

Profiling and data access patterns in Allen

Douglas Li

Stanford University

October 2021

Table of Contents

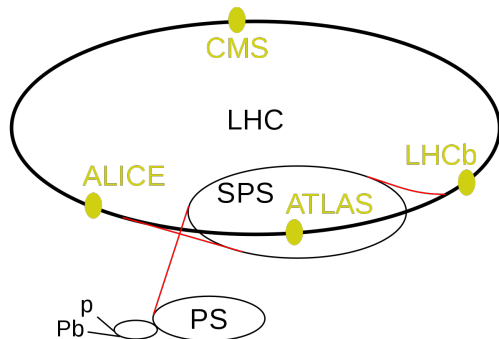
- 1 A quick introduction to LHCb
- 2 Allen and GPUs
- 3 Profiling and data access patterns

Table of Contents

- 1 A quick introduction to LHCb
- 2 Allen and GPUs
- 3 Profiling and data access patterns

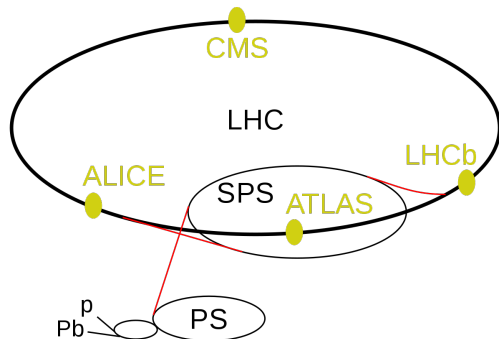
LHCb detector

- LHCb is one of four large detectors at the Large Hadron Collider (LHC)



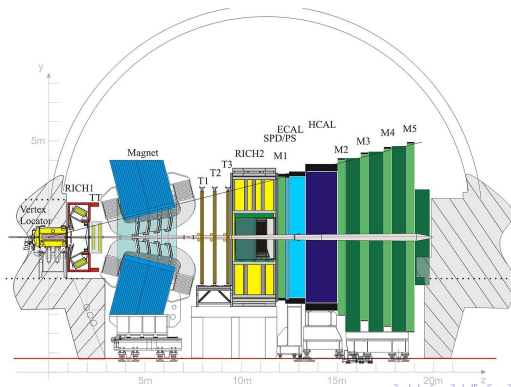
LHCb detector

- LHCb is one of four large detectors at the Large Hadron Collider (LHC)
- The LHC collides high energy beams of protons into each other to study the particles produced from such collisions.



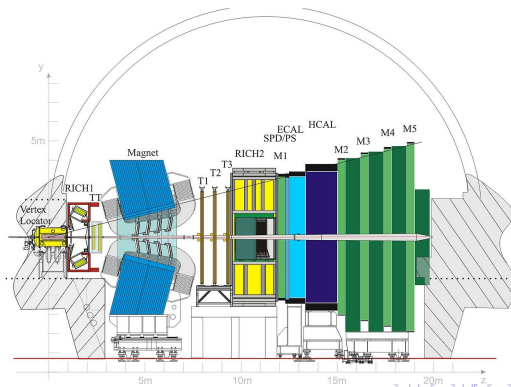
LHCb upgrades

- The LHCb is currently shut down for upgrades before Run 3 of data taking. Importantly, Run 3 will not use a hardware trigger.



LHCb upgrades

- The LHCb is currently shut down for upgrades before Run 3 of data taking. Importantly, Run 3 will not use a hardware trigger.
- Only a small fraction of the data read out from subdetectors can be kept for study: the trigger selects which events to keep.



LHCb upgrades

- The LHCb is currently shut down for upgrades before Run 3 of data taking. Importantly, Run 3 will not use a hardware trigger.
- Only a small fraction of the data read out from subdetectors can be kept for study: the trigger selects which events to keep.
- Without a hardware trigger, the software trigger must be redesigned to handle the full event rate of 30 MHz.

