October 10-12, 2023

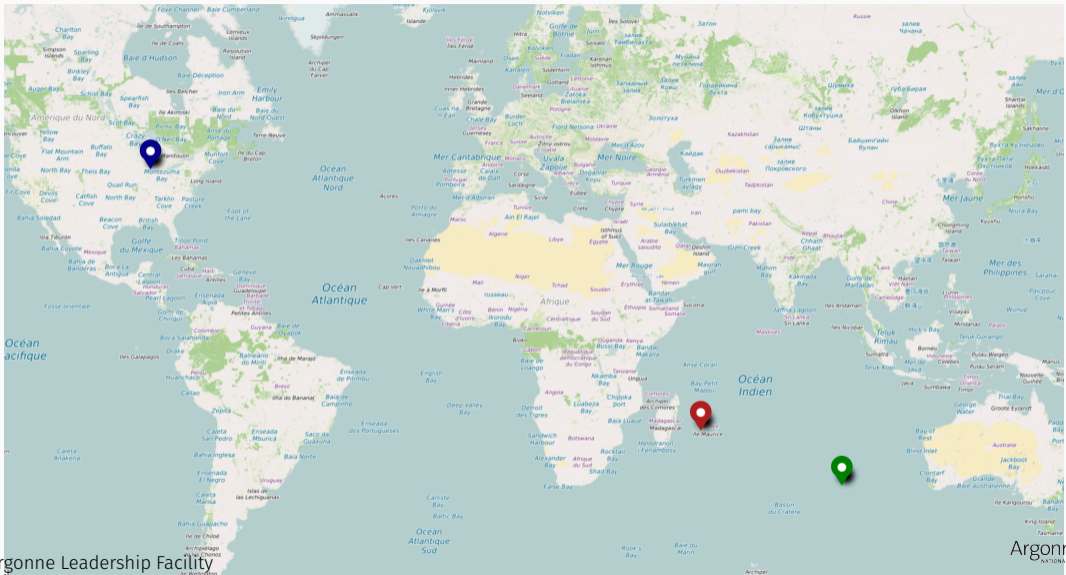# ALCF Hands-on HPC Workshop

# SYCL Overview

Thomas Applencourt

October 10, 2023

- Member of the Performance Engineering Team
- "Specialist" on Runtime[1]
- Argonne Representative to the SYCL committee

---
[1]At least I'm interested by them

Argonne
NATIONAL LABORATORY

- SYCL is a specification developed by the Khronos Group (OpenCL, SPIR, Vulkan, OpenGL)
- C++ 17 API
  - No language extension, No pragmas, No attribute
  - Lot of lambda, lot of template
- Borrow lot of *concept* from battle tested OpenCL (platform, device, work-group, range)
- Single Source
- Portable (HIP, Cuda Driver, L0 Backend exists)

Argonne
NATIONAL LABORATORY

```cpp
1   #include <sycl/sycl.hpp>
2   int main() {
3     // Create a "stream" aka object used to submit work to the device
4     sycl::queue Q;
5     // Some instrospection!
6     std::cout << "Running on "
7               << Q.get_device().get_info<sycl::info::device::name>()
8               << std::endl;
9     const int size = 10;
10    // Allocate memory who can be access from the host and the device
11    int *A = sycl::malloc_shared<int>(size,Q);
12    // Submit a kernel to the GPU, lambda! <3
13    Q.parallel_for(global_range, [=](sycl::id<1> idx) { A[idx] = idx; }).wait();
14    for (size_t i = 0; i < size; i++)
15      std::cout << "A[ " << i << " ] = " << A[i] << std::endl;
16    return 0;
17  }
```

Argonne ▲
NATIONAL LABORATORY

- Reduction, Atomic, Linear Algebra (via oneMKL)

```
1   sycl:: Q;
2   [...] // Allocated a left an an exercise to the reader
3   oneapi::mkl::blas::gemm(Q, oneapi::mkl::transpose::nontrans,
4                              oneapi::mkl::transpose::nontrans,
5                              m, n, k, alpha, A, ldA, B, ldB, beta, C, ldC).wait();
```
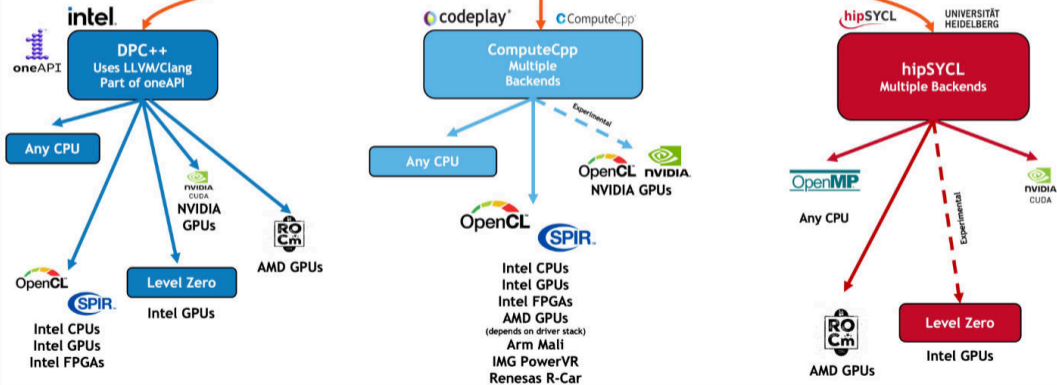
- Backend Interopt! (i.e can get the cuda::stream associated with a sycl::queue, or create a sycl::queue from a cuda::stream)

- Can mix OpenMP and SYCL

# Implementer of SYCL

- Easy peasy: *https://github.com/oneapi-src/SYCLomatic*
- Can translate full projects (Abhi is our home expert)

Argonne
NATIONAL LABORATORY

Controversial Take:

- SYCL and CUDA Runtime sit at the same level in the ladder abstraction: No fundamental reason performance should be different[2]
- Code gen is less important than a good runtime

---

[2]Minus event creation, batching...

Argonne
NATIONAL LABORATORY

# Data (we will add yours at the of the day)



BabelStream Performance



Lulesh Performance



DSlash Performance



RSBench Performance

# Pro and Cons of SYCL

Pro

- Close to C++ (familiarity to C++ programmer)
- Backed by Industry (a lot of man power)
- Portable (deployed at ALCF, NERSC and other)
- Simple
- Performant (please repport bugs if not)

Cons

- Close to C++ (lambda are intelligible and errors 1000 lines long)
- Backed by Industry (will they abandon it?)

Argonne
NATIONAL LABORATORY

- All of SYCL, aka enough to get your started!
- Know enough to know if SYCL if a good fit for your project
- *https: //docs.alcf.anl.gov/polaris/programming-models/sycl-polaris/*
- */eagle/projects/fallwkshp23/SYCL*

Argonne
NATIONAL LABORATORY