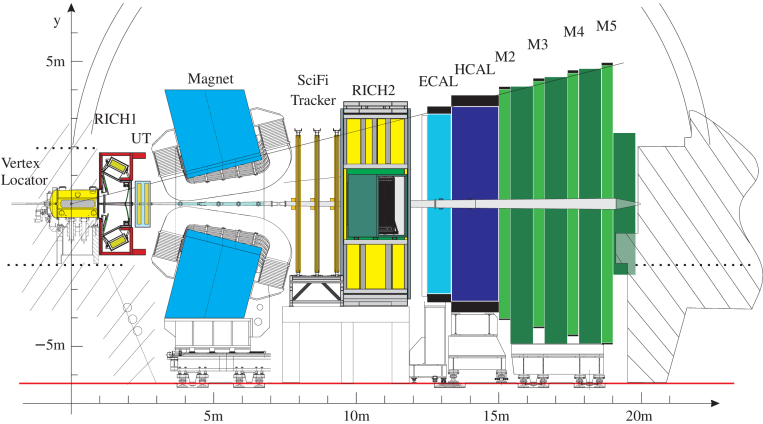
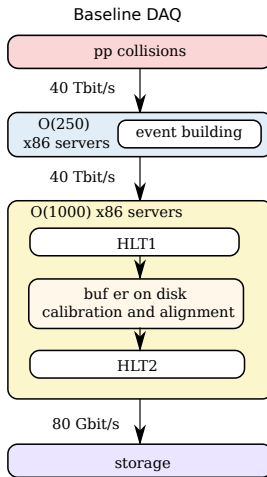
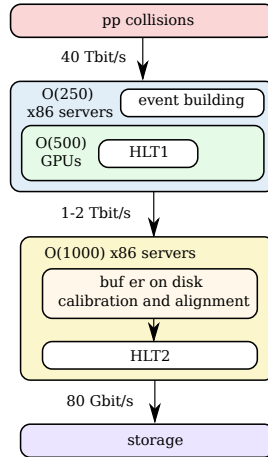
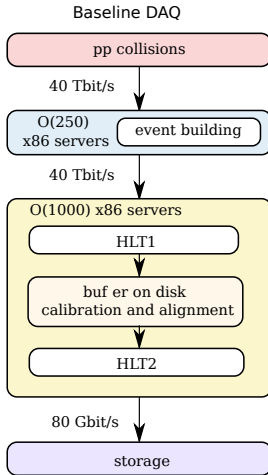


Experience in LHCb

Run 3 upgrade in LHCb



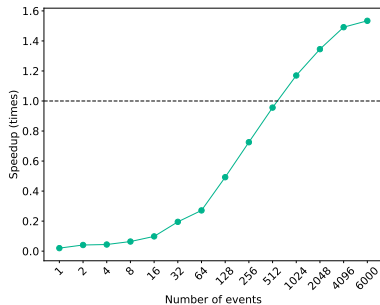




Running efficiently on GPUs

Previous attempts in LHCb were done using a service to offload part of the work to GPUs. Single events were submitted to an service, which accumulated and offloaded several of them. However, the LHCb event size is too small (60-100 kB), not enough to occupy efficiently the GPU.

The number of events is a key factor in LHCb workload:



LHCb GPU core infrastructure

The Allen framework

The *Allen* framework is a compact, scalable and modular framework, built for running the LHCb HLT1 on GPUs.

Requirements

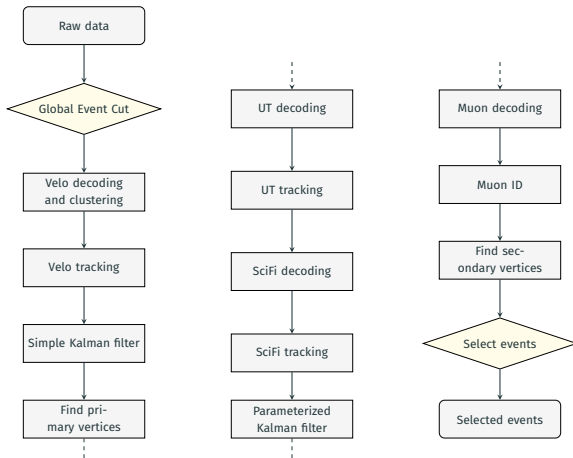
- A C++17 compliant compiler, boost, zeromq
- CUDA v10.0 (currently C++14 is supported)