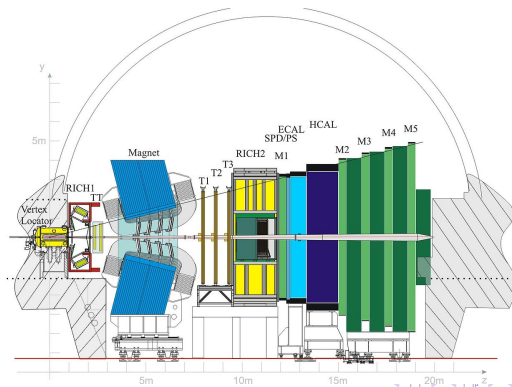


Table of Contents

- 1 A quick introduction to LHCb
- 2 Allen and GPUs**
- 3 Profiling and data access patterns

What is Allen?

The screenshot shows the GitHub repository page for 'Allen' by LHCb. The repository is described as 'Full software HLT1 reconstruction sequence on GPU.' It has 5,717 commits, 605 branches, 16 tags, 6.3 MB files, 330.8 GB storage, and 16 releases. A recent merge commit by Christoph Hasse is shown, along with a table of files and their last commit details.

LHCb > Allen

Allen Project ID: 38633 🔔 ☆ Star 31 🍴 Fork 18

↔ 5,717 Commits 🌿 605 Branches 🏷️ 16 Tags 📄 6.3 MB Files 💾 330.8 GB Storage 📅 16 Releases

Full software HLT1 reconstruction sequence on GPU.

master Allen / + History Find file Web IDE 📄 Clone

Merge branch 'dcampora_default_to_building_single_sequence' into 'master' ... a613e2f3

Christoph Hasse authored 10 hours ago

📄 README 📄 Apache License 2.0 📄 CONTRIBUTING 📄 CI/CD configuration

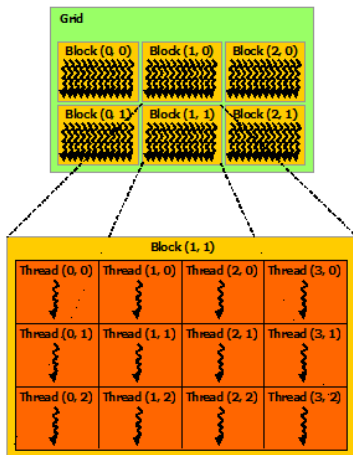
Name	Last commit	Last update
📄 Dumpers	Electron lines	3 days ago
📄 Rec/Allen	Electron lines	3 days ago
📄 ReleaseNotes	Release notes and dependencies for Allen v1...	1 week ago
📄 backend	Drop CUDAFLAG target.	1 month ago

What is Allen?

- "HLT1 reconstruction sequence": Stage 1 of software trigger. A sequence of algorithms that does the following tasks:
 - "Hits" are identified in each of the VELO, UT, and SciFi detectors.
 - (Each detector has several layers, and as particles pass through the layers, they leave signals in each layer. These signals are hits.)
 - These hits are combined to form tracks (trajectories of a particle).
 - The tracks are then used to reconstruct *primary vertices* (p-p collision points) and *secondary vertices* (decay points)
- "GPUs": Graphical Processing Units

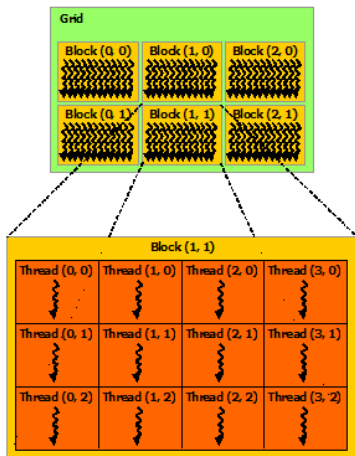
The layout of a GPU

- Grid-block-thread hierarchy.



