



Commands which can be scheduled in parallel should/may run faster than the same commands serialized

A **command** is a request to execute work that is submitted to a queue such as the invocation of a SYCL kernel function, the invocation of a host task or an asynchronous copy. (SYCL 2020 Specification, Glossary)

• Overlaps is and optimization and is **not** required by the specification.

Which Commands can run Concurrently?

- 1. Multiple Low occupancy kernels
- In this paper we use a traditional "clpeak" kernel (chain of FMA on float (C)) scheduled to use one work-item 2. Host to Device data-transfer + Device to Host data-transfer
- In this paper data-transfer occur between host memory (either allocated via malloc (M), or pinned memorysycl::malloc_host (H)) and Device Memory (via sycl::malloc_device (D)). e.g.: H2D, D2M, ...
- 3. Low occupancy kernel + Host to/from Device data-transfer

This poster aims to report on with conditions different SYCL runtimes achieve concurrency. Deep explanations of the results are out of the scope of this poster. But we will be more than happy to discuss low-level details!

How to Achieve Concurrency?

The SYCL specification allow concurrent execution in two main scenarios:

- Out Of Order Queue Multiple SYCL commands submitted in a out of order queue can be executed concurrently ("OpenCL way")
- sycl::queue Q{sycl::gpu_selector()};
- ² for (auto& command: commands)
- do_work(Q,command);
- 4 Q.wait();
- Multiple In Order Queues. Multiple SYCL commands can be submitted in a multiple in-order queues. Each queue can be executed concurrently ("CUDA way")
- 1 const sycl::device D{sycl::gpu_selector()};
- 2 const sycl::context C(D);

```
3 std::vector<sycl::queue> Qs;
```

```
4 for (auto& : commands)
```

```
Qs.push_back(sycl::queue(C, D,sycl::property::queue::in_order{}));
```

6 for (int i = 0; i < commands.size(); i++)</pre>

```
do_work(Qs[i], commands[i]);
```

```
_{8} for (auto &Q : Qs)
```

```
9 Q.wait();
```

Compiler & Drivers Versions Used

ID Drivers Name

II CUDA

III ROCM

ID	Compiler Name	Commit
i	DPCPP	2022.1.0
 II	Intel/LLVM	03ff41f
 	Intel/LLVM	1fe5eaa
i∨	hipSYCL	258dc87

(a) Compiler Versions

(b) Drivers Verions

4.5.2

SYCL Concurrency on GPU Platforms: Empirical Measurement

Thomas Applencourt¹ Abhishek Bagusetty¹ Ajay Panyala² Aksel Alpay³

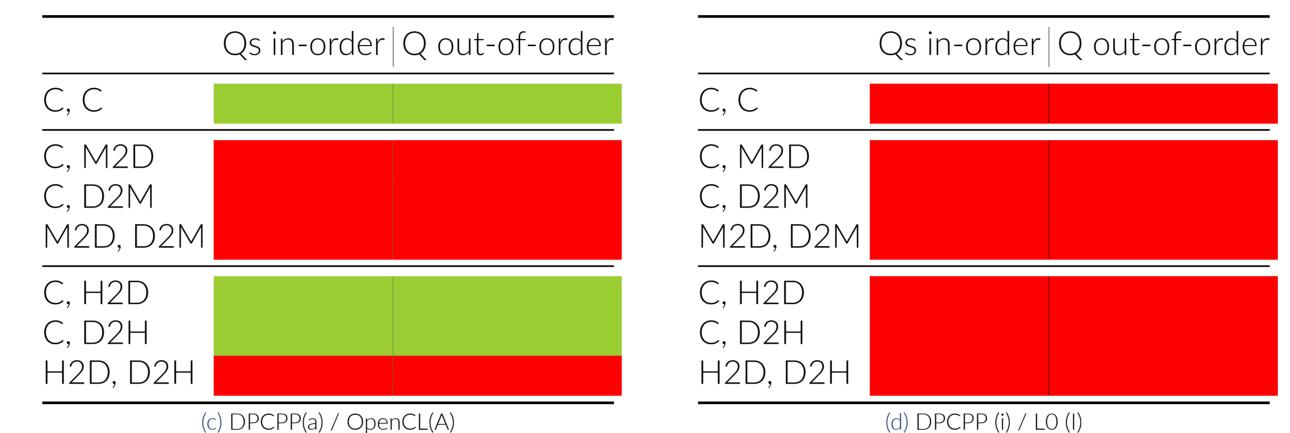
¹Argonne National Laboratory, ²Pacific Northwest National Laboratory, ³Heidelberg University



Concurrent Execution Concurrent Execution but Serial Execution when using Profiling enable Queues Serial Execution