Features

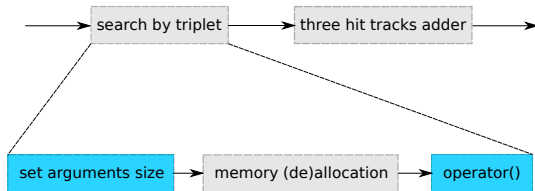
- Multithreaded, multistream framework.
- Configurable static sequences.
- Pipelined stream sequence.
- Custom memory manager, no dynamic allocations, SOA datatypes.
- Built-in validation with Monte Carlo.
- Optional compilation with ROOT for generation of graphs.
- Integration with Gaudi build system.
- Cross-architecture compatibility.



Current Allen HLT1 sequence



Algorithm sequence



An algorithm definition (1)

```
1 __global__ void velo_three_hit_tracks_adder(  
2     uint32_t* dev_VELO_cluster_container,  
3     uint* dev_module_cluster_start,  
4     Velo::TrackHits* dev_tracks,  
5     Velo::TrackletHits* dev_three_hit_tracks,  
6     bool* dev_hit_used,  
7     uint* dev_atomics_VELO);  
8  
9 struct velo_three_hit_tracks_adder_t : public GpuAlgorithm {  
10     constexpr static auto name {"velo_three_hit_tracks_adder_t"};  
11     decltype(gpu.function(velo_three_hit_tracks_adder)) function {velo_three_hit_tracks_adder};  
12     using Arguments = std::tuple<  
13         dev_VELO_cluster_container,  
14         dev_estimated_input_size,  
15         dev_tracks,  
16         dev_three_hit_tracks,  
17         dev_hit_used,  
18         dev_atomics_VELO>;  
19  
20     void set_arguments_size(  
21         ArgumentRefManager<Arguments> arguments,  
22         const RuntimeOptions& runtime_options,  
23         const Constants& constants,  
24         const HostBuffers& host_buffers) const;  
25  
26     void operator()(  
27         const ArgumentRefManager<Arguments>& arguments,  
28         const RuntimeOptions& runtime_options,  
29         const Constants& constants,  
30         HostBuffers& host_buffers,  
31         cudaStream_t& cuda_stream,  
32         cudaEvent_t& cuda_generic_event) const;  
33 };
```

An algorithm definition (2)

```
1 void velo_three_hit_tracks_adder_t::operator()(
2     const ArgumentRefManager<Arguments>& arguments,
3     const RuntimeOptions& runtime_options,
4     const Constants& constants,
5     HostBuffers& host_buffers,
6     cudaStream_t& cuda_stream,
7     cudaEvent_t& cuda_generic_event) const
8 {
9     function.invoke(dim3(host_buffers.host_number_of_selected_events[0]), block_dimension(), cuda_stream)(
10         arguments.offset<dev_VELO_cluster_container >(),
11         arguments.offset<dev_estimated_input_size >(),
12         arguments.offset<dev_tracks >(),
13         arguments.offset<dev_three_hit_tracks >(),
14         arguments.offset<dev_hit_used >(),
15         arguments.offset<dev_atomics_VELO >());
16 }
17
18 __global__ void velo_three_hit_tracks_adder(
19     uint32_t* dev_VELO_cluster_container,
20     uint* dev_module_cluster_start,
21     Velo::TrackHits* dev_tracks,
22     Velo::TrackletHits* dev_three_hit_tracks,
23     bool* dev_hit_used,
24     uint* dev_atomics_VELO)
25 {
26     ...
27 }
```

All algorithms in Allen are by design **SIMD**. Conventional GPU algorithms require parallelism at two levels:

- Grid dimension - Independent groups of work
- Block dimension - Threads sharing a common cache; they can be synchronized

All algorithms in Allen are by design **SIMD**. Conventional GPU algorithms require parallelism at two levels:

- Grid dimension - Independent groups of work
- Block dimension - Threads sharing a common cache; they can be synchronized