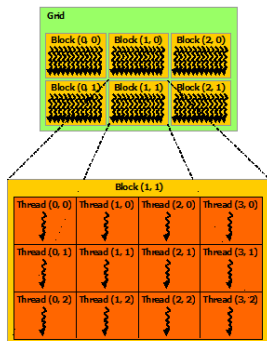
The layout of a GPU

- Grid-block-thread hierarchy.
- In CUDA, functions are executed N times by N different threads. Threads are grouped into thread blocks which are grouped into grids.



The layout of a GPU

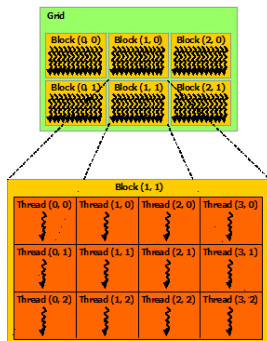
- We can take advantage of this hierarchy as follows. Events can be handled by thread blocks, and individual data in each event (hits/tracks) can be handled by threads.



CUDA C++ Programming Guide

The layout of a GPU

- We can take advantage of this hierarchy as follows. Events can be handled by thread blocks, and individual data in each event (hits/tracks) can be handled by threads.
- So each event can be processed independently, and the hits/tracks within that event can also be processed in parallel.



CUDA C++ Programming Guide

Table of Contents

- 1 A quick introduction to LHCb
- 2 Allen and GPUs
- 3 Profiling and data access patterns**

Profiling Allen

- We can measure the performance of each algorithm using the profiling tool NVIDIA Nsight Compute.

The screenshot displays the NVIDIA Nsight Compute interface for a kernel named '616 - fit_secondary_vertices'. The top navigation bar includes options for 'Page: Details', 'Launch: 47 - 616 - fit_secondary_vertices', and 'Add Baseline'. The main header shows 'Current' launch details: '616 - fit_secondary_vertices (935, 1, 1)x(16, 16, 1)' with a time of 358.94 usecond, 498,246 cycles, 64 registers, and a GeForce RTX 2080 Ti GPU. Performance metrics include SM Frequency (1.39 cycle/usecond) and CC Process (7.5 [133860] Allen).

The 'GPU Speed Of Light' section provides a high-level overview of GPU utilization. It includes a table of metrics:

Metric	Value	Duration [usecond]	Value
SOL SM [%]	11.43	Duration [usecond]	358.94
SOL Memory [%]	34.48	Elapsed Cycles [cycle]	498,246
SOL L1/TEX Cache [%]	36.48	SM Active Cycles [cycle]	398,031.51
SOL L2 Cache [%]	29.87	SM Frequency [cycle/usecond]	1.39
SOL DRAM [%]	34.48	DRAM Frequency [cycle/usecond]	6.91

A 'Bottleneck' warning is present, indicating low compute throughput and memory bandwidth utilization relative to the peak performance of the device. It suggests looking at [Scheduler Statistics](#) and [Warp State Statistics](#) for potential reasons.

The 'Compute Workload Analysis' section details the analysis of the compute resources of the streaming multiprocessors (SM). It includes a table of metrics:

Metric	Value	SM Busy [%]	Value
Executed Ipc Elapsed [inst/cycle]	0.19	SM Busy [%]	6.01
Executed Ipc Active [inst/cycle]	0.24	Issue Slots Busy [%]	6.01
Issued Ipc Active [inst/cycle]	0.24		

A 'High Pipe Utilization' warning indicates that all pipelines are under-utilized, suggesting the kernel is very small or doesn't issue enough warps per scheduler.

The 'Memory Workload Analysis' section details the analysis of the memory resources of the GPU. It includes a table of metrics:

Metric	Value	Mem Busy [%]	Value
Memory Throughput [Gbyte/second]	228.67	Mem Busy [%]	24.46
L1/TEX Hit Rate [%]	72.61	Max Bandwidth [%]	34.48
L2 Hit Rate [%]	93.19	Mem Pipes Busy [%]	11.43

The 'Scheduler Statistics' section provides a summary of the activity of the schedulers issuing instructions. It notes that the upper bound of warps in the pool (Theoretical Warps) is limited by the launch configuration, and that the scheduler selects a single warp from which to issue one or more instructions (Issued Warp).

