



GPU Performance Portability Using Standard C++ with SYCL

15th February 2023

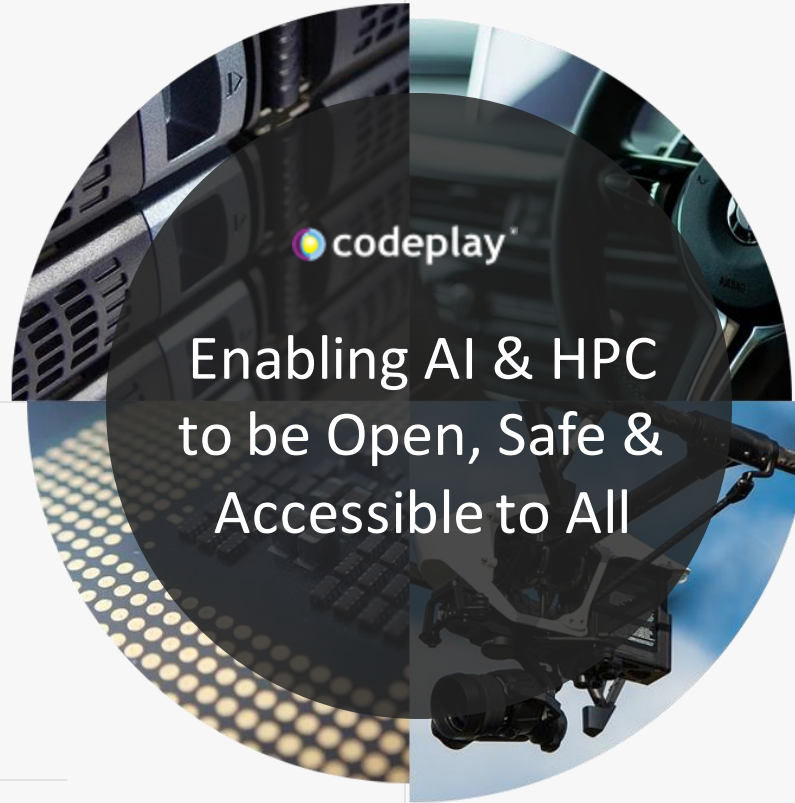
WAMTA 23 - Hugh Delaney

Company

Leaders in enabling high-performance software solutions for new AI processing systems

Enabling the toughest processors with tools and middleware based on open standards

Established 2002 in Scotland, acquired by Intel in 2022 and now ~90 employees.



codeplay
Enabling AI & HPC
to be Open, Safe &
Accessible to All

intel.

BROADCOM.

SYNOPSYS®
CEVA

Partners

Imagination

RENESAS

KMC
Kyoto Microcomputer Co., Ltd.

NSI-TEXE

BERKELEY LAB

OAK RIDGE
National Laboratory

Argonne
NATIONAL LABORATORY

And many more!

Supported Solutions



oneAPI

An open, cross-industry, SYCL based, unified, multiarchitecture, multi-vendor programming model that delivers a common developer experience across accelerator architectures

ComputeCpp™

C++ platform via the SYCL™ open standard, enabling vision & machine learning e.g. TensorFlow™

ComputeAorta™

The heart of Codeplay's compute technology enabling OpenCL™, SPIR-V™, HSA™ and Vulkan™

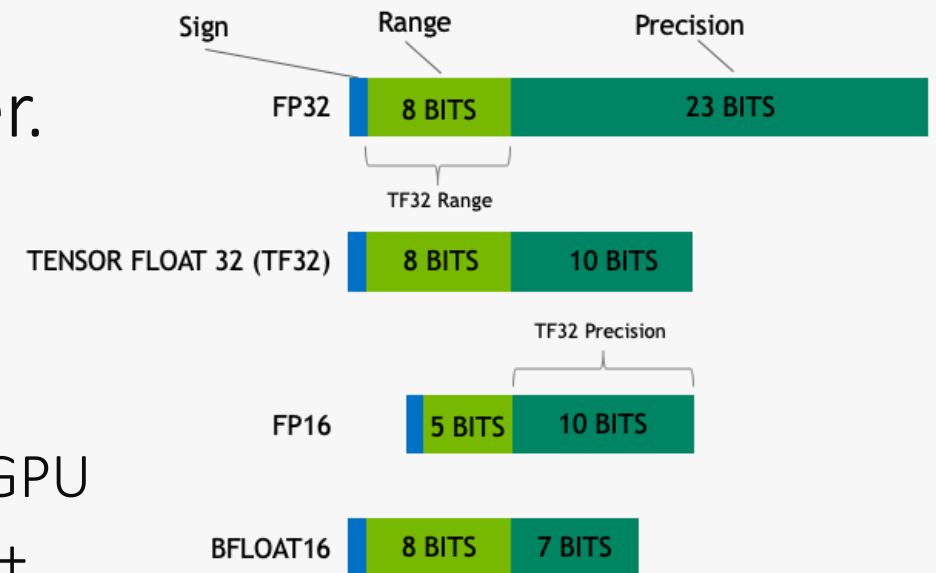
Markets

High Performance Compute (HPC)
Automotive ADAS, IoT, Cloud Compute
Smartphones & Tablets
Medical & Industrial

Technologies: Artificial Intelligence
Vision Processing
Machine Learning
Big Data Compute

Hugh Delaney – Software Engineer

- Working on the HIP and CUDA backend of the DPC++ SYCL compiler.
- Recent work:
 - Add support for the CXX standard library for the CUDA backend.
 - Add support for native atomic ops for AMDGPU
 - Add support for tf32 data type for the DPC++ joint matrix API using CUDA tensor cores.