Every event in LHCb is an independent physics event. Within each event, some algorithms exhibit higher parallelizability than others. Most times, algorithms benefit from the following convention:

- Grid dimension - *Events under execution*
- Block dimension - *Intra-event parallelism*

Cross-architecture compatible

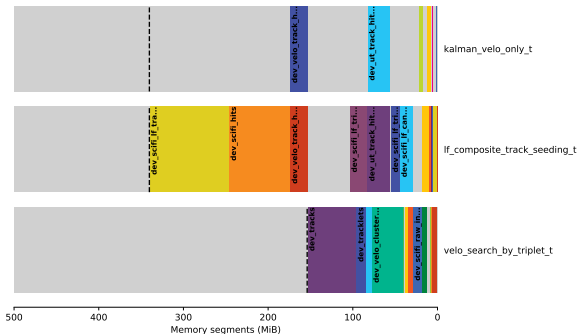
If the code has *block-dimension strided for loops*, and all `if` statements for a single thread refer to threads of index 0, then with some macros and function definitions it is possible to compile the code for CPUs:

```
1 // Definitions
2 #define __global__
3 #define __syncthreads()
4 struct GridDimensions { uint x, y, z; };
5 struct BlockIndices { uint x, y, z; };
6 struct BlockDimensions { uint x=1, y=1, z=1; };
7 struct ThreadIndices { uint x=0, y=0, z=0; };
8 thread_local GridDimensions gridDim;
9 thread_local BlockIndices blockIdx;
10 thread_local BlockDimensions blockDim;
11 thread_local ThreadIndices threadIdx;
12
13 ...
14
15 // Kernel call excerpt
16 gridDim = {num_blocks.x, num_blocks.y, num_blocks.z};
17 for (unsigned int i = 0; i < num_blocks.x; ++i) {
18     for (unsigned int j = 0; j < num_blocks.y; ++j) {
19         for (unsigned int k = 0; k < num_blocks.z; ++k) {
20             blockIdx = {i, j, k};
21             function(std::get<I>(invoke_arguments)...);
22         }
23     }
24 }
```

Memory management

We allocate memory at the **startup** of the application. A custom memory manager assigns memory segments on demand. This is essentially the same ALICE is doing. Good practices:

- No dynamic memory allocations (within *operator()*).
- Consolidate memory into compact (AO)SOAs.
- Prefer coalesced accesses where possible.



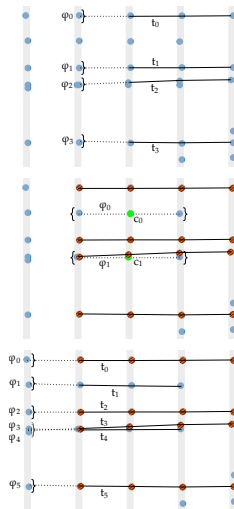
Algorithms

Example on spatial and temporal locality, data access patterns (VELO tracking)

- SIMD architectures benefit from coalesced and contiguous data access patterns
- Cache memory is limited in size
- Locality: Access patterns should restrict to a portion of memory

Search by triplet employs an SOA data structure for the VELO reconstruction, so that every access to memory has an increased probability of returning several required data in a cache line.

Additionally, modules in the VELO subdetector are visited only once, interleaving *seeding* and *forwarding* for all building tracks, maximizing spatial and temporal locality.

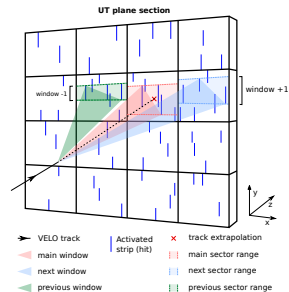


Example on spatial reductions (UT tracking)

