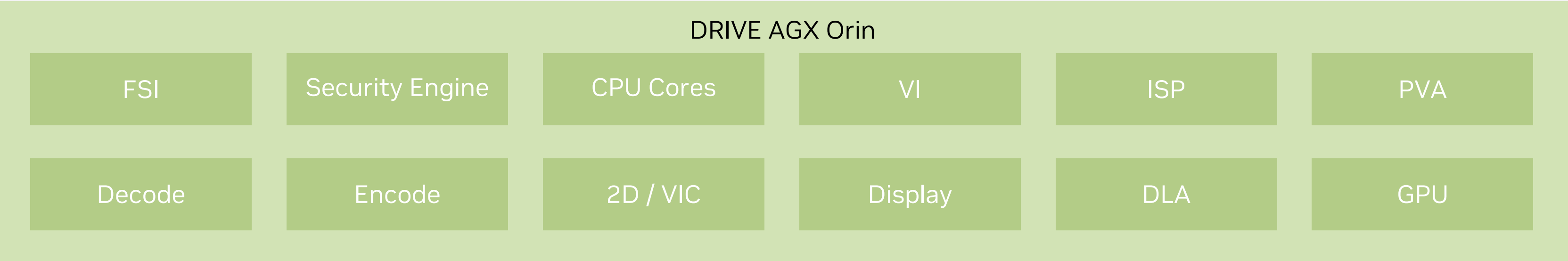
Operating System



Resource Isolation &
Freedom of Interference



Hardware



DriveWorks — Comprehensive Middleware Solution

**Rich Library of
Algorithms and Tools**
to accelerate your applications

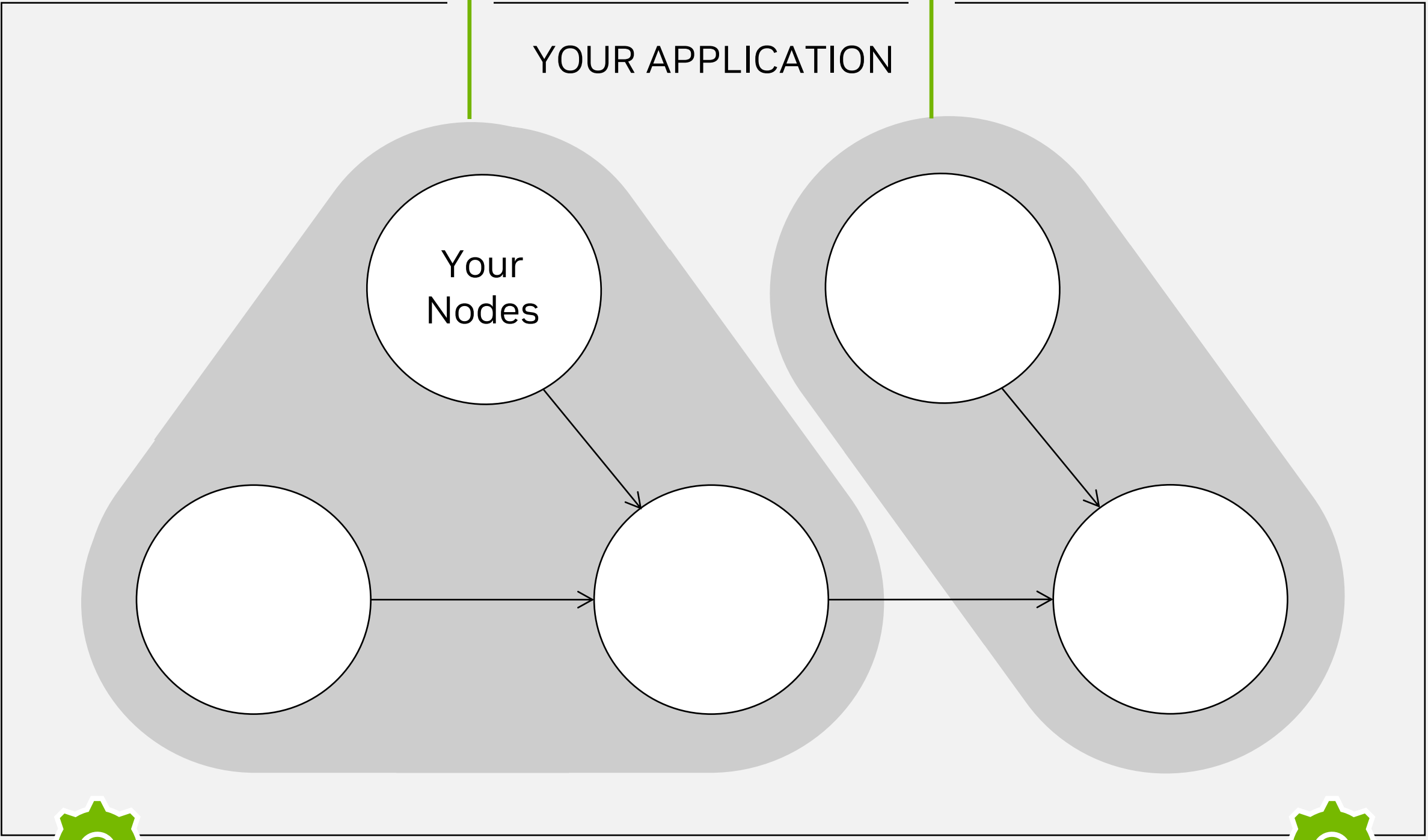
Compute Graph Framework
to simplify complex
computational pipelines

System Task Manager
to deterministically schedule
application tasks

DRIVEWORKS ACCELERATION LIBRARIES & TOOLS

Calibration	Visualization	Data Manipulation	TensorRT Optimizer
Vehicle Abstraction Layer	Egomotion Estimation	Data Recording	High-Speed Transport Framework
Sensor Abstraction Layer	Image Processing	Point Cloud Processing	Logging & Diagnostics