Warp Stalls

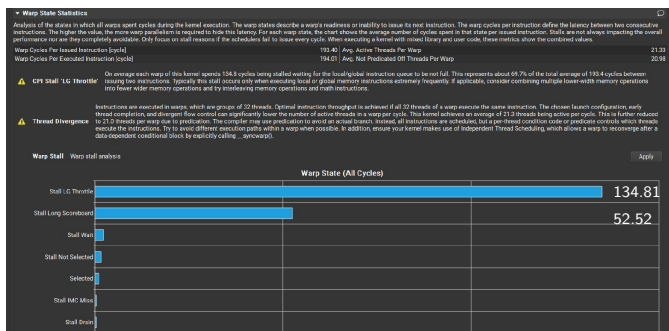
- The hardware partitions each thread block into *warps*, which are groups of 32 threads.

Warp Stalls

- The hardware partitions each thread block into *warps*, which are groups of 32 threads.
- A *warp scheduler* schedules warps for execution, but can stall for a variety of reasons.

Warp Stalls

Profiler indicates that latency is a limiting factor to performance in e.g. `fit_secondary_vertices` algorithm.



The x-axis measures the number of cycles per instruction a warp spends in a given stall state. We see that the LG Throttle stall, with 134.81 cycles/instruction, is the bottleneck, and that this is related to frequent accesses to local or global memory. So this hints that we may want to look at how our code accesses memory.

Memory Coalescing

- The Source Counters section also indicates uncoalesced global accesses, again suggesting memory accesses as a problem.

The screenshot shows the 'Source Counters' section of a profiler. It displays a table of metrics and a list of uncoalesced global accesses. The table includes 'Branch Instructions [inst]' (121,248) and 'Branch Instructions Ratio [%]' (0.03). The list of accesses shows 10 entries, each with a warning icon, the text 'Uncoalesced Global Accesses', and details about the access, including the expected sectors, the number of sectors actually accessed (3.64x), and the PC address.

Source metrics, including branch efficiency and sampled warp stall reasons. Sampling Data metrics are periodically sampled over the kernel runtime. They indicate when warps were stalled and couldn't be scheduled. See the documentation for a description of all stall reasons. Only focus on stalls if the schedulers fail to issue every cycle.

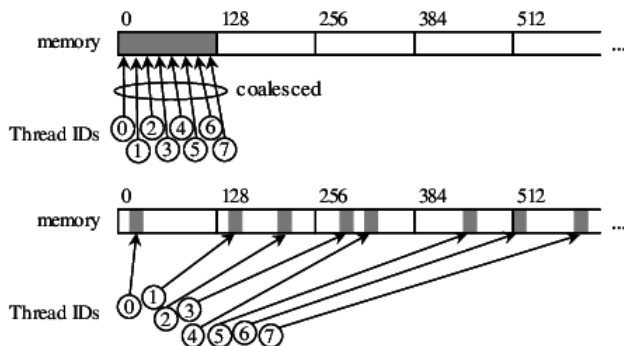
Branch Instructions [inst]	121,248	Branch Efficiency [%]	98.05
Branch Instructions Ratio [%]	0.03	Avg. Divergent Branches	6.91

- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda95a0](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/src/VertexFitter.cu:70
- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda95a0](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/src/VertexFitter.cu:71
- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda9a00](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/include/VertexFitter.icc:39
- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda9a00](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/include/VertexFitter.icc:38
- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda9a00](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/include/VertexFitter.icc:37
- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda9a00](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/include/VertexFitter.icc:30
- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda9a00](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/include/VertexFitter.icc:84
- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda9a00](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/include/VertexFitter.icc:170
- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda9a00](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/include/VertexFitter.icc:170
- ⚠ **Uncoalesced Global Accesses** Uncoalesced global access, expected 27896 sectors, got 101552 (3.64x) at PC [0x707e6dda9a00](#) at /home/doli/test/Allen/device/vertex_fit/vertex_fitter/include/VertexFitter.icc:170

PC sampling data PC sampling data Apply

Memory Coalescing

- Threads can only access global memory in *fixed-size transactions*, which are 32 bytes in size.