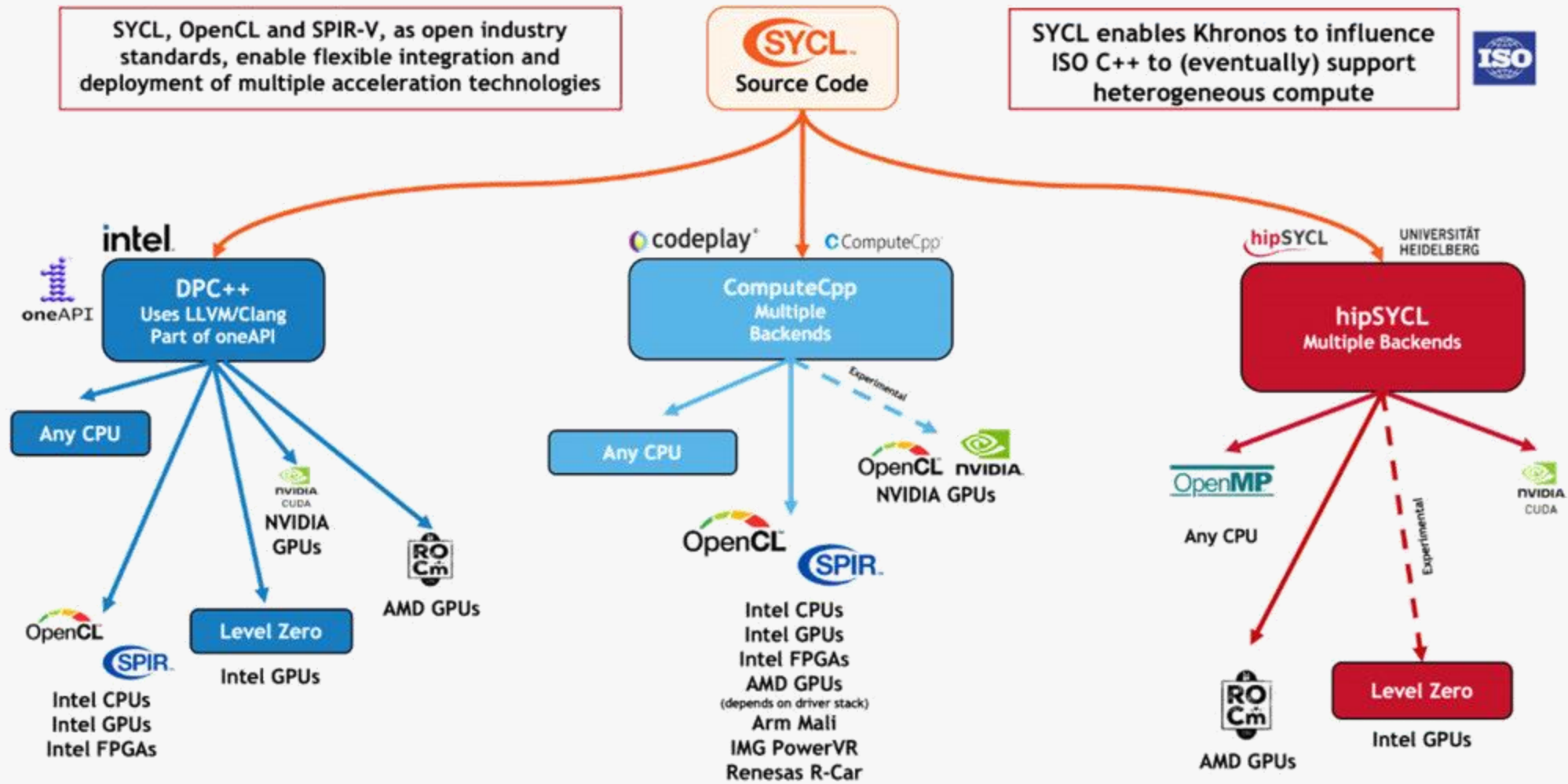
What is SYCL?