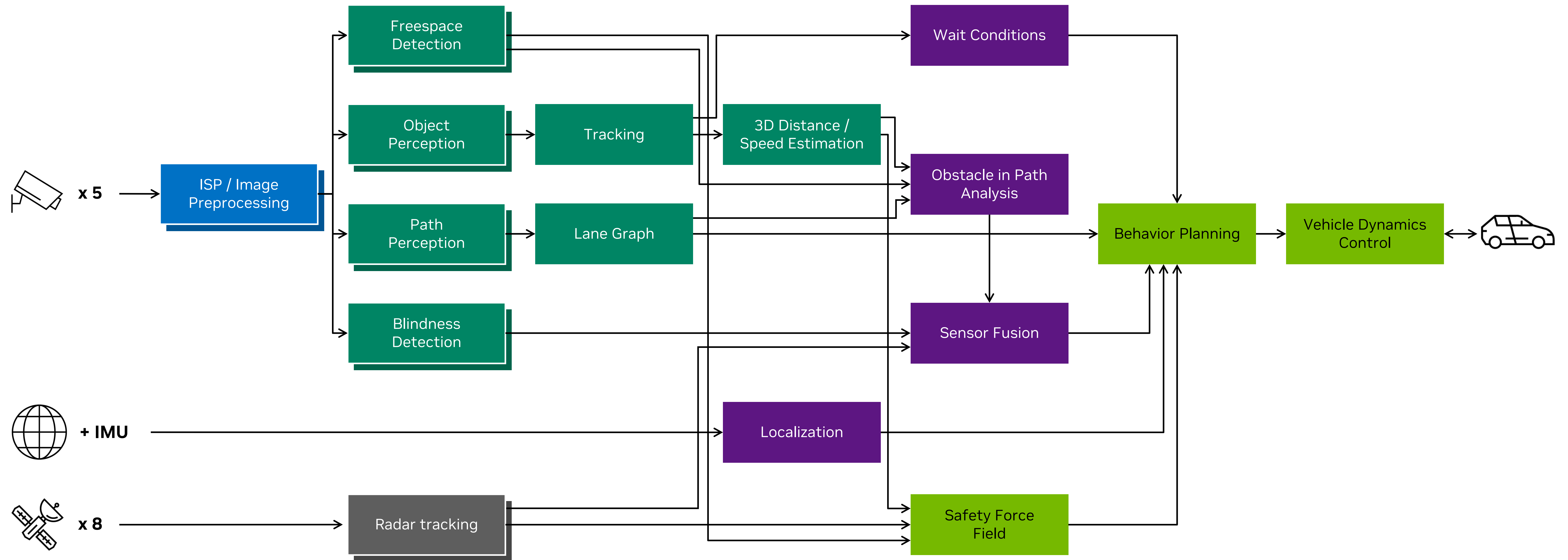
Compute Graph Framework CGF



STM Scheduler

AV Workload

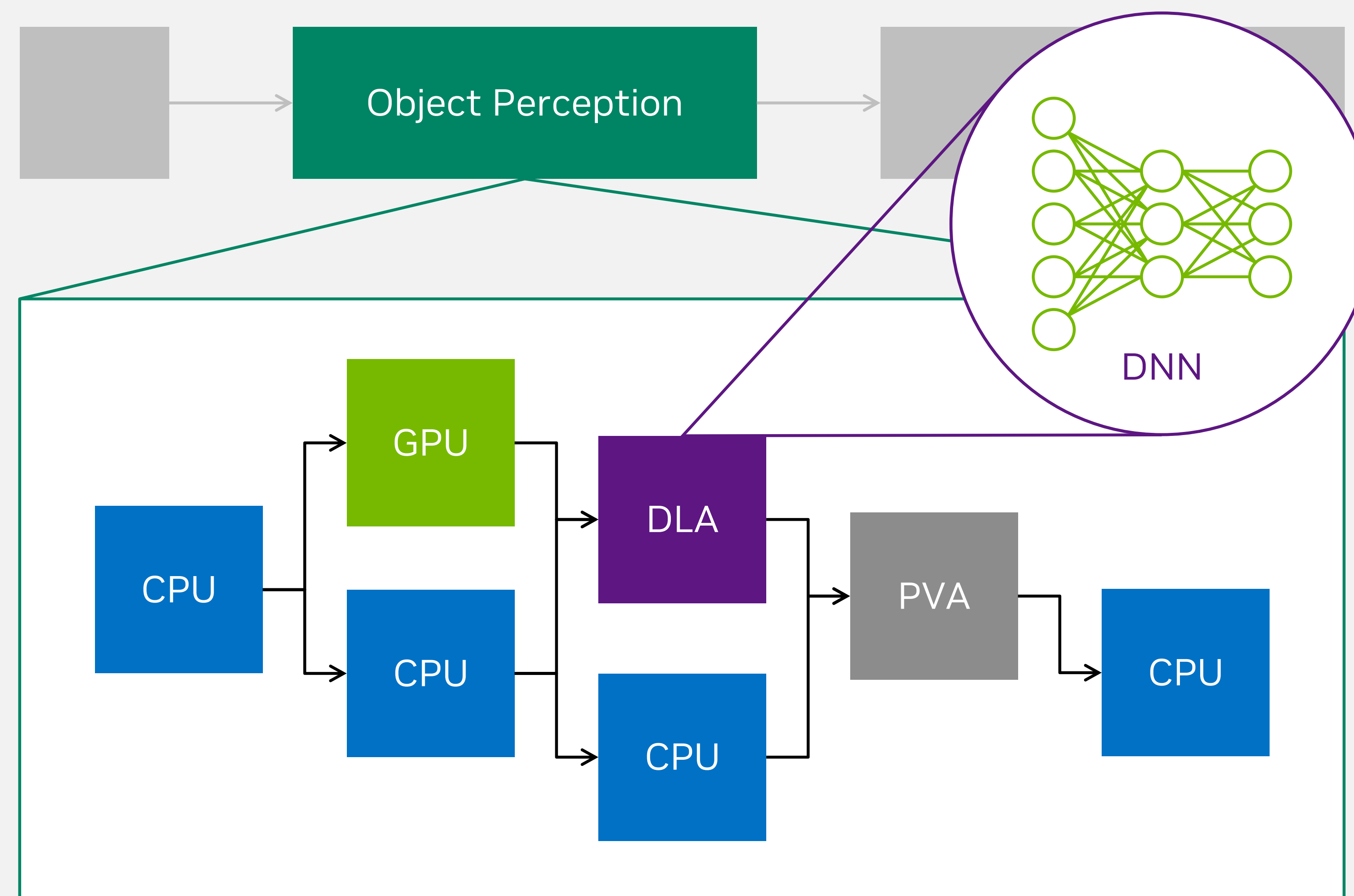
Complex (heterogenous) compute graph



**cartoon pipeline only for illustrative purposes*

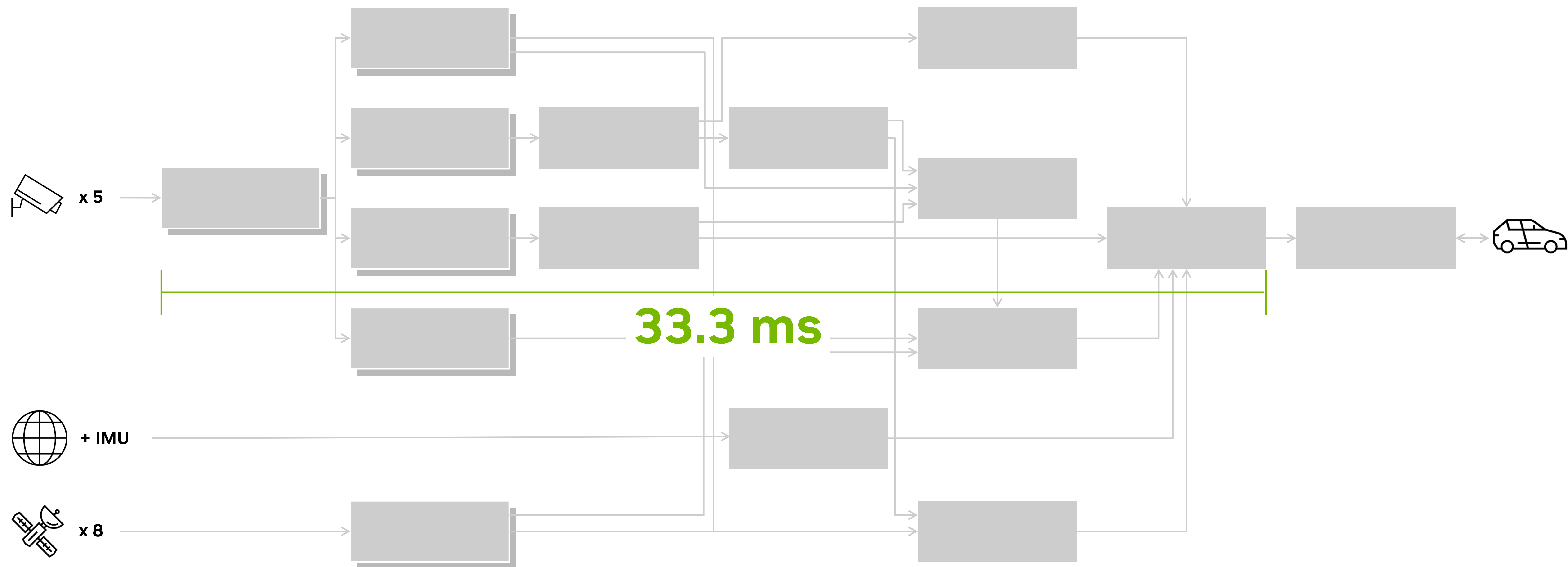
AV Workload

Each of these again a compute graph



AV Workload

In real-time...

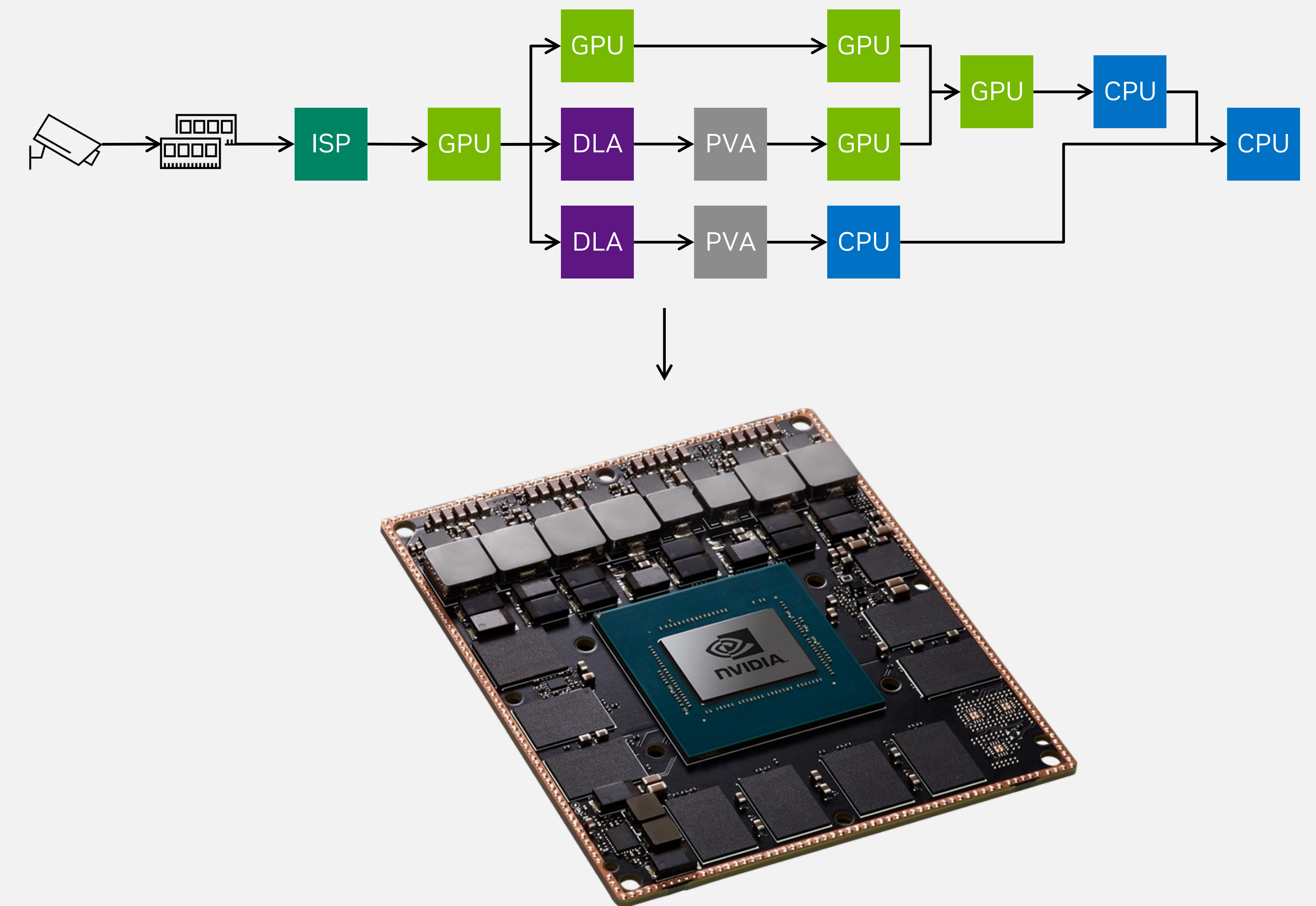


**latency bound only for illustrative purposes*

Scheduling Problem

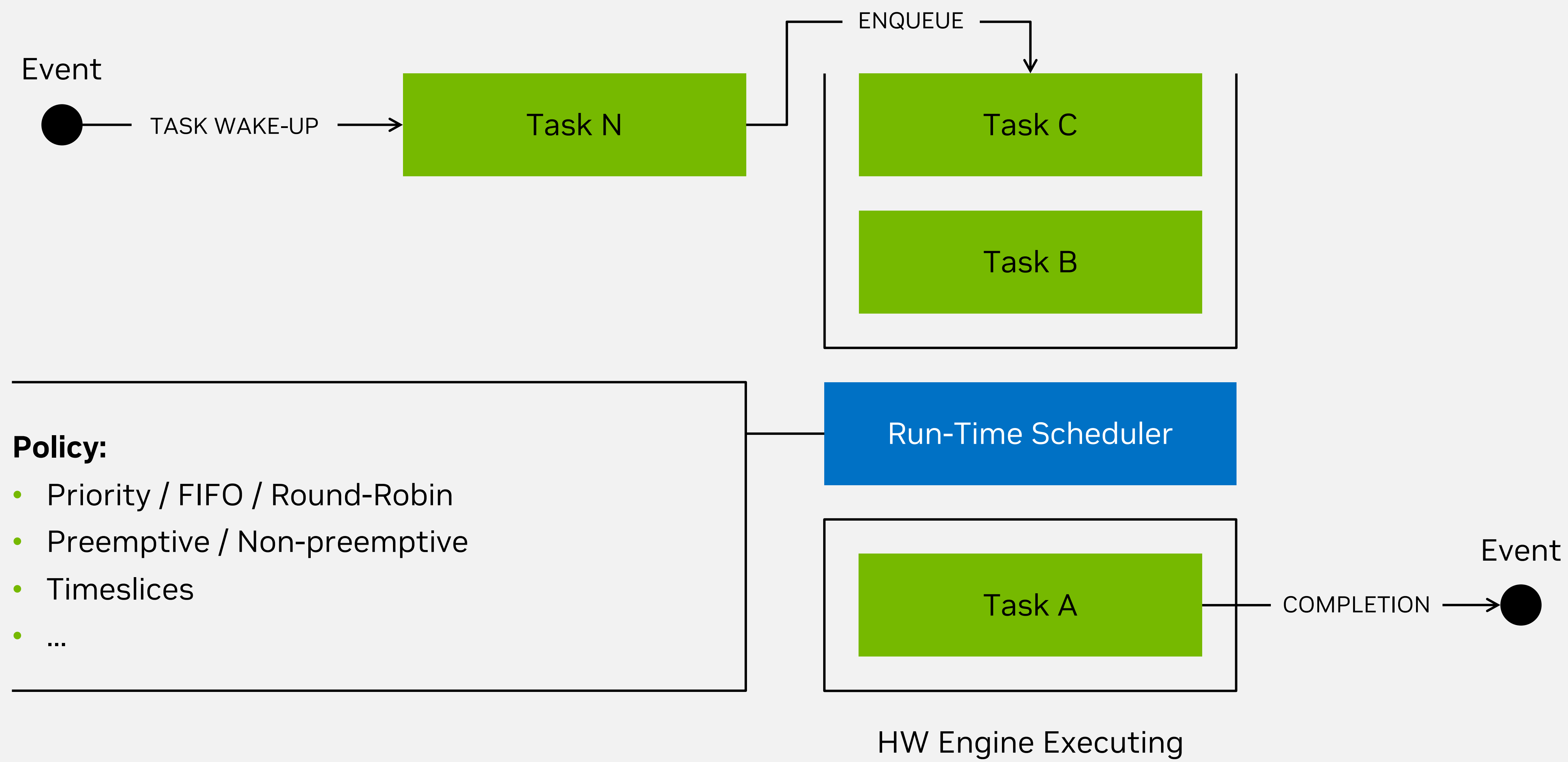
How to map...

