

CUDA -> SYCL

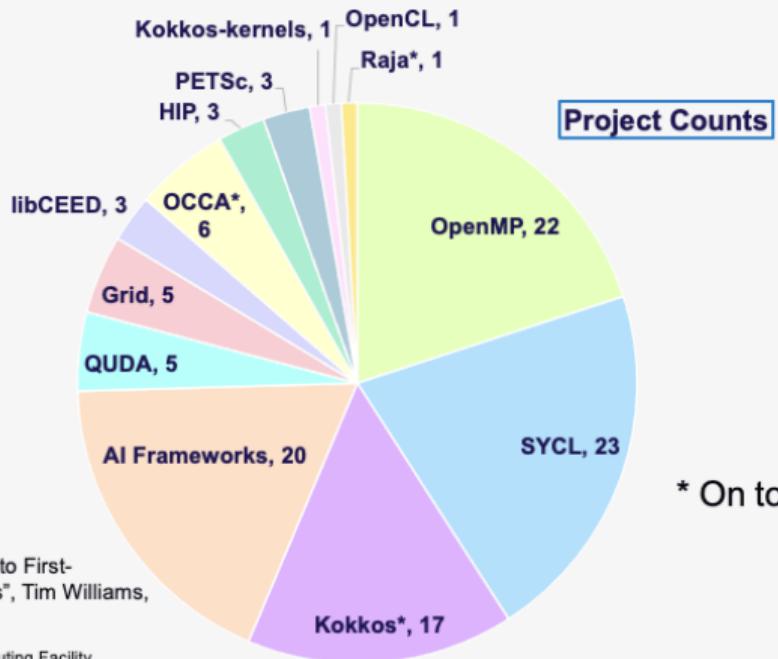
Thomas Applencourt, Abhishek Bagusetty
April 29, 2025

Overview of SYCL

SYCLomatic and cO

Overview of SYCL

Programming Model Choices by Y1 Aurora Projects + ECP



From: "Targeting Applications to First-Generation Exascale Systems", Tim Williams, Feb. 2025

Code Example

```
1 #include <sycl/sycl.hpp>
2
3 int main(int argc, char **argv) {
4     sycl::queue Q;
5     std::cout << "Running on "
6             << Q.get_device().get_info<sycl::info::device::name>() << "\n";
7
8     // Allocate Device Memory
9     int *A = sycl::malloc_device<int>(global_range, Q);
10    // Submit blocking kernel who use the memory
11    Q.parallel_for(global_range, [=](auto id) { A[id] = id; });
12    Q.wait();
13    // Allocate Host Memory
14    std::vector<int> A_host(global_range);
15    // Copy the device memory to the host
16    Q.copy(A, A_host.data(), global_range);
17    Q.wait();
18    sycl::free(A, Q);
19
20 }
```

In a nutshell

- At worse, just see SYCL as a "nice"¹ C++ wrapper around CUDA API ²
- Kernel code, will look similar³
- Run on Intel, NVIDIA, and AMD machine⁴

¹As much as C++ syntax can be nice []();

²At least no state machine, far cleaner API than CUDA Runtime

³Abhi and Steve are far more knowledgeable than me

⁴In the order of love received by the backend

- Stream == In Order Queue (default `sycl queue` out-of-order)
- Need to pass the ID to kernel
- Thread == Work-Item, Warp == Work-group
- Same synchronization primitive for kernel side (group barrier, shuffle, ...)
- Faster Indexing follow C++ (the reverse from CUDA – blame C++ / Nevin)
- Dependency via Event (DAG)
- Interop with CUDA (getting native object, or mix and match)

The SYCL Working Group listened

Official Khronos (KHR) Extension

- Free Function Query (no need to pass the ID)
- Queue Empty query
- Graph Extension to lower latency
- printf (help wanted)

SYCLomatic and cO