- SYCL is a C++ interface enabling heterogeneous acceleration.
- SYCL is an open standard under the Khronos Group.
- SYCL can be used to run single source code on AMD, Nvidia, or Intel GPUs, and more.
- Performance of SYCL code is comparable to performance of "native" APIs.
- SYCL is the "native" heterogeneous API for Intel GPUs.

```
array_view<float> a, b, c;
std::vector<float> a, b, c;
for(<2> idx) restrict(amp) {
#pragma parallel_for
for(int i = 0; i < a.size(); i++) {
    c[i]
}
__global__ vec_add(float *a, float *b, float *c) {
    return c[i] = a[i] + b[i];
}
float *a, *b, *c;
vec_add<<<range>>>(a, b, c);

cgh.parallel_for(range, [=](cl::sycl::id<2> idx) {
    c[idx] = a[idx] + b[idx];
});
```

SYCL Implementations



What Does SYCL Look Like?

```

1 #include <CL/sycl.hpp>
2
3 using T = float;
4
5 constexpr size_t n = 1638400;
6
7 int main() {
8     std::vector<T> A(n);
9     std::vector<T> B(n);
10
11     for (auto i = 0; i < n; ++i) {
12         A[i] = i;
13     }
14
15     sycl::queue q{};
16     std::cout << "Running on device: "
17               << q.get_device().get_info<sycl::info::device::name>() << std::endl;
18
19     try {
20         sycl::buffer bufA{A};
21         sycl::buffer bufB{B};
22
23         q.submit([&](sycl::handler &cgh) {
24             sycl::accessor accA{bufA, cgh};
25             sycl::accessor accB{bufB, cgh};
26             cgh.parallel_for(sycl::range{n},
27                             [=](sycl::id<1> i) { accB[i] = accA[i] * accA[i]; });
28         });
29     } catch (sycl::exception &e) {
30         std::cerr << "Caught exception: " << e.what();
31     }
32
33     std::cout << "Results: " << std::endl;
34     std::for_each(B.begin(), B.begin() + 10, [](T &c) { std::cout << c << " "; });
35     std::cout << "\n\n";
36 }
  
```



Device code

```

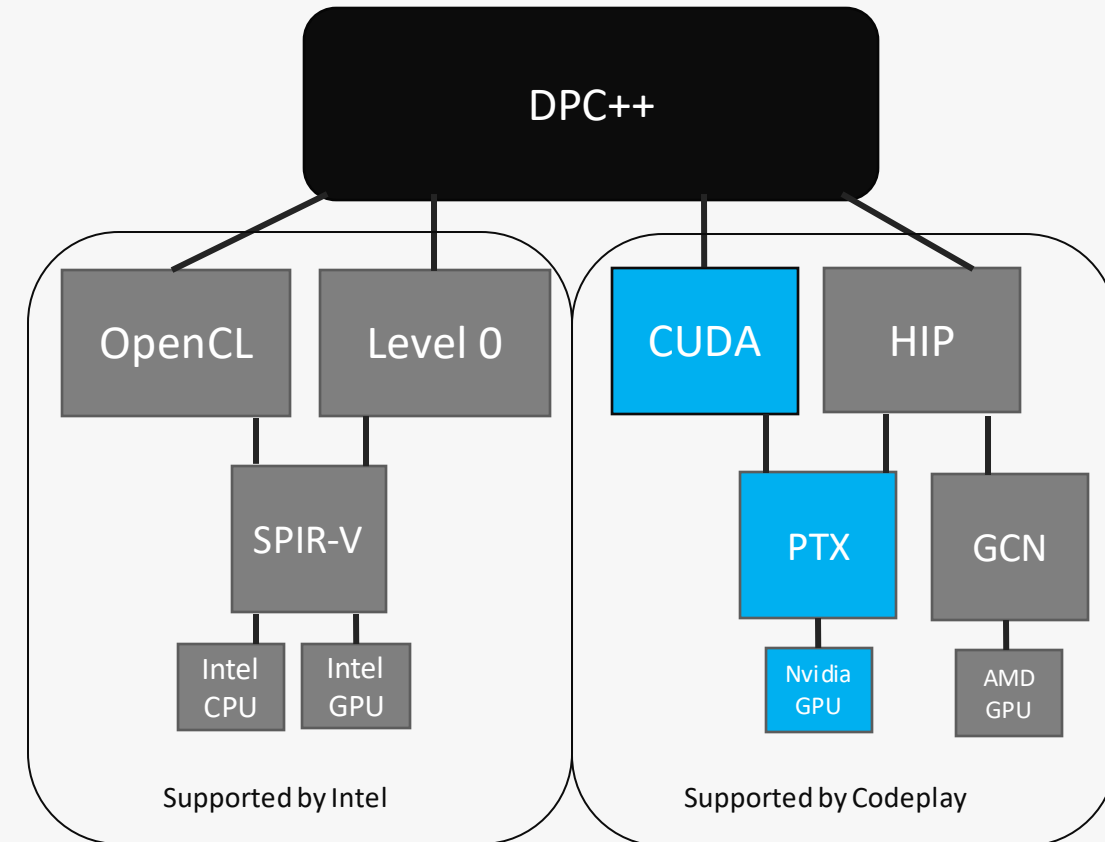
Executor x86-64 icx 2022.1.0 (C++, Editor #1)
A- [x] Wrap lines [x] [x] [x] [x]
x86-64 icx 2022.1.0 [x] -fsycl
Program returned: 0
Program stdout
Running on device: AMD EPYC 7R32
Results:
0 1 4 9 16 25 36 49 64 81
  
```

- Modern C++
- Uses exceptions, lambdas, templates, CTAD.
- Aside: Library-only implementations are possible.
- Experiment on <https://godbolt.org/z/6aon4Evsz>



What is DPC++?