- Constraints:
 - Latency / Worst-Case Response Time
 - Throughput / TOPS
 - Determinism
- Safety Considerations:
 - Program Flow-Monitoring
 - System Verification
 - ISO26262: “Manage / Restrict” use of interrupts, shared resources, concurrency, ...



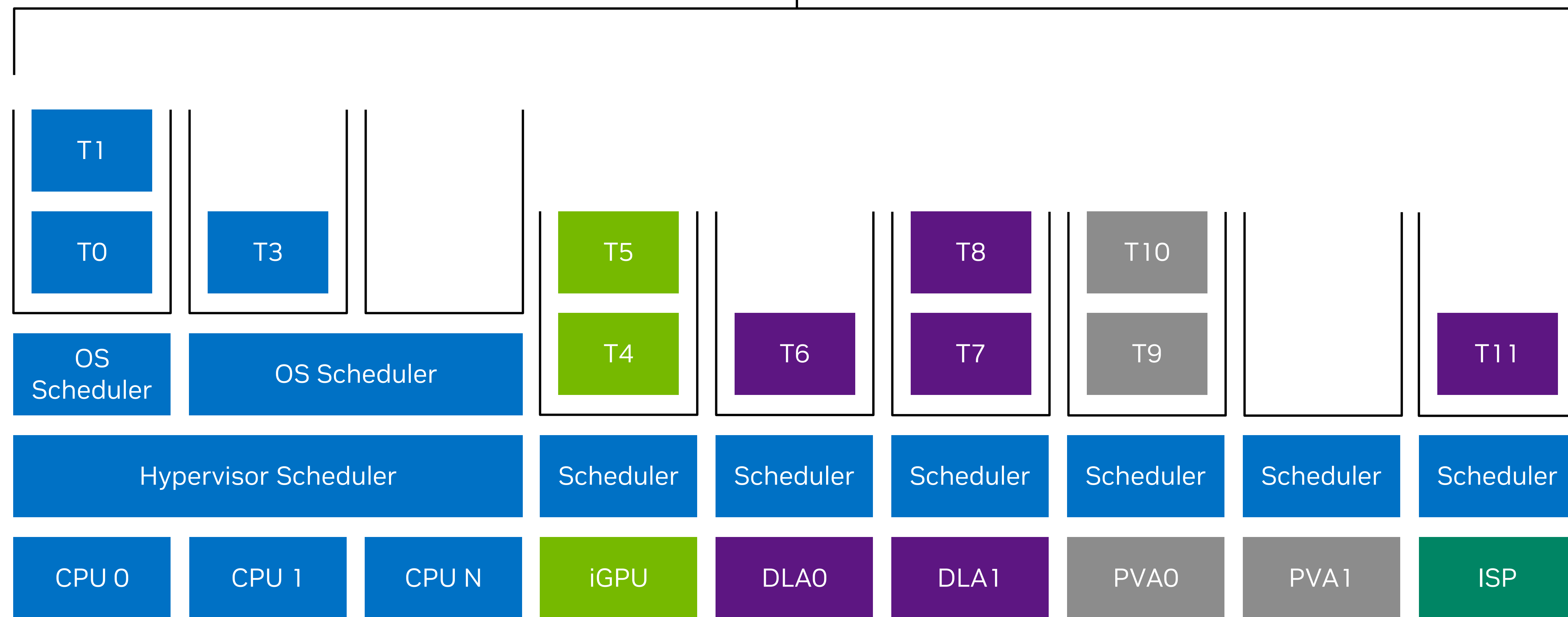
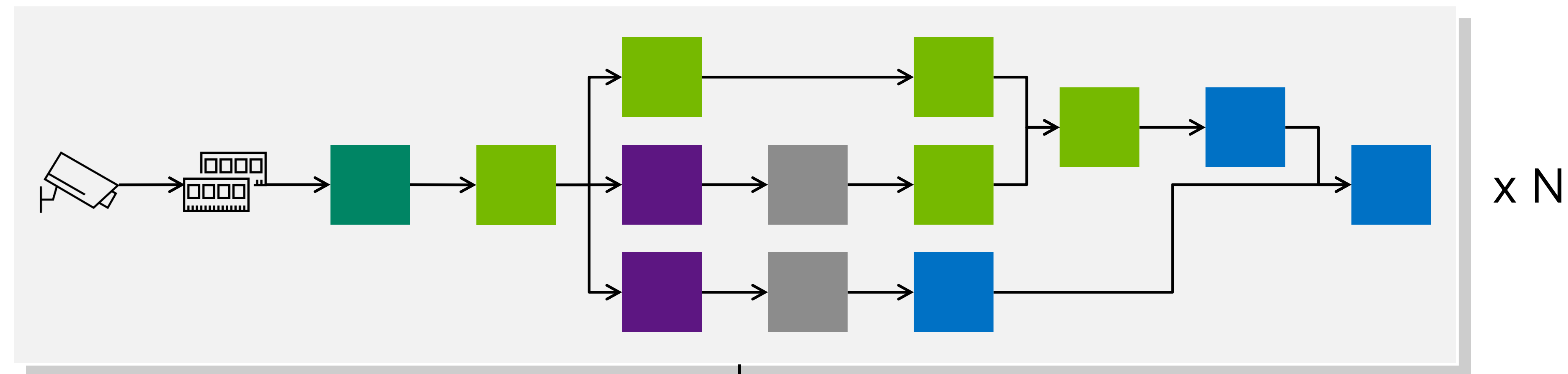
Traditional Scheduling View

Event-driven system



System-Level Scheduling Problem

How to guarantee determinism across many schedulers?



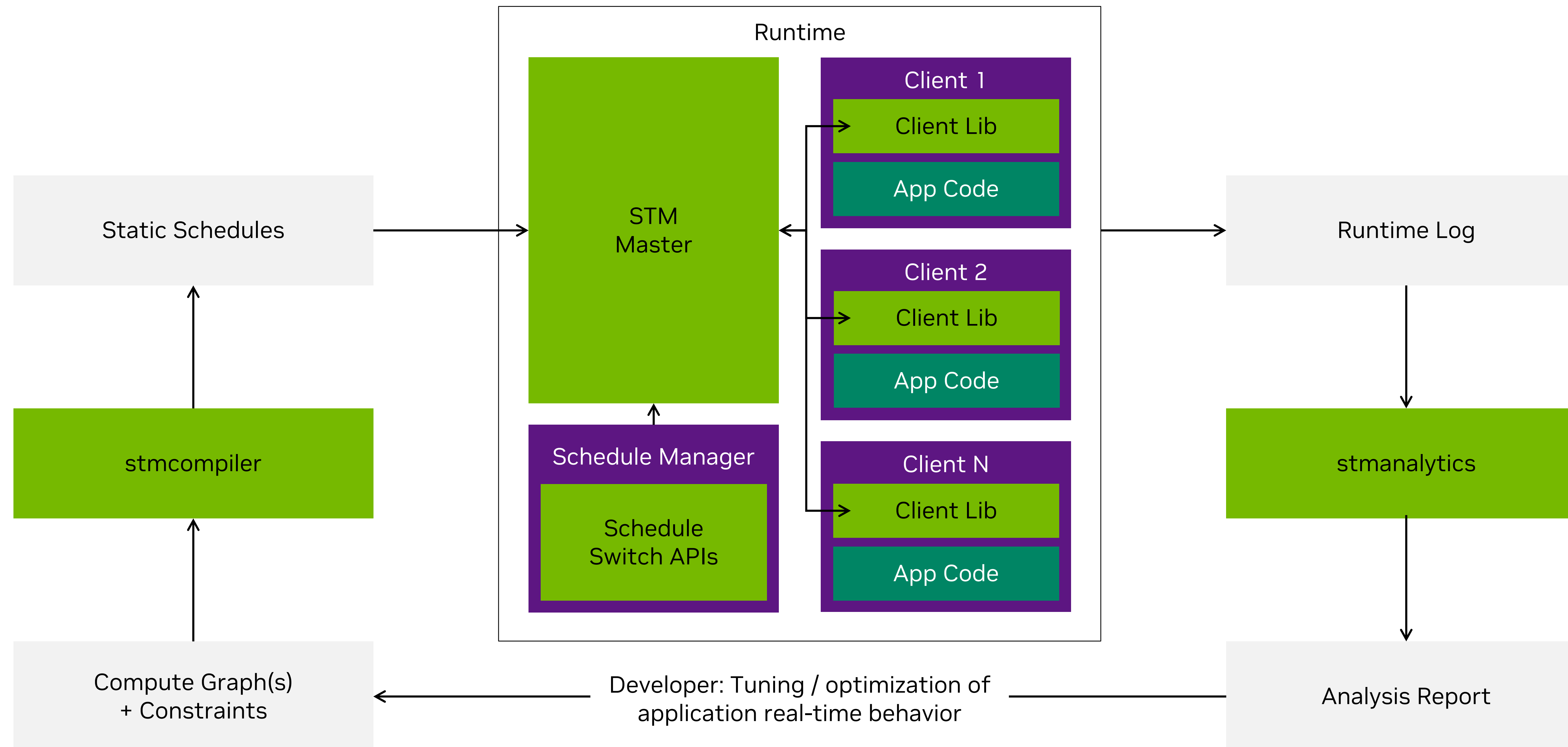
Introducing NVIDIA System Task Manager

Solution as part of NVIDIA DriveWorks

- User-Space
- Non-Preemptive
- Static Ordering
- Deterministic Performance
- Multi-Process
- Supports Heterogeneous Platforms