Keskin, Cetin, Kocak 2015

Here we can think of each threadID as getting the data for some track or hit. Then if the data for tracks 1-8 is coalesced, only one global access is needed. If not, then multiple global accesses (in this case 5) are needed to retrieve the data. This leads to more accesses so warps will stall waiting in the global instruction queue.

Struct of Arrays (SOA) vs. Array of Structs (AOS)

- The question of whether to define data structures as structs of arrays or arrays of structs appears frequently in parallel programming.
- Array of structs:

```
1 struct coords {  
2     float x;  
3     float y;  
4     float z;  
5 }  
6  
7 coords* array_of_coords;  
8
```

Struct of Arrays vs. Array of Structs

- The question of whether to define data structures as structs of arrays or arrays of structs appears frequently in parallel programming.
- Struct of arrays:

```
1struct CoordsArray {
2private:
3    const float * coords;
4    const unsigned num_hits;
5public:
6    // Constructor
7    __host__ __device__ CoordsArray(float * base_pointer, const unsigned offset, const unsigned total_number_of_hits) :
8    m_base_pointer(base_pointer + offset), num_hits(total_number_of_hits) {}
9
10   // Getter functions.
11   __host__ __device__ float x(unsigned index) const { return coords[index]; }
12   __host__ __device__ float y(unsigned index) const { return coords[index + num_hits]; }
13   __host__ __device__ float z(unsigned index) const { return coords[index + 2 * num_hits]; }
14
15   // Setter functions.
16   __host__ __device__ void set_x(unsigned index, float value) { coords[index] = value; }
17   __host__ __device__ void set_y(unsigned index, float value) { coords[index + num_hits] = value; }
18   __host__ __device__ void set_z(unsigned index, float value) { coords[index + 2 * num_hits] = value; }
19}
20
```


- A general heuristic is to use SOA to keep global memory accesses coalesced.

- A general heuristic is to use SOA to keep global memory accesses coalesced.
- However, this intuition is often wrong, and there is no substitute for trying out both structures and comparing performance.

- A general heuristic is to use SOA to keep global memory accesses coalesced.
- However, this intuition is often wrong, and there is no substitute for trying out both structures and comparing performance.
- Also possible to use combinations of SOA and AOS. And different variables can be split up into different SOA's depending on where they are used.

- The fit_secondary_vertices algorithm uses an array of TrackMVAVertex structs (AOS) to store all the information needed for a secondary vertex fit.

```
// Secondary vertices.
VertexFit::TrackMVAVertex* event_secondary_vertices = parameters.dev_consolidated_svs + sv_offset;

// Loop over sv's.
for (unsigned i_sv = threadIdx.x; i_sv < n_svs; i_sv += blockDim.x) {
    event_secondary_vertices[i_sv].chi2 = -1;
    event_secondary_vertices[i_sv].minipchi2 = 0;
    auto i_track = event_svs_trk1_idx[i_sv];
    auto j_track = event_svs_trk2_idx[i_sv];
    const ParKalmanFilter::FittedTrack trackA = event_tracks[i_track];
    const ParKalmanFilter::FittedTrack trackB = event_tracks[j_track];

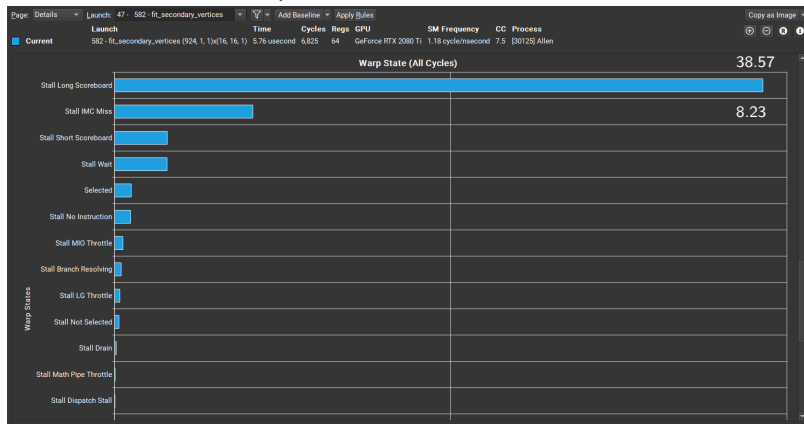
    // Do the fit.
    // TODO: In case doFit returns false, what should happen?
    if (doFit(trackA, trackB, event_secondary_vertices[i_sv])) {
        event_secondary_vertices[i_sv].trk1 = i_track;
        event_secondary_vertices[i_sv].trk2 = j_track;
    }
}
```