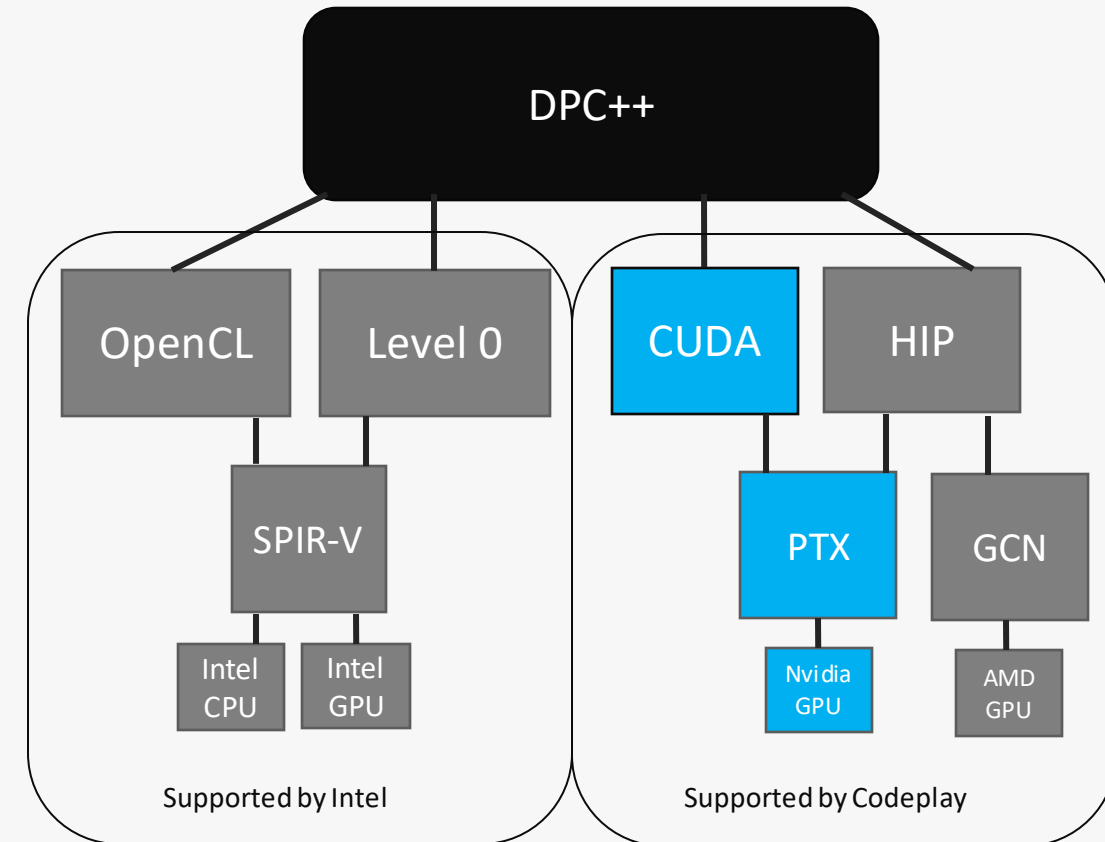
- DPC++ is the open source Intel-backed SYCL implementation.
- DPC++ is an LLVM fork.
 - Shared infrastructure with wider LLVM project
- CUDA-clang / openMPTarget / AMDGPU



How Does the DPC++ Toolchain Work?



- Kernel code is compiled for any platform specified to the compiler.
- Device code is bundled within the final binary.



Demo



- Compilation in action!

Portability



Examples of Specialization:

- Writing code for a specific architecture (i.e. GPUs)
- Writing code using vendor-specific APIs. (i.e. CUDA)
- Optimizing code for the specific features of a particular device.

Examples of Abstraction:

- Using APIs that map to different backends. (i.e. SYCL, openMP, Kokkos)
- Using libraries that support multiple backends. (i.e. oneDNN, oneMKL)
- Allowing an API to abstract away details of data transfer and heterogeneous execution.



Portability



Benefits of Specialization:

- Optimization

Benefits of Abstraction:

- Code portability.
- High level APIs (usually) increase productivity.
- No need to learn platform-specific APIs.