NVIDIA System Task Manager

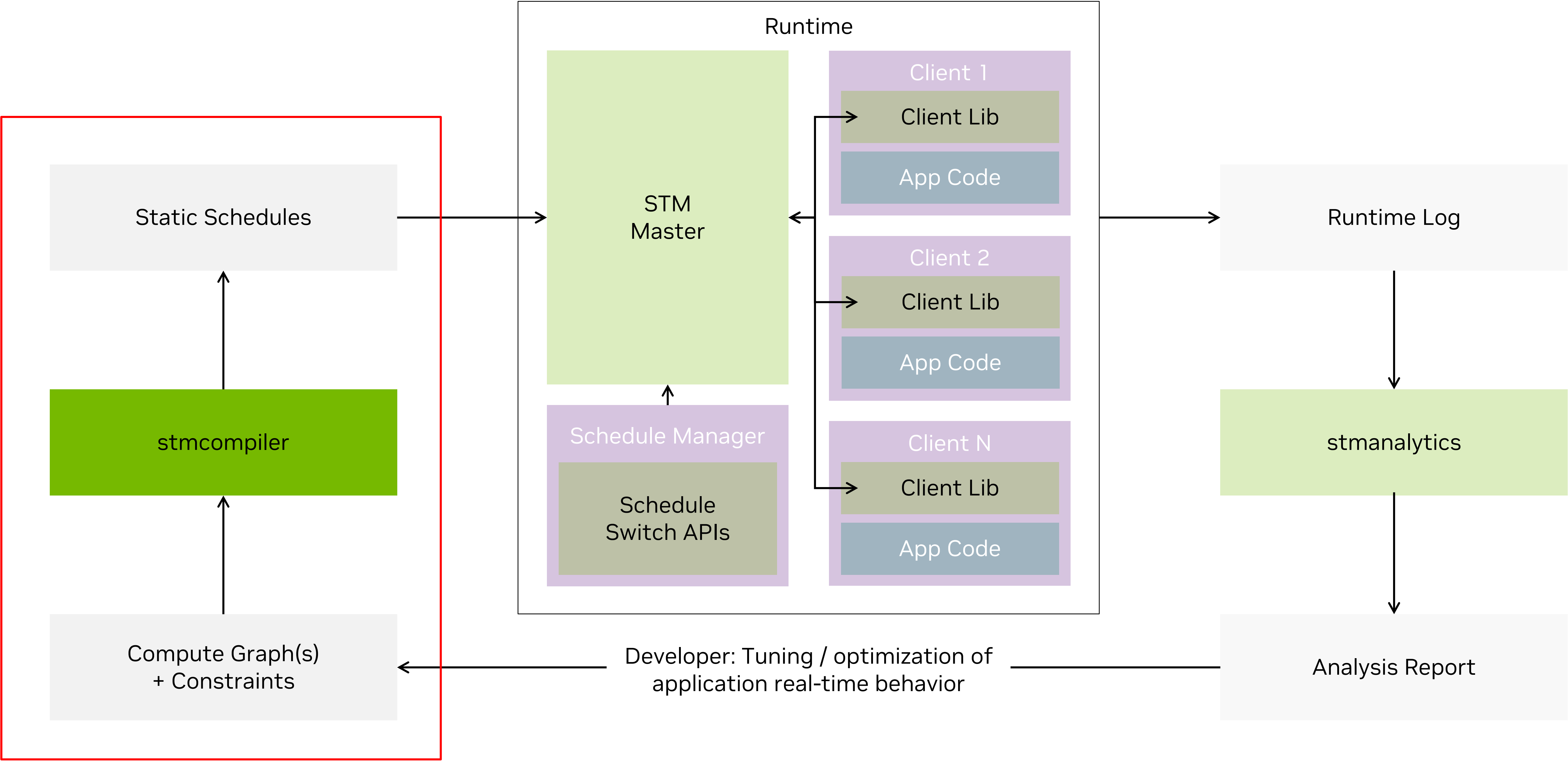
Framework



■ Tools / Binaries / APIs provided by STM ■ Applications created by developer that utilize STM

Generating Static Schedules

STM Compiler



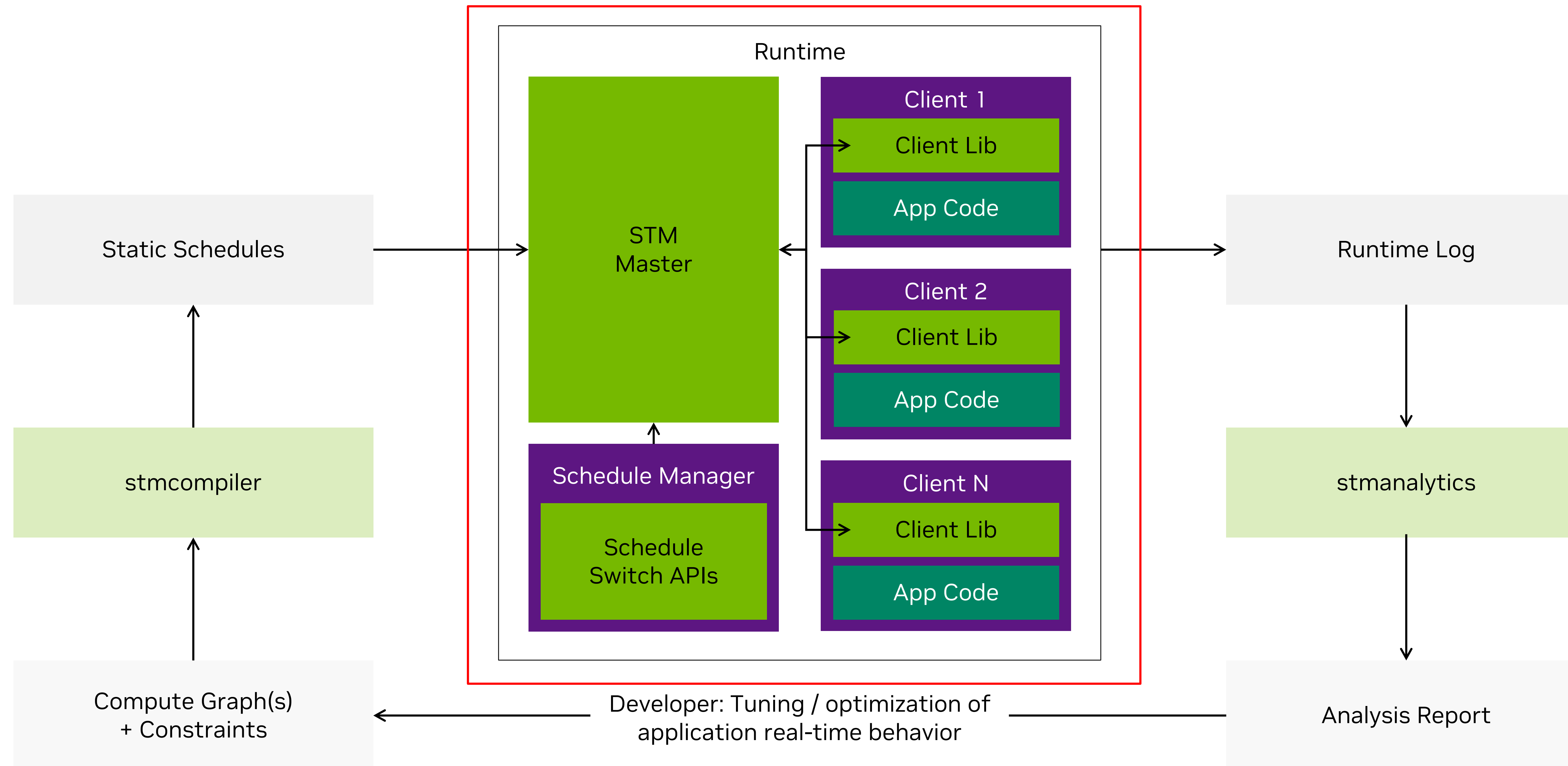
Tools / Binaries / APIs provided by STM



Applications created by developer that utilize STM

Executing Static Schedules

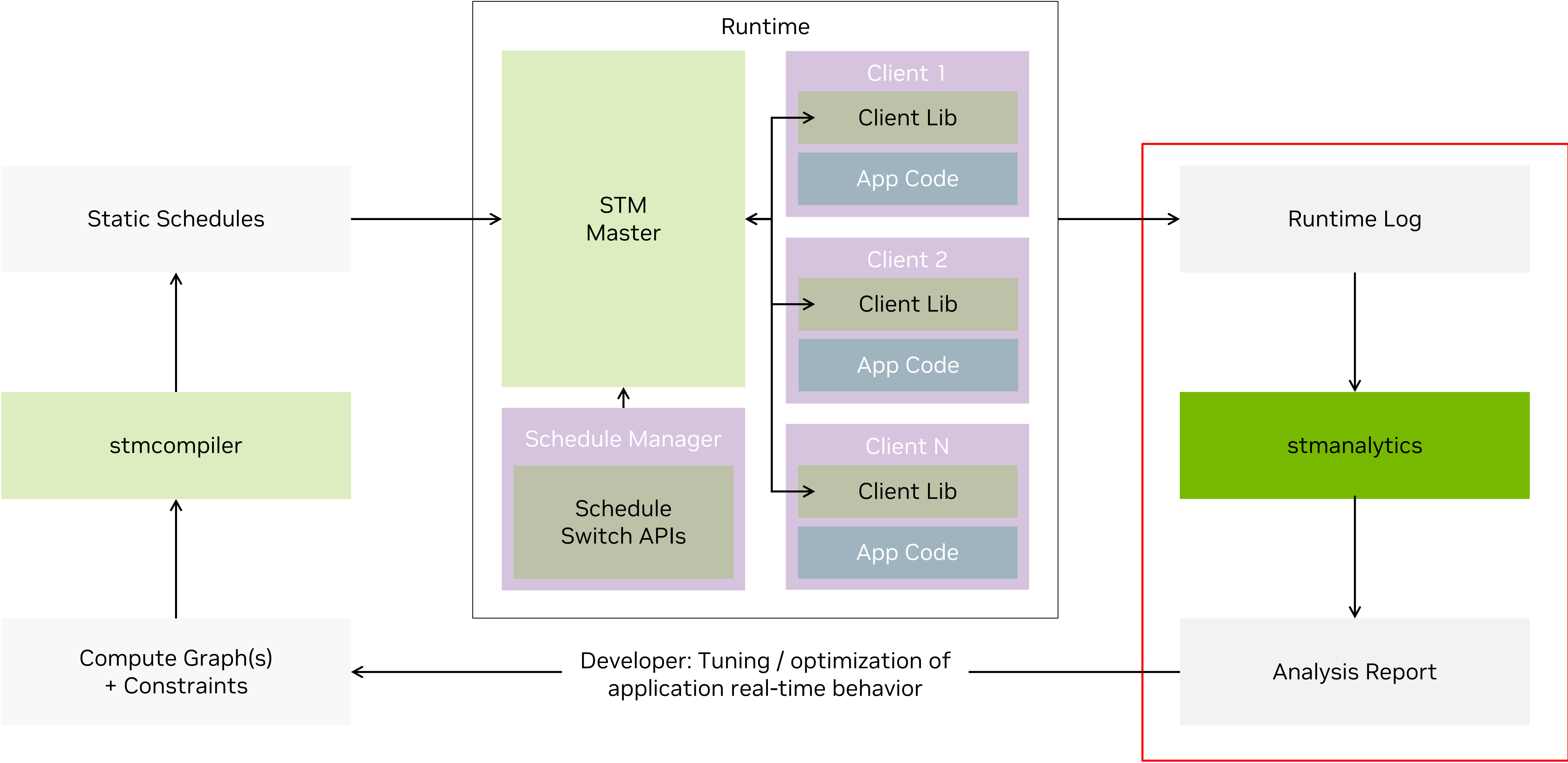
STM Runtime



■ Tools / Binaries / APIs provided by STM ■ Applications created by developer that utilize STM