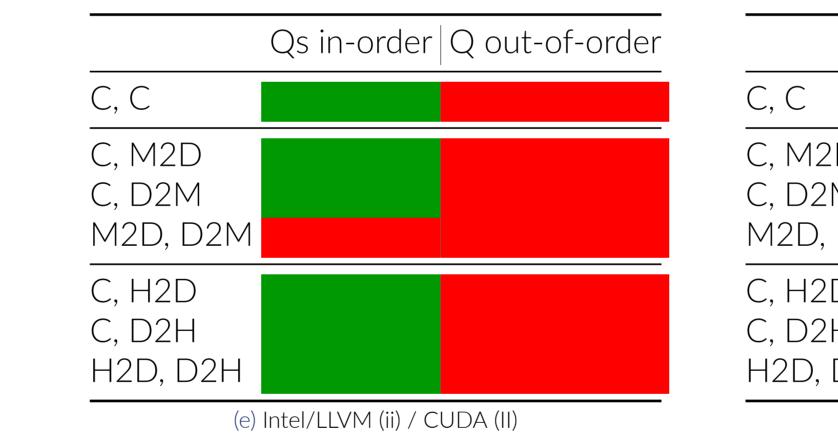
Table 1. Color Scheme Legend

Intel GPU Platform (Iris Pro Graphics 580)

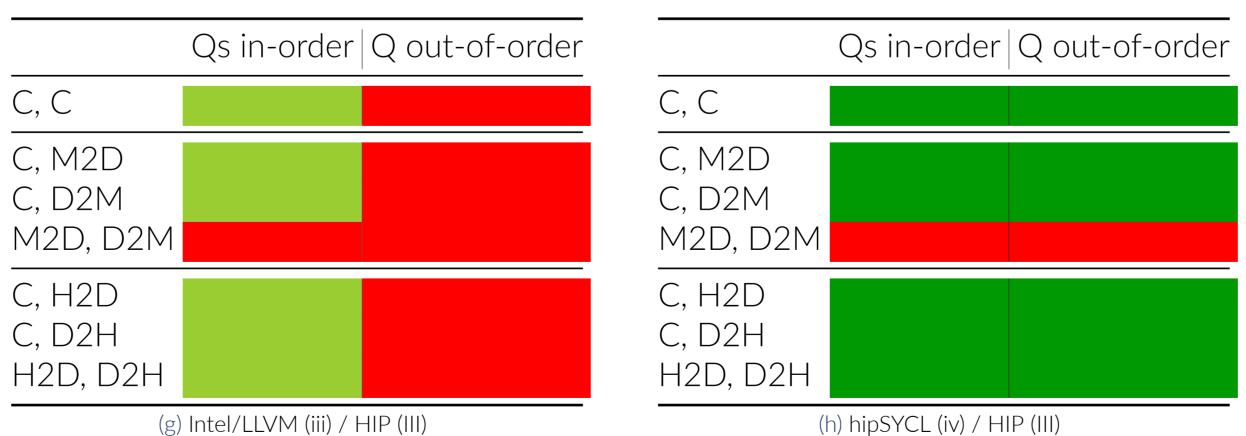


(c) DPCPP(a) / OpenCL(A)

NVDIA Platform (A100)



AMD Platform (MI100)



(g) Intel/LLVM (iii) / HIP (III)

Version Compute Runtime 21.40.21182 11.0.2

Result Color Code

	Qs in-order	Q out-of-order	
2D 2M , D2M			
2D 2H D2H			
(f) hipSYCL (iv) / CUDA (II)			

We run each step N times, and take the minimum time to reduce measurement noise

- max_command_time/total_time_serial
 - roughly the same time
- By default memory transfers are "as big as possible" (\approx D.get_info<sycl::info::device::max_mem_alloc_size>())
- total_time_concurrent/total_time_serial
- 3. Verify that Theoretical speedup and Empirical speedup roughly match

- 1 \$./sycl_con out_of_order C D2H ² Performing Autotuning ³ Parameters tunned: tripcount_C: 466765 globalsize_D2H: 1073739776 6 Best Total Time Serial: 610271us Best Time Command 0 (C): 303298us Best Time Command 1 (D2H): 306973us 9 Maximum Theoretical Speedup: 1.98803x 10 Best Total Time //: 355451us ¹¹ Speedup Relative to Serial: 1.71689x 12 SUCCESS: Close from Theoretical Speedup
- See url in the bottom right of the poster
- 3. On Intel Platform, LO need a little more love is many case.
- Intel/LLVM has a hard time with out-of-order queue **sycl::queues**. This Should help out-of-order queues performance.
- concurrency







Methodology

Run commands lists serially. Compute the maximum Theoretical Speedup

• Commands parameters (number of FMA, size of the data-transfer) are auto-tuned so that each command take

2. Run the list of commands in a mode who allow concurrency. Compute Empirical Speedup

Tool Output Example

Summary

1. We developed a empirical concurrency testing framework. People are encouraged to use it!

2. Using "pinned memory", sycl::malloc_host, may be required by drivers for concurrency

• Future version Intel/LLVM LO backend is planned to use immediate command list. Should allow concurrency

4. On NVDIA and AMD Platform, hipSYCL delivers most reliable overlap of operations.

• Discussion on-going for Intel/LLVM CUDA, HIP back-end to implement an M:N mapping between streams and

5. Using profiling queue, via **sycl::property::queue::enable_profiling**, can impact