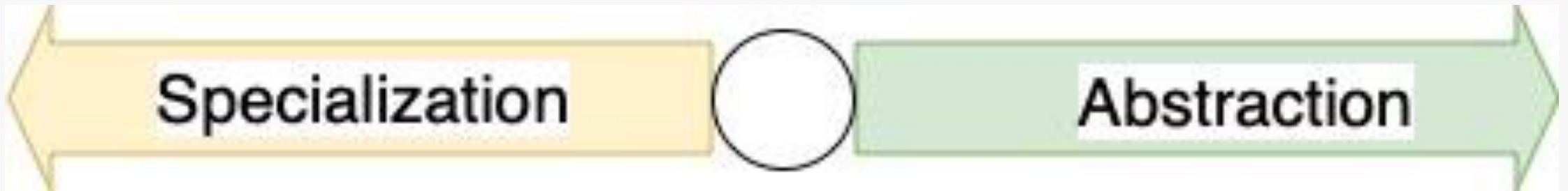


Portability



More programming effort
Lower productivity

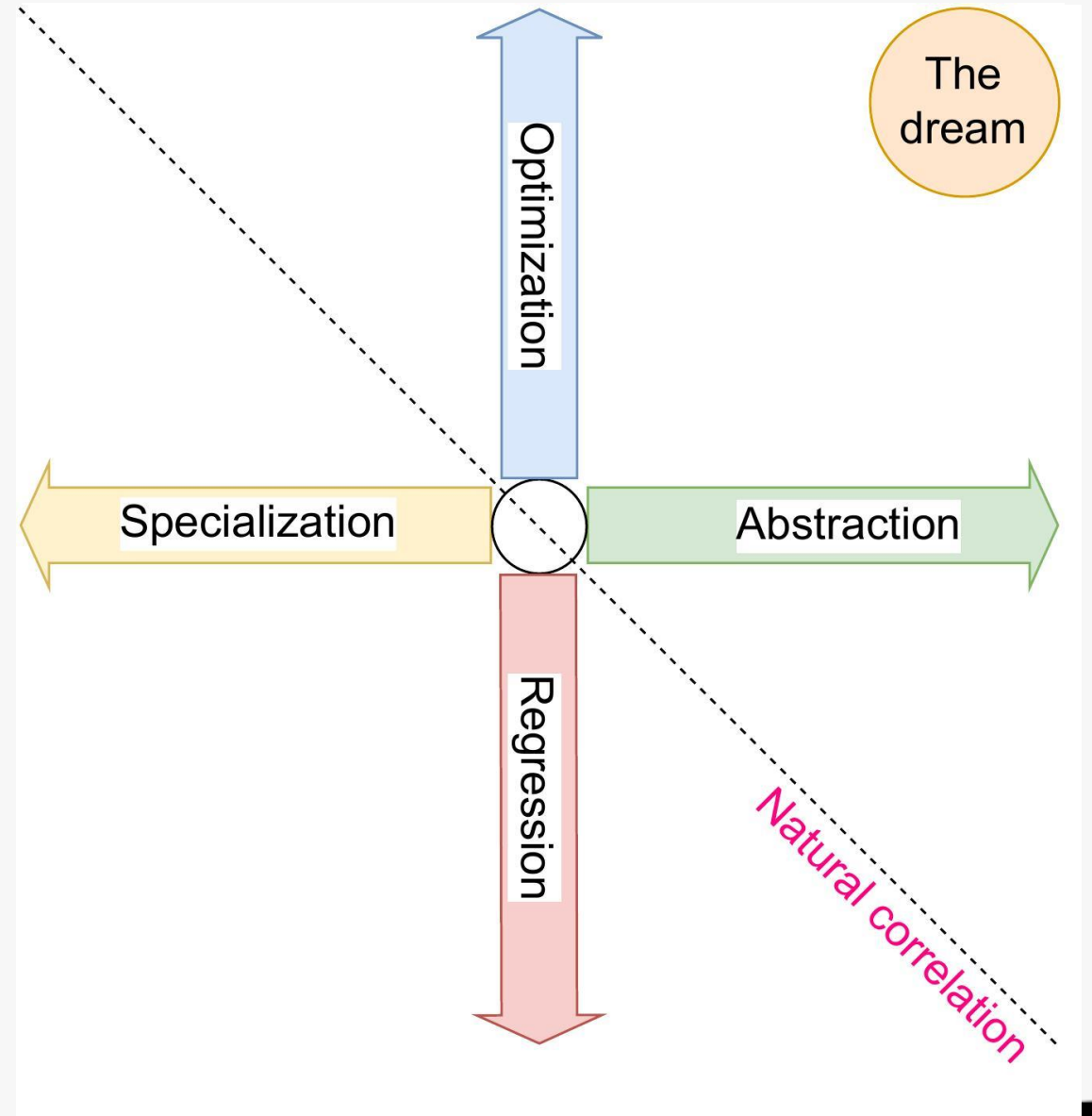
Less programming effort
Higher productivity