Analyzing Execution Metrics

STM Analytics



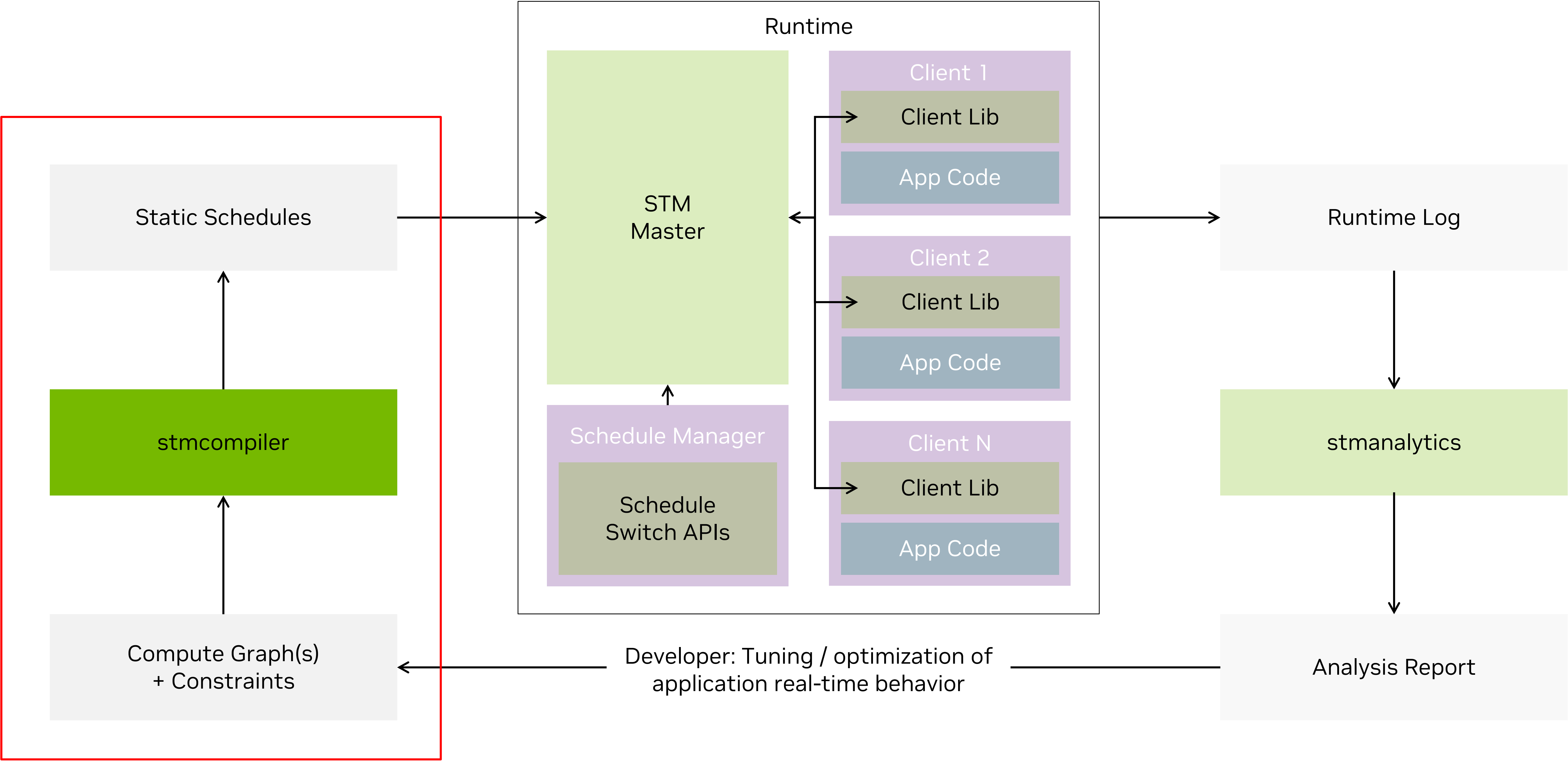
Tools / Binaries / APIs provided by STM



Applications created by developer that utilize STM

Generating Static Schedules

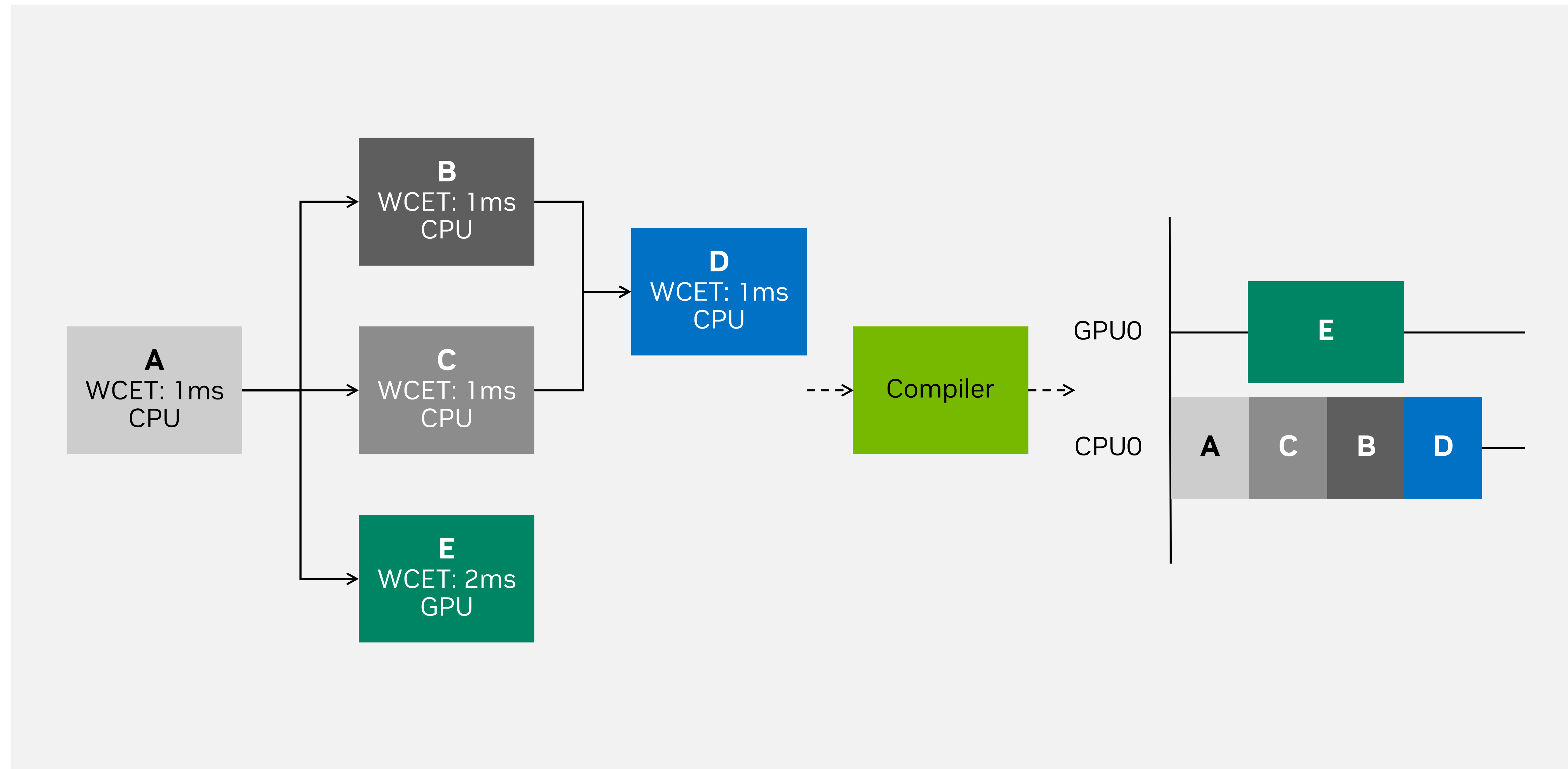
STM Compiler



■ Tools / Binaries / APIs provided by STM ■ Applications created by developer that utilize STM

Schedule Compilation

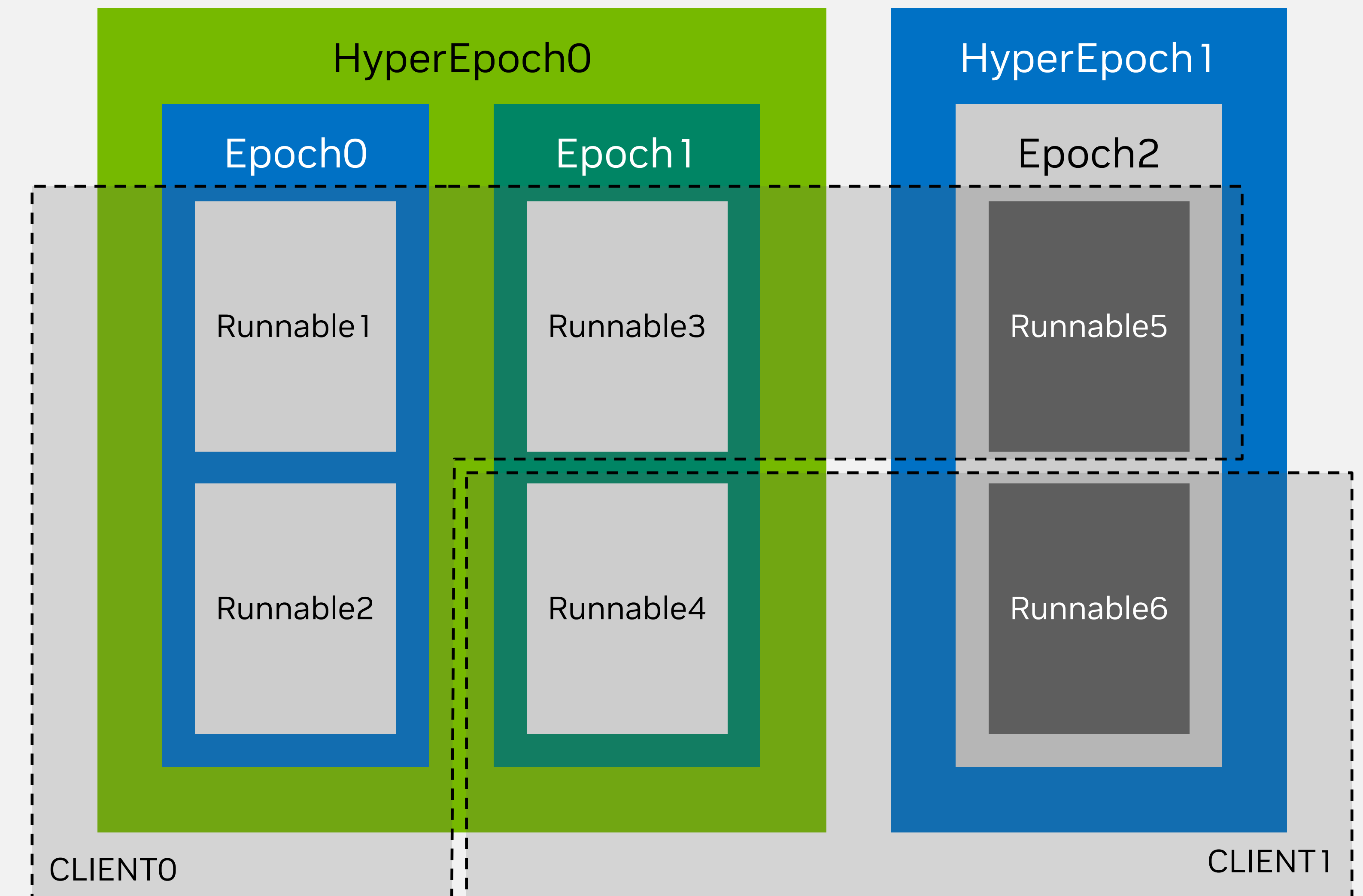
Converting requirements into static, recurring schedule