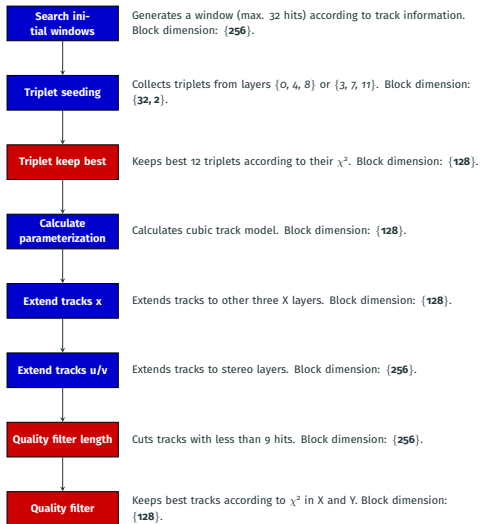
- Track reconstruction typically presents a high multiplicity of hit candidates
- Spatial reductions like KD-tree structures or search windows help reduce the dimensionality of hits under consideration

The UT subdetector is partially decoded into *sector groups*, aka groups of sectors sharing the same starting x coordinate. Within each sector group, hits are ordered by their y coordinate.

CompassUT determines search windows for each incoming Velo track. A configurable number of windows is determined, and binary searches are performed over x and y.

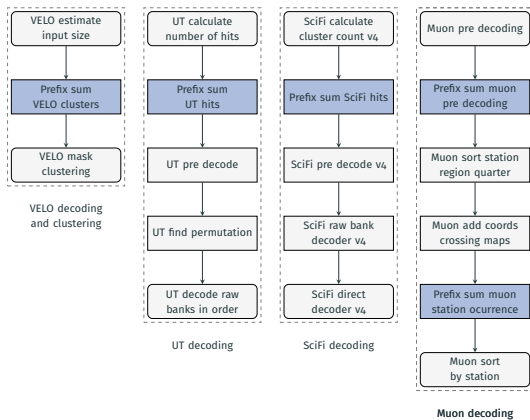


Each algorithm with its own grid and block size (Forward tracking)

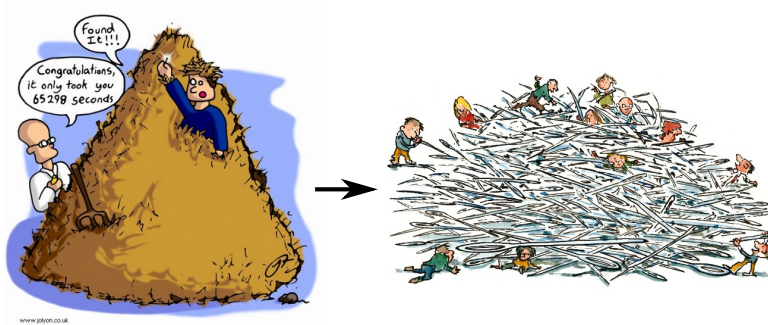


Calculate size, prefix sum (CPU), allocate

The sizes of all buffers are de/allocated with a custom memory manager. Prefix sums (*accumulated sums*) are used to determine offset / size of data buffers per event. We *offload* them to CPU.



Due to the increase in arithmetic capacity from GPUs, it is possible to tune for better physics while fulfilling the HLT1 throughput requirement of LHCb (more in a moment).



Integration

Baseline scenario without GPUs

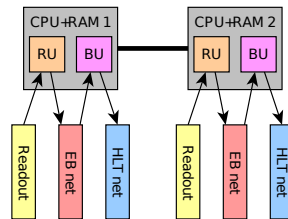
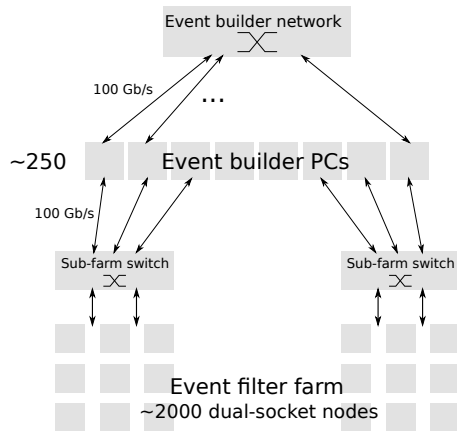


Figure 1: Event builder PC.