Performance

Portability

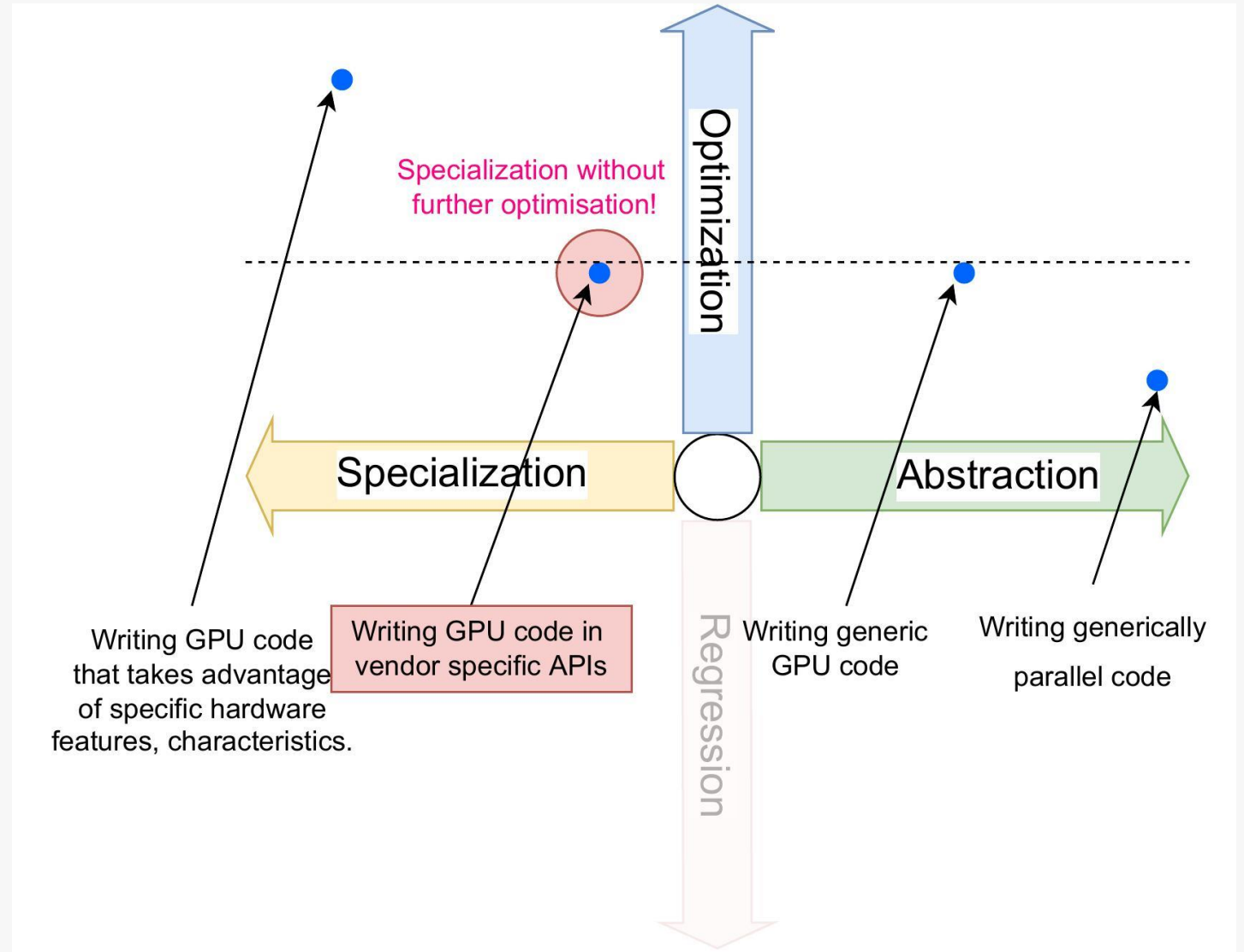
- Caveat: this is not a numerically precise scale. For a more rigorous presentation:
 - *"Navigating Performance, Portability, and Productivity", Pennycook, Deakin et al.*
- In general, the more specialized our code, the more optimization we should expect.
- We should not accept greater specialization, unless it offers greater performance.