*WCET: Worst Case Execution Time

Scheduling Primitives

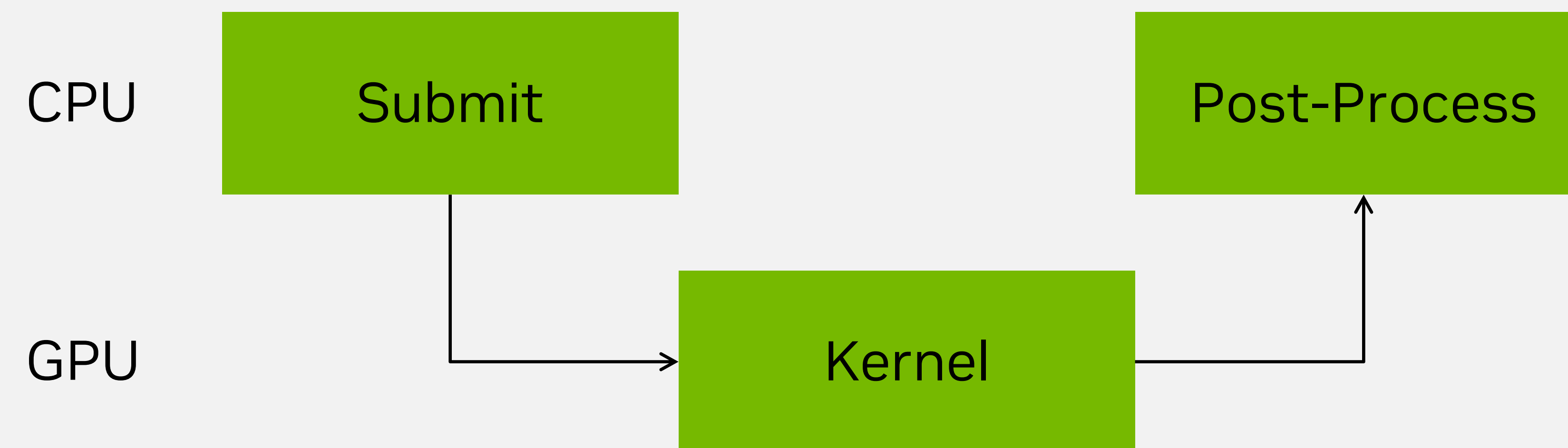
- Runnable
- HyperEpoch
- Epoch
- Client
- WCET
- Dependency
- Resource
- Refer [Documentation](#) for details



Boundaries in NVIDIA DRIVE Scheduler

Tutorial

Simple GPU application with pre and post processing on CPU



Questions

- How many runnables?
- How many hyperepochs?
- How many epochs?
- How many clients?
- What are the wcets for each runnable?
- What are the dependencies between the runnables?

Graph Specifications, Global Resources and HyperEpochs

cpu_gpu_single_process.yml

```
---
Version: 3.0.0          # Version number of this input yaml format
cpu_gpu_single_process: # Name of the DAG
  Identifier: 101        # Unique identifier of this schedule. Used when there are multiple schedules
  Resources:            # Global resources available in the system
    CPU: [CPU0]         # We define single CPU
    GPU: [iGPU]         # We define single GPU
  Hyperepochs:          # Hyperepoch definitions
    - hyperepoch0:      # We define single hyperepoch
      Period: 100ms     # The period for this hyperepoch is 100ms
      Epochs:           # Epoch definitions
        - epoch0:       # We define single epoch
```


Clients — Runnables

cpu_gpu_single_process.yml

```
Clients:
- GpuX:
  Resources:
    CUDA_STREAM:
      - CUDA_STREAMX: iGPU
  Epochs:
    - hyperepoch0.epoch0:
      Runnables:
        - submit:
          WCET: 10ms
          Dependencies: []
          Submits: GpuX.kernel
          Resources: [CPU0, CUDA_STREAM]
        - kernel:
          WCET: 10ms
          Dependencies: []
          Resources: [iGPU]
        - post_process:
          WCET: 10ms
          Dependencies: [GpuX.kernel]
          Resources: [CPU0]
```

Client definitions
We define single client "GpuX"
Resource definitions for this client
Cuda Streams used by "GpuX"
We define single cuda stream and map it to "iGPU"
Epochs definitions for "GpuX"
Epoch 0 definition for "GpuX"
Runnables present in "GpuX"
We define a pre-process runnable "submit"
We define a wcet of 10 ms for "submit"
"submit" does not depend on any other runnable
"submit" submits another runnable "kernel"
"submit" uses resources "CPU0" and "CUDA_STREAM"
We define the GPU runnable "kernel"
We define a wcet of 10 ms for "kernel"
"kernel" does not depend on any other runnable
"kernel" implicitly depends on "submit"
"kernel" requires "iGPU" to run
we define another runnable "post_process"
we define a wcet of 10 ms for "post_process"
"post_process" depends on "GpuX.kernel"
"post_process" requires "CPU0" to run

STM Compiler Usage

```
stmcompiler.py -i cpu_gpu_single_process.yml
```

Produces compiled schedule in cpu_gpu_single_process.stm


```
[webinar] 0: bash*
```

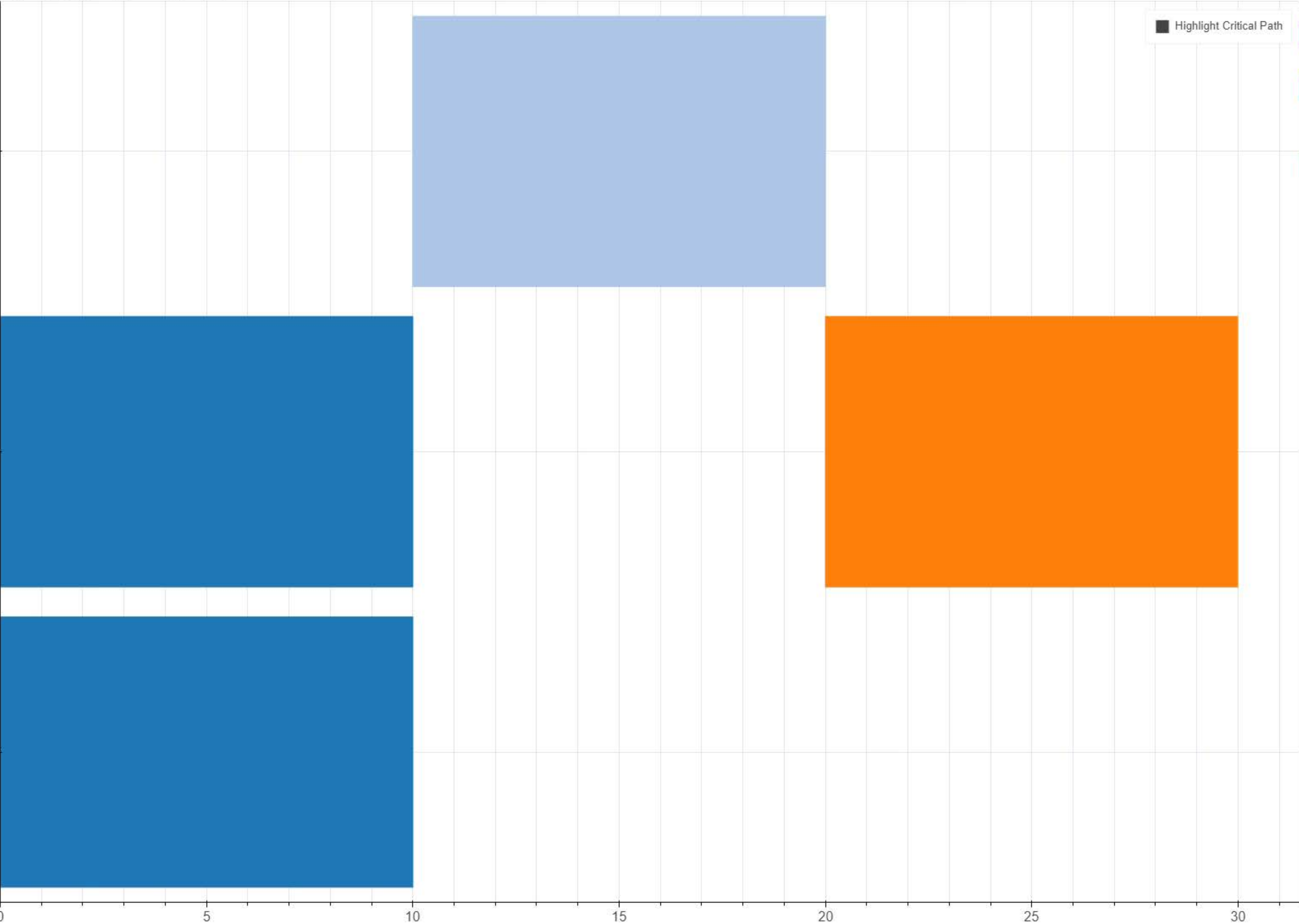
cpu_gpu_single_process.stm

Resource

default.iGPU

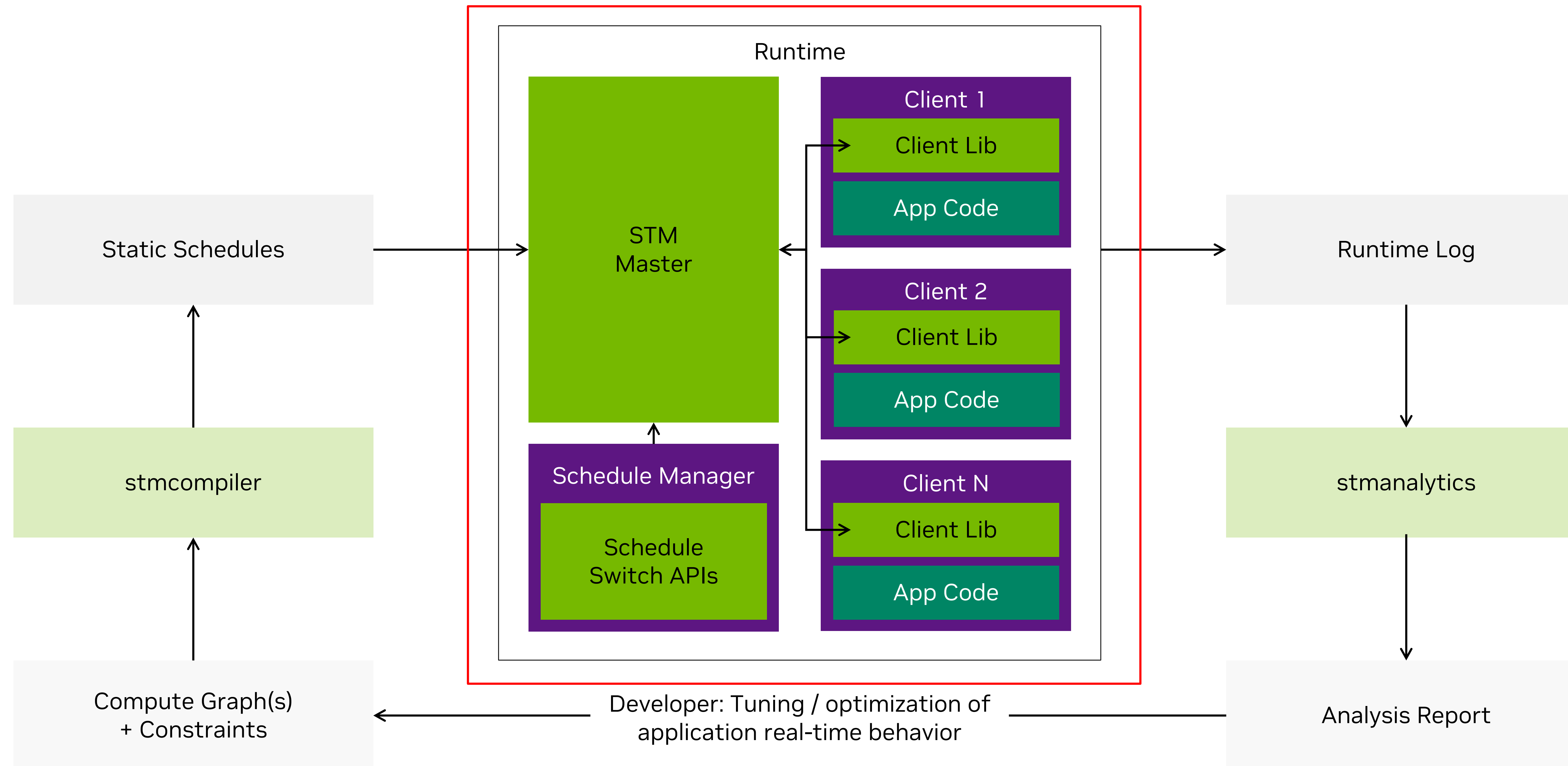
default.CPU0

GpuX.default.CUDA_STREAMX



Executing Static Schedules

STM Runtime



■ Tools / Binaries / APIs provided by STM ■ Applications created by developer that utilize STM