Event builders with GPUs

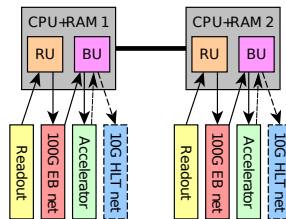
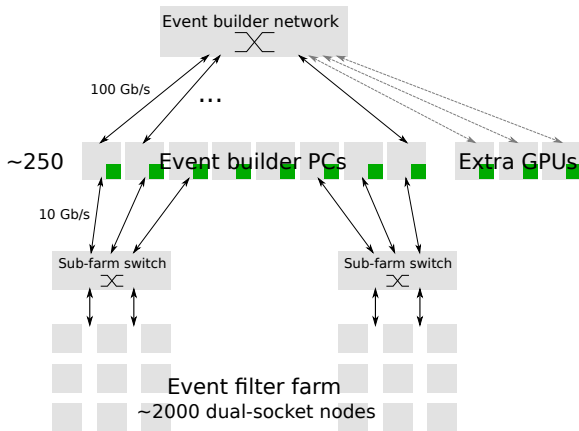
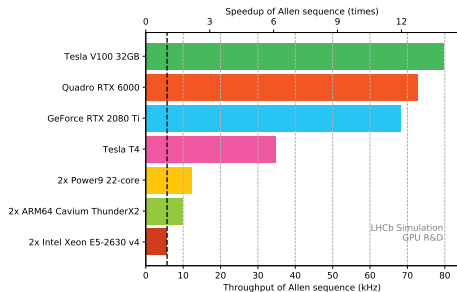


Figure 2: GPU-equipped event builder PC.

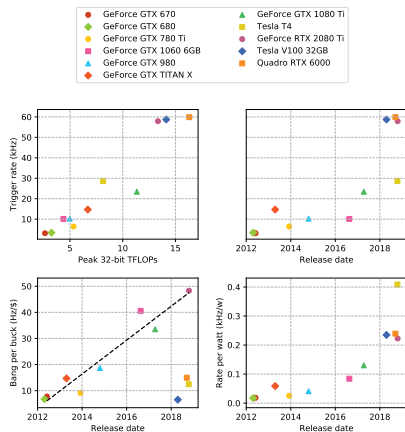
Target processing rate

In order to be able to perform the HLT1 filter inside the event builder with GPUs, the full throughput of collisions must be processed in near-time.



Opportunistic GPU usage: When there is no data-taking, GPUs could be use for something else.

Trigger rate versus peak TFLOPS, cost and power envelope across graphics cards:



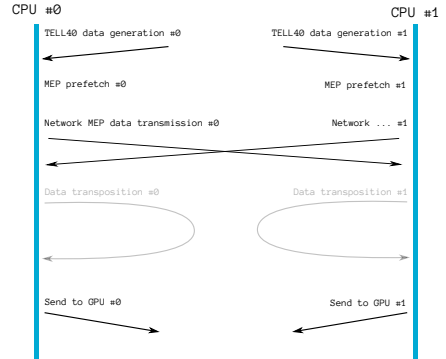
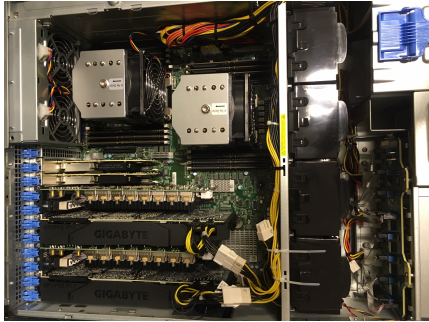
1

¹Optimization of high-throughput real-time processes in physics reconstruction, PhD thesis, Daniel Hugo Cámpora Pérez.

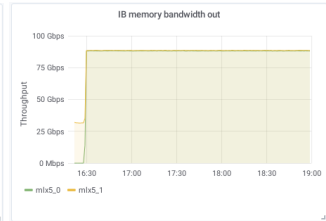
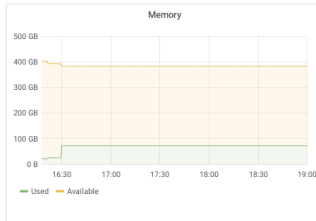
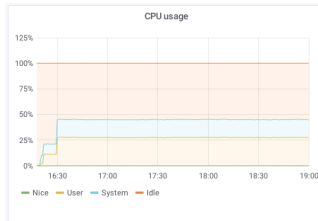
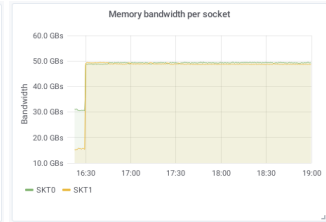
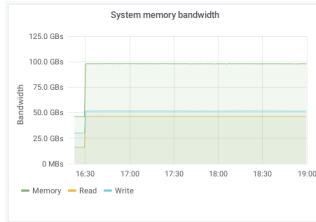
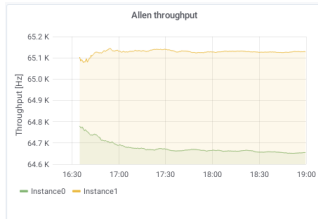
Integration