Performance

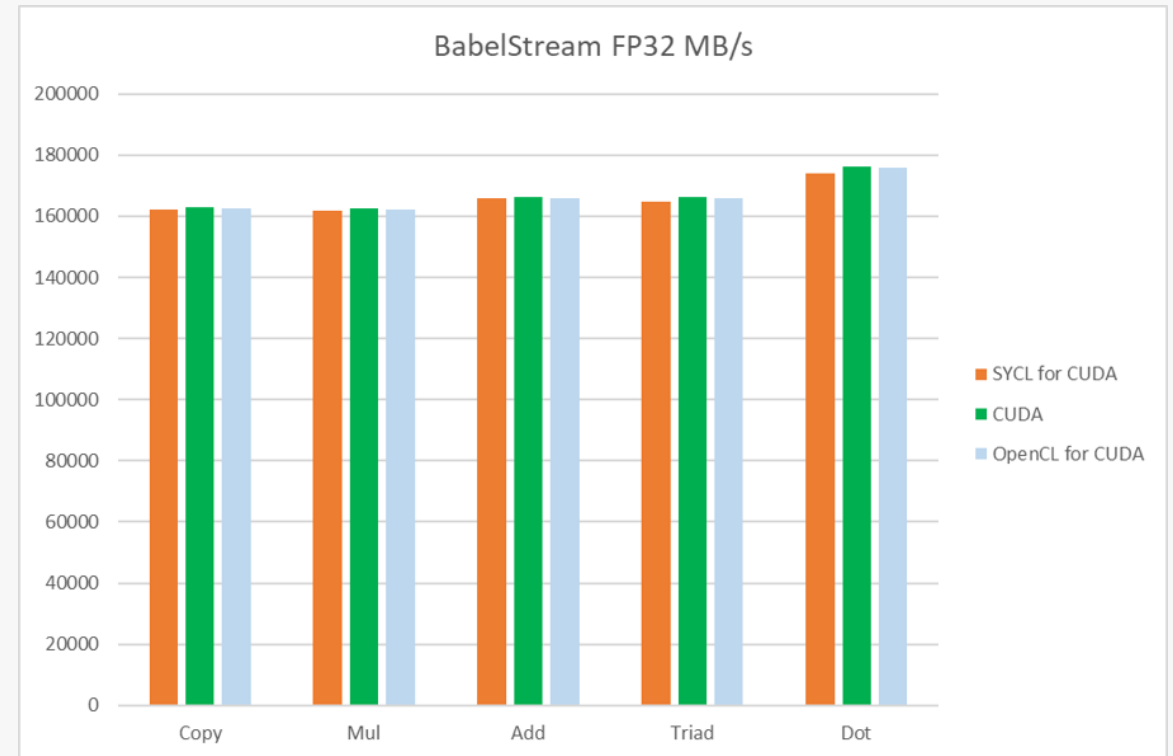
Portability

- Can GPU code written in SYCL really match the performance of native APIs?



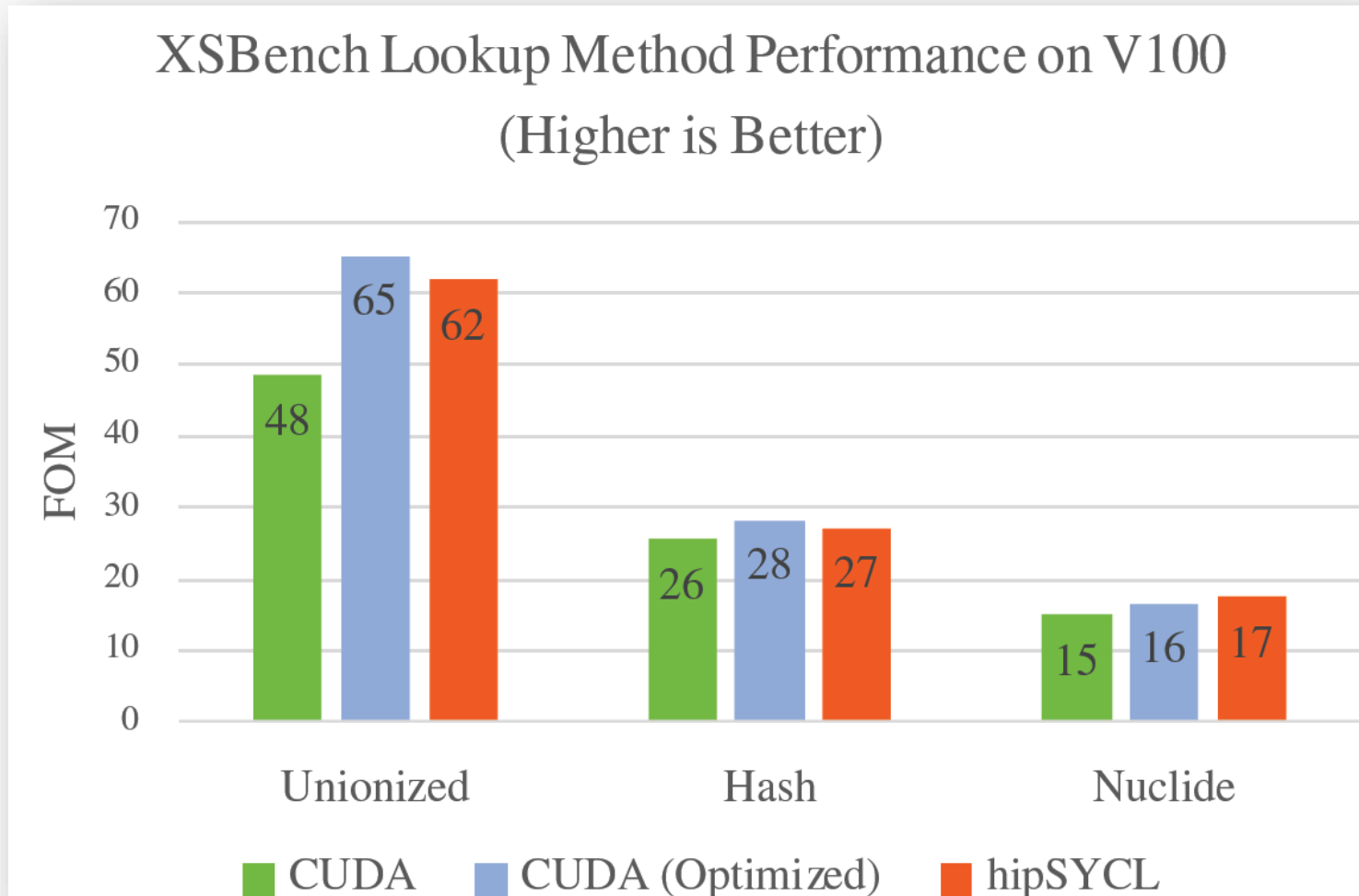
DPC++ performance on Nvidia GPUs

- This graph compares the BabelStream benchmarks results for:
 - Native CUDA code
 - OpenCL code
 - SYCL code using the CUDA backend
- Run on GeForce GTX 980 with CUDA 10.1
- Minimal SYCL overhead