Runtime APIs

Init and deinit

`void stmClientInit(const char* clientName);`

- First call
- Initializes STM state
- Blocks for communication with stm_master

`void stmClientExit(void);`

- Cleans up STM state
- Last call

```
int main(int argc, const char** argv)
{
    (void)argc;
    (void)argv;

    cudaError_t cuErr = cudaStreamCreateWithFlags(&m_stream,
    cudaStreamNonBlocking); assert(cuErr == cudaSuccess);

    // Init
    stmClientInit("GpuX");

    // Register runnables
    stmRegisterCudaSubmitter(submit, "submit", NULL);
    stmRegisterCpuRunnable(post_process, "post_process", NULL);

    // Register resources
    stmRegisterCudaResource("CUDA_STREAMX", m_stream);

    // Execute
    stmErrorCode_t stmErr = stmEnterScheduler();
    assert(stmErr == STM_SUCCESS);

    // Cleanup
    stmClientExit();
}
```

Runtime APIs

Registration

```
stmErrorCode_t stmRegisterCpuRunnable (stmRunnable_t  
func, const char* const runnableId, void* userdata);
```

```
stmErrorCode_t  
stmRegisterCudaSubmitter(stmCudaSubmitter_t func,  
const char* const runnableId, void* userdata);
```

- Register function pointers and map them to string names
- Last param is arg to pass while invoking the function pointers
- Names must be same as in the yaml

```
stmErrorCode_t stmRegisterCudaResource(const char*  
resourceName, cudaStream_t stream);
```

- Similarly, map created resources to string names
 - Names must be same as in the yaml
- + Similar for many other engines!

```
int main(int argc, const char** argv)
{
    (void)argc;
    (void)argv;

    cudaError_t cuErr = cudaStreamCreateWithFlags(&m_stream,
    cudaStreamNonBlocking); assert(cuErr == cudaSuccess);

    // Init
    stmClientInit("GpuX");

    // Register runnables
    stmRegisterCudaSubmitter(submit, "submit", NULL);
    stmRegisterCpuRunnable(post_process, "post_process", NULL);

    // Register resources
    stmRegisterCudaResource("CUDA_STREAMX", m_stream);

    // Execute
    stmErrorCode_t stmErr = stmEnterScheduler();
    assert(stmErr == STM_SUCCESS);

    // Cleanup
    stmClientExit();
}
```

Runtime APIs

Execution

`stmErrorCode_t stmEnterScheduler(void);`

- Main function within which execution happens
- Blocking call
- To exit, either send SIGTERM to stm_master, or an asynchronous call to stmExitScheduler() from within the same process, or specify max number of hyperepoch frames as param to stm_master

```
int main(int argc, const char** argv)
{
    (void)argc;
    (void)argv;

    cudaError_t cuErr = cudaStreamCreateWithFlags(&m_stream,
    cudaStreamNonBlocking); assert(cuErr == cudaSuccess);

    // Init
    stmClientInit("GpuX");

    // Register runnables
    stmRegisterCudaSubmitter(submit, "submit", NULL);
    stmRegisterCpuRunnable(post_process, "post_process", NULL);

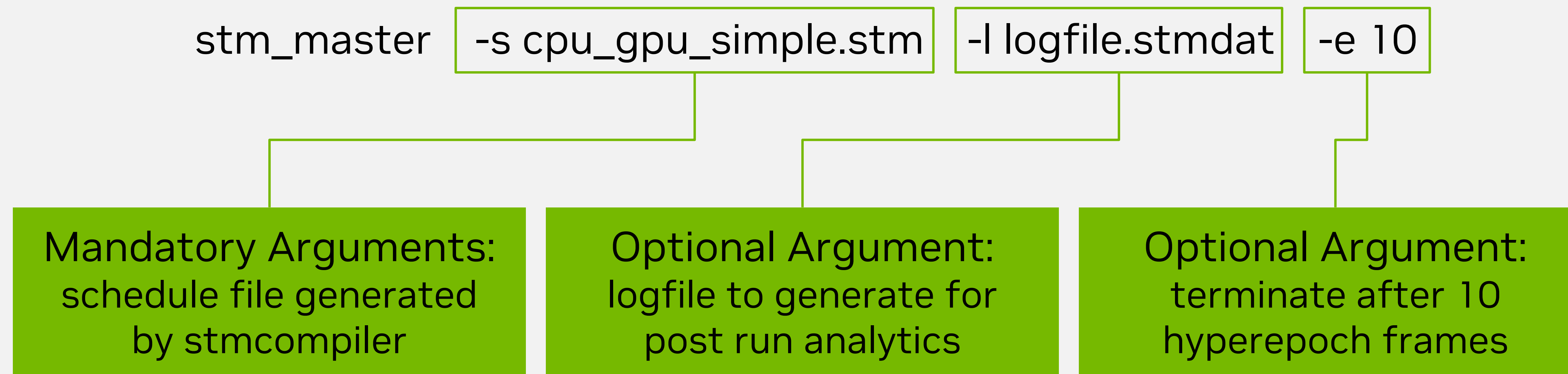
    // Register resources
    stmRegisterCudaResource("CUDA_STREAMX", m_stream);

    // Execute
    stmErrorCode_t stmErr = stmEnterScheduler();
    assert(stmErr == STM_SUCCESS);

    // Cleanup
    stmClientExit();
}
```

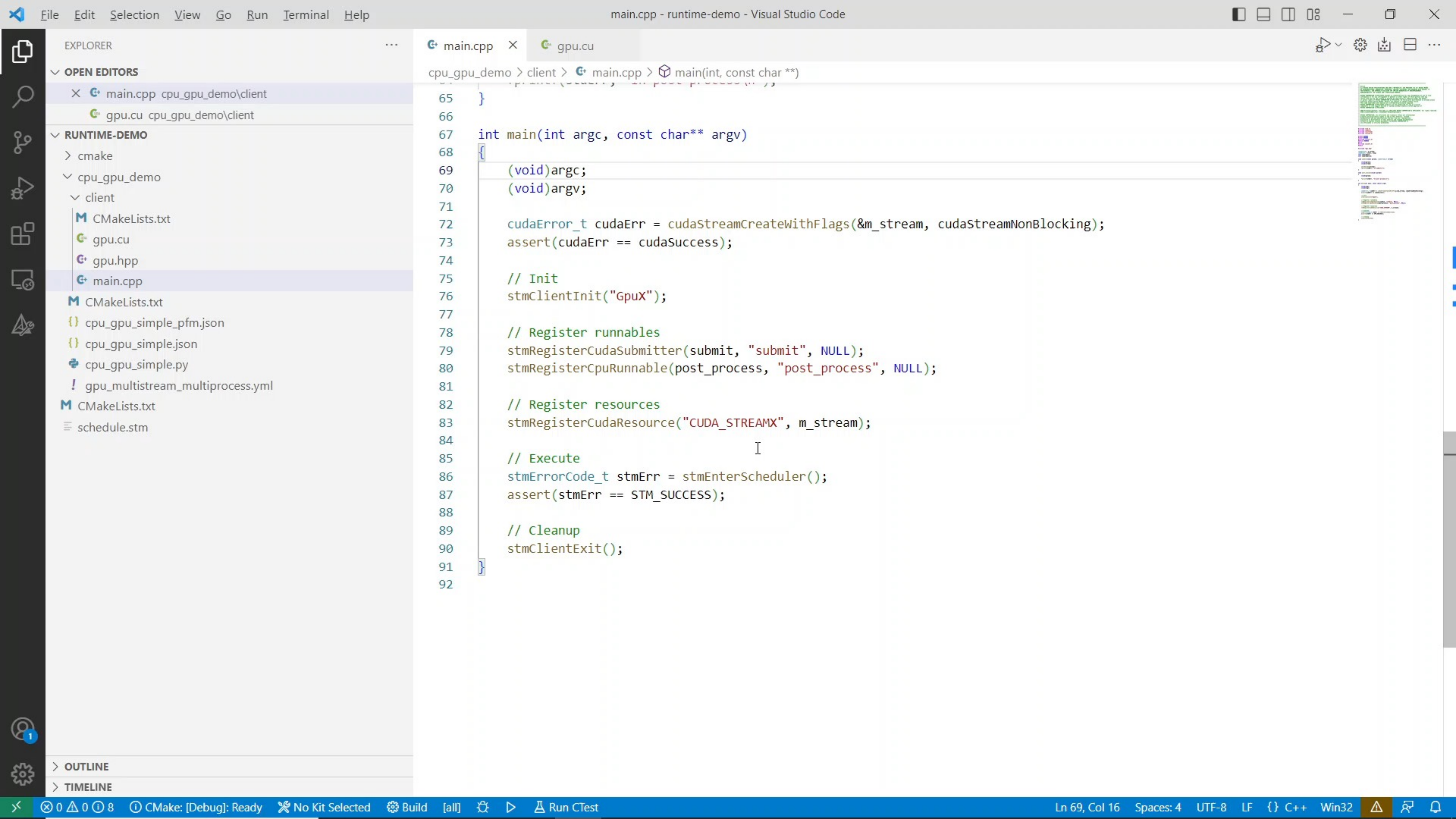

Runtime Execution

Execution of the application



`stm_client:`

No arguments since this is user created and we created one without any arguments required