The most attractive realization of a GPU HLT1 in LHCb would be in the Event Builders.

Many aspects need to be demonstrated, such as CPU consumption, memory consumption and throughput, airflow, thermal stability, GPU performance stability, network throughput...



Selected integration test results



Conclusions

- ALICE and LHCb will have **full software** data processing in Run 3.
- Both experiments have developed GPU-based systems that are ready for Run 3.

- ALICE expects **100**× more data, and do an online compression of events.
- Most of the code exist for GPUs (baseline), and there is room for more.
- Data management incurs overhead, more consecutive components on the GPU would remove this overhead.

- LHCb expects about **40**× more data, and does a two-stage trigger.
- The entire sequence of the first stage *HLT1* has been developed for GPUs.
- The GPU HLT1 is being studied as an alternative solution to the baseline (CPUs).

- Both ALICE and LHCb use **custom memory managers** for the GPU to avoid the cost of allocating / deallocating.
- Both have **cross-architecture compatibility** by a switch at compile time.

In terms of integration,

- ALICE has more experience with a DAQ system with GPUs.
- Next steps are to have a better GPU usage in synchronous / asynchronous stages.

- LHCb needs proof that we are not breaking the experiment.
- Integration and burnout tests are ongoing.
- In terms of hardware, the system would occupy space in the Event Builders.
- In terms of software, HLT1 is its own application, and needs to be able to speak to certain parts of Gaudi and produce output understandable by HLT2 application.

Both have shown efficient implementations of GPU tracking algorithms and others:

- In ALICE, a *global method* is used to find track seeds, which maps naturally to parallel architectures.
- In LHCb, a *local method* exploits parallelism in modules (VELO), and track seeds are worked on in the following algorithms (PV, UT, Forward tracking, Kalman).
- LHCb requires processing many events in parallel to saturate the GPU, due to the tiny event size compared to ALICE (100 kB versus 50 MB).
- ALICE processes several time frames in parallel to saturate modern GPUs.

Good practices for GPU

In general, good practices of SIMD are useful in GPUs, and GPUs have some additional ones:

- Division of work into parallelizable independent tasks; blocks of threads can communicate (ie. wider vector widths).
- Data structures AOS vs SOA vs AOSOA, temporal and spatial locality.
- Dynamic memory allocation is costly, and more so on GPUs.

- GPUs are good for arithmetic workload with few branches.
- On the other hand, they are not just salvation: too high multiplicities will still run slow, and have to be worked out.
- If possible, it is better to run contiguous reconstruction steps on the GPU.

It has been amply shown that GPUs are efficient architectures for HEP workloads.

- They require hard work, and they provide potential better physics and throughput.
- Heterogeneous computing is becoming more of a reality.

In particular, the experience in ALICE and LHCb show that this effort scales to newer architectures.

- All the lessons learned in modern CPUs (SIMD) transfer to GPUs.
- Learning curve is not steep, many tools to get started.

Looking forward to seeing HEP GPU applications grow in the next few years.

Backup

Consider the following CUDA code:

```
1 constexpr int N = 32;
2 __global__ void saxpy_plus(float* x, float* y, const float a) {
3     y[threadIdx.x] = x[threadIdx.x] * a + y[threadIdx.x];
4     __syncthreads();
5     if (threadIdx.x < 10) {
6         y[i] += 1;
7     }
8     if (threadIdx.x == 11) {
9         y[threadIdx.x] += 20;
10    }
11 }
12 ...
13 saxpy_plus<<< /*blocks*/ M, /*threads*/ N>>>(x, y, a);
```

- The number of threads is set statically to N=32.
- The statement in line 3 makes assumptions of the number of threads.
- The two `if` statements also make assumptions of the number of threads (they require at least 11 threads).