- We change this algorithm and the Kalman filtering algorithm to use multiple SOA's. We also cut down on global memory accesses by declaring data members in register memory.

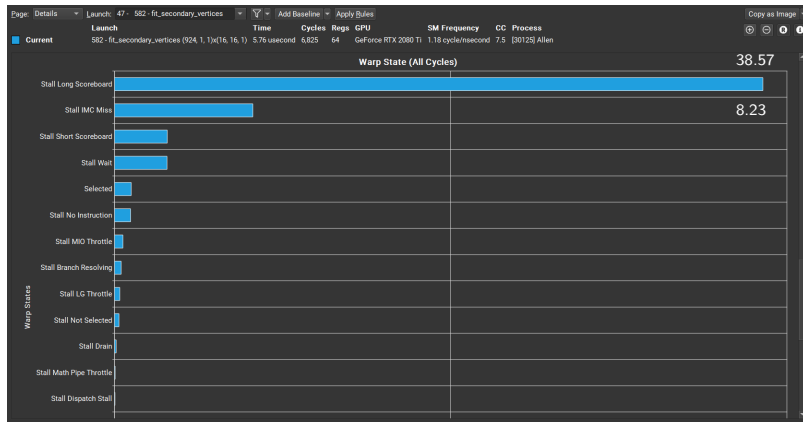
- We change this algorithm and the Kalman filtering algorithm to use multiple SOA's. We also cut down on global memory accesses by declaring data members in register memory.
- Does the number of warp stalls decrease?

fit_secondary_vertices

- We change this algorithm and the Kalman filtering algorithm to use multiple SOA's. We also cut down on global memory accesses by declaring data members in register memory.
- Does the number of warp stalls decrease?



- The LG Throttle stall is reduced to 0.32 cycles per instruction. The largest stall is Long Scoreboard, with 38.57 cycles per instruction.



- In the long run, this view-based model can be integrated into the entire pipeline. This is being worked on by Tom and Daniel.

- In the long run, this view-based model can be integrated into the entire pipeline. This is being worked on by Tom and Daniel.
- Ultimate test will be the effect on throughput.

- There are many optimizations that can be made to Allen.

Takeaways

- There are many optimizations that can be made to Allen.
- Profiling tools give a lot of information and can offer a guided look into speedup opportunities. However it is up to the programmer to decide which information is relevant and how this information can inform how we design the code.

Takeaways

- There are many optimizations that can be made to Allen.
- Profiling tools give a lot of information and can offer a guided look into speedup opportunities. However it is up to the programmer to decide which information is relevant and how this information can inform how we design the code.
- Getting maximum performance out of Allen will require carefully integrating the software with constraints of the hardware, specifically the CUDA programming model.

Acknowledgments

- Tom Boettcher and Daniel Campora

Acknowledgments

- Tom Boettcher and Daniel Campora
- University of Cincinnati and DIANA fellows program

- Any questions?