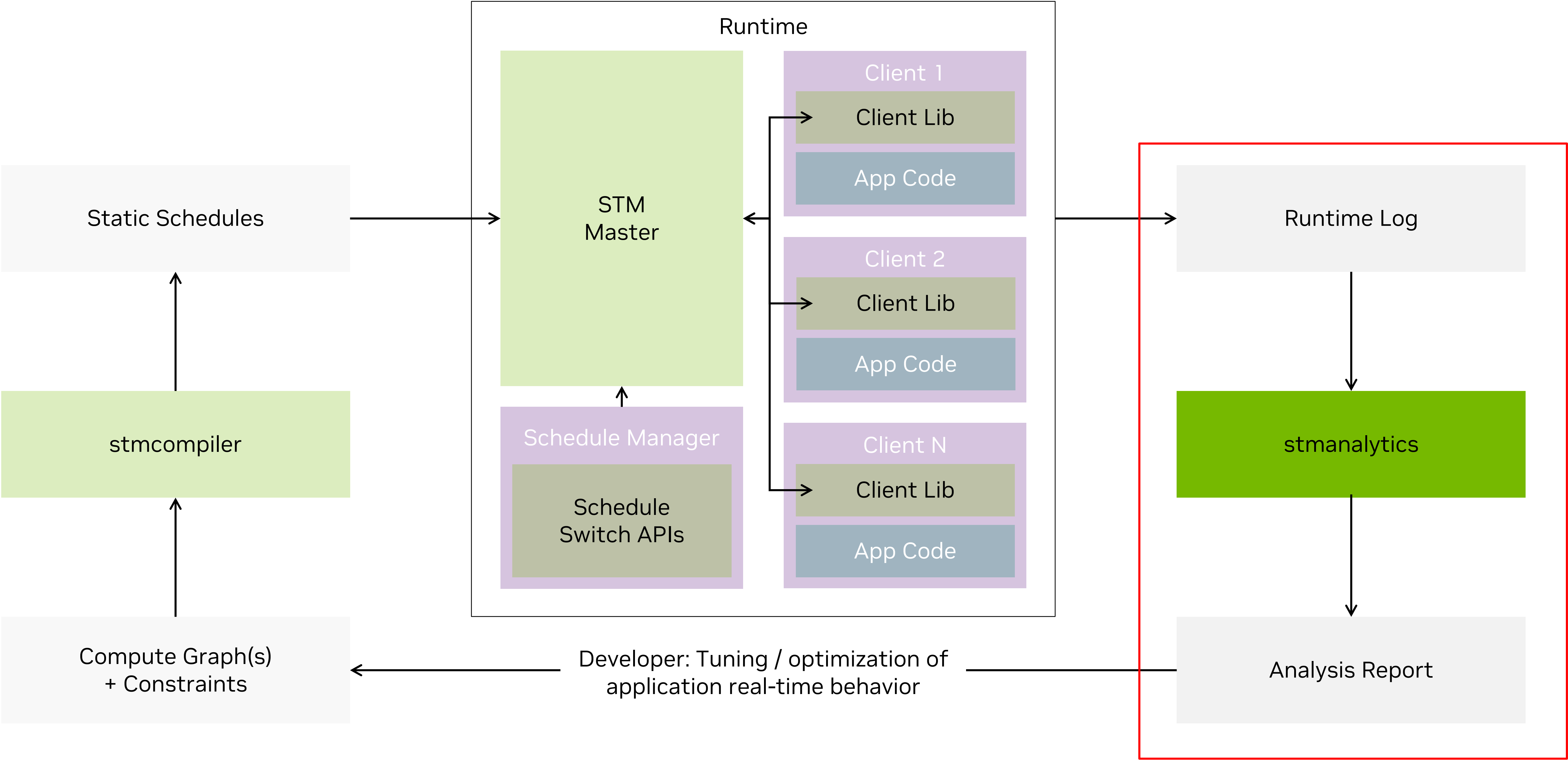
```
main.cpp - runtime-demo - Visual Studio Code

cpu_gpu_demo > client > main.cpp > main(int, const char **)

65 }
66
67 int main(int argc, const char** argv)
68 {
69     (void)argc;
70     (void)argv;
71
72     cudaError_t cudaErr = cudaStreamCreateWithFlags(&m_stream, cudaStreamNonBlocking);
73     assert(cudaErr == cudaSuccess);
74
75     // Init
76     stmClientInit("GpuX");
77
78     // Register runnables
79     stmRegisterCudaSubmitter(submit, "submit", NULL);
80     stmRegisterCpuRunnable(post_process, "post_process", NULL);
81
82     // Register resources
83     stmRegisterCudaResource("CUDA_STREAMX", m_stream);
84
85     // Execute
86     stmErrorCode_t stmErr = stmEnterScheduler();
87     assert(stmErr == STM_SUCCESS);
88
89     // Cleanup
90     stmClientExit();
91 }
92
```


Analyzing Execution Metrics

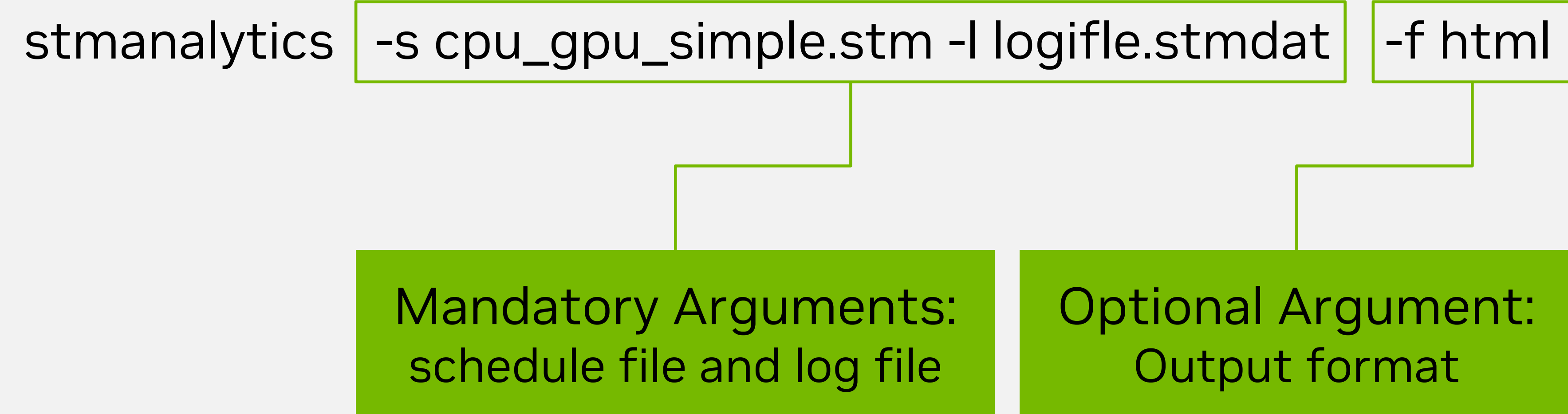
STM Analytics



Analytics

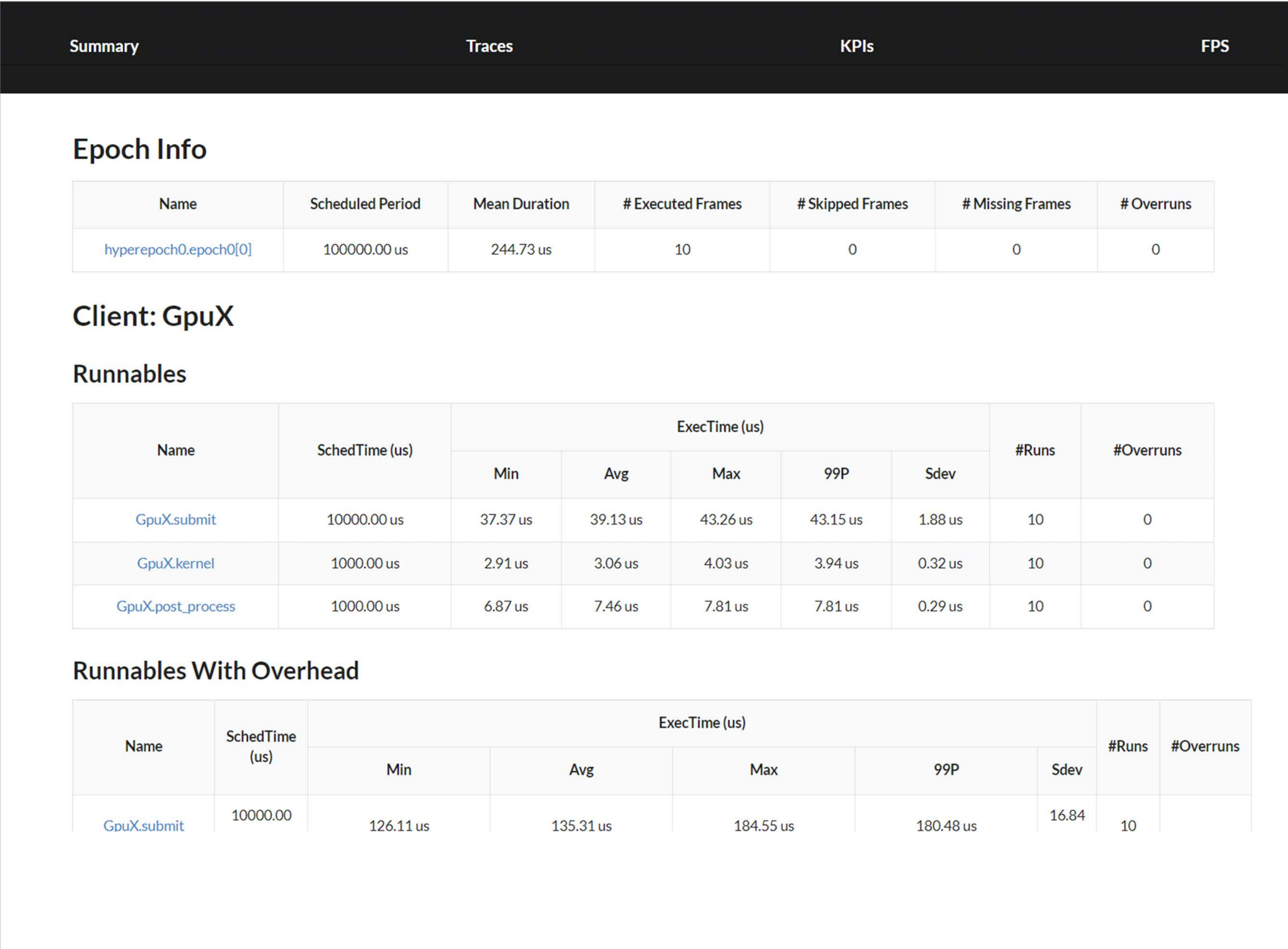
- Post run info on how long each runnable and the entire frame took
- Helpful to get insight on what engines are being over / under budgeted and direct dev work to those
- Also provides info on overhead due to scheduling via STM when possible
- Necessary for getting accurate WCETs to provide as input to stmcompiler

Analytics



x1@NV-3Q5KNN3: ~\$

Analytics Example



Summary Information



Per-Runnable Information

Get Started with DRIVE SDK

Extensive documentation and training material available on NVIDIA Developer

Learn More

- Visit the [DRIVE training](#) page for webinars and other resources
- Check out information related to [DRIVE Hyperion](#), [DRIVE AGX Orin](#) and [DRIVE SDK](#)

Get Access

- Join the [DRIVE AGX SDK Program](#) on NVIDIA Developer
- [Read the docs](#) for DRIVE OS and DriveWorks documentation
- [Download DRIVE OS](#) which includes DriveWorks, NvMedia, CUDA, cuDNN and TensorRT

Contact Us

- Contact your distributor or NVIDIA's [Automotive Team](#)