In contrast, consider this code:

```
1 constexpr int N = 32;
2 __global__ void saxpy_plus(float* x, float* y, const float a) {
3     for (int i=threadIdx.x; i<N; i+=blockDim.x) {
4         y[i] = x[i] * a + y[i];
5     }
6     __syncthreads();
7     for (int i=threadIdx.x; i<10; i+=blockDim.x) {
8         y[i] += 1;
9     }
10    if (threadIdx.x == 0) {
11        y[11] += 20;
12    }
13 }
14 ...
15 saxpy_plus<<<< /*blocks*/ M, /*threads*/ T >>>(x, y, a);
```

- A call to `saxpy_plus` with any number of threads will produce the same result.

CPU support: A CPU version

If the code has *block-dimension strided for loops*, and all `if` statements for a single thread refer to threads of index 0, then with some macros and function definitions it is possible to compile the code for CPUs:

```
1 // Definitions
2 #define __global__
3 #define __syncthreads()
4 struct GridDimensions { uint x, y, z; };
5 struct BlockIndices { uint x, y, z; };
6 struct BlockDimensions { uint x=1, y=1, z=1; };
7 struct ThreadIndices { uint x=0, y=0, z=0; };
8 thread_local GridDimensions gridDim;
9 thread_local BlockIndices blockIdx;
10 thread_local BlockDimensions blockDim;
11 thread_local ThreadIndices threadIdx;
12
13 ...
14
15 // Kernel call excerpt
16 gridDim = {num_blocks.x, num_blocks.y, num_blocks.z};
17 for (unsigned int i = 0; i < num_blocks.x; ++i) {
18     for (unsigned int j = 0; j < num_blocks.y; ++j) {
19         for (unsigned int k = 0; k < num_blocks.z; ++k) {
20             blockIdx = {i, j, k};
21             function(std::get<I>(invoke_arguments)...);
22         }
23     }
24 }
```

Integration test setup (1)

For our first test, we setup a single server with:

- Supermicro server
- 2 × Intel Xeon Silver 4114
- 376 GB of memory
- Differences wrt. candidate server: Cascade Lake (better PCIe performance), different chassis (better thermals)

It has three PCIe Gen3 16x slots per socket. Two of those are double width. Configuration on each socket:

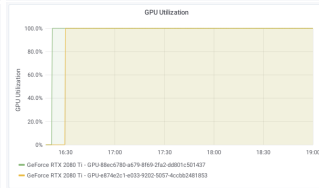
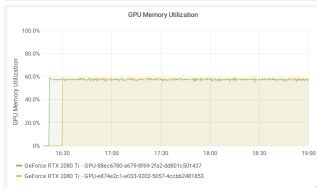
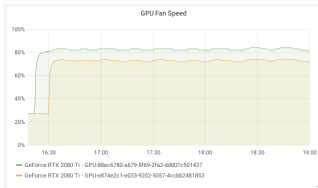
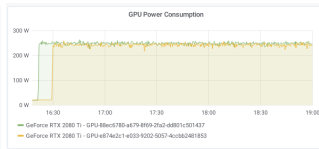
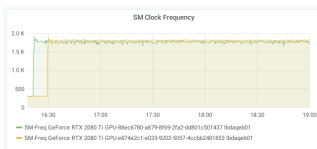
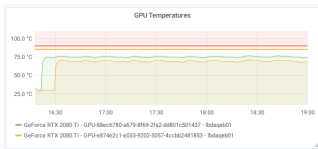
- Infiniband EDR card (100 Gbps)
- TELL40
- Gigabyte GeForce RTX 2080 Ti

Integration test setup (2)

Notes:

- The TELL40 can generate data into the server memory on each socket.
- Both network cards are connected back to back. A flow can then be simulated as if coming from the event building application.
- Each GPU can process data independently from each other. Two GPU applications are run, each one attached to a different GPU.

More integration test results (1)



More integration test results (2)