<http://uob-hpc.github.io/BabelStream>

Argonne National Labs Comparison



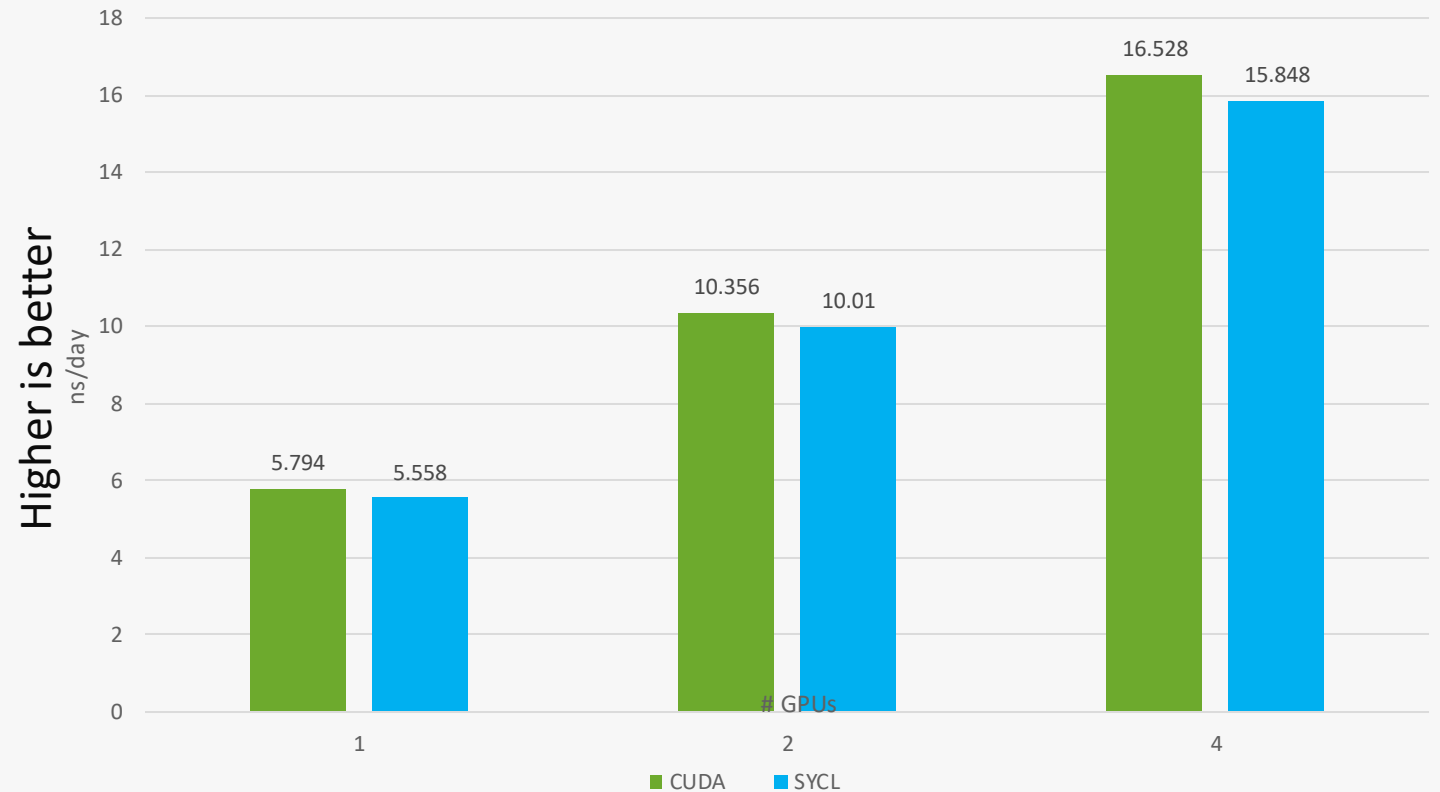
SYCL achieves equivalent performance to Optimized CUDA

[Presented at IWOC, April'20](#)

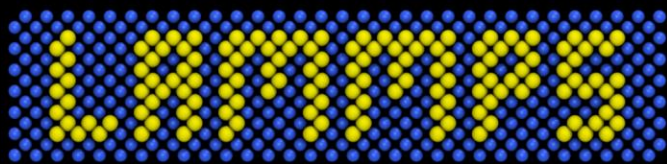
LAMMPS on DPC++

6.5M Atom Particle Simulation

- Molecular/Particle Simulator
- Distributed heterogenous computing over multiple GPUs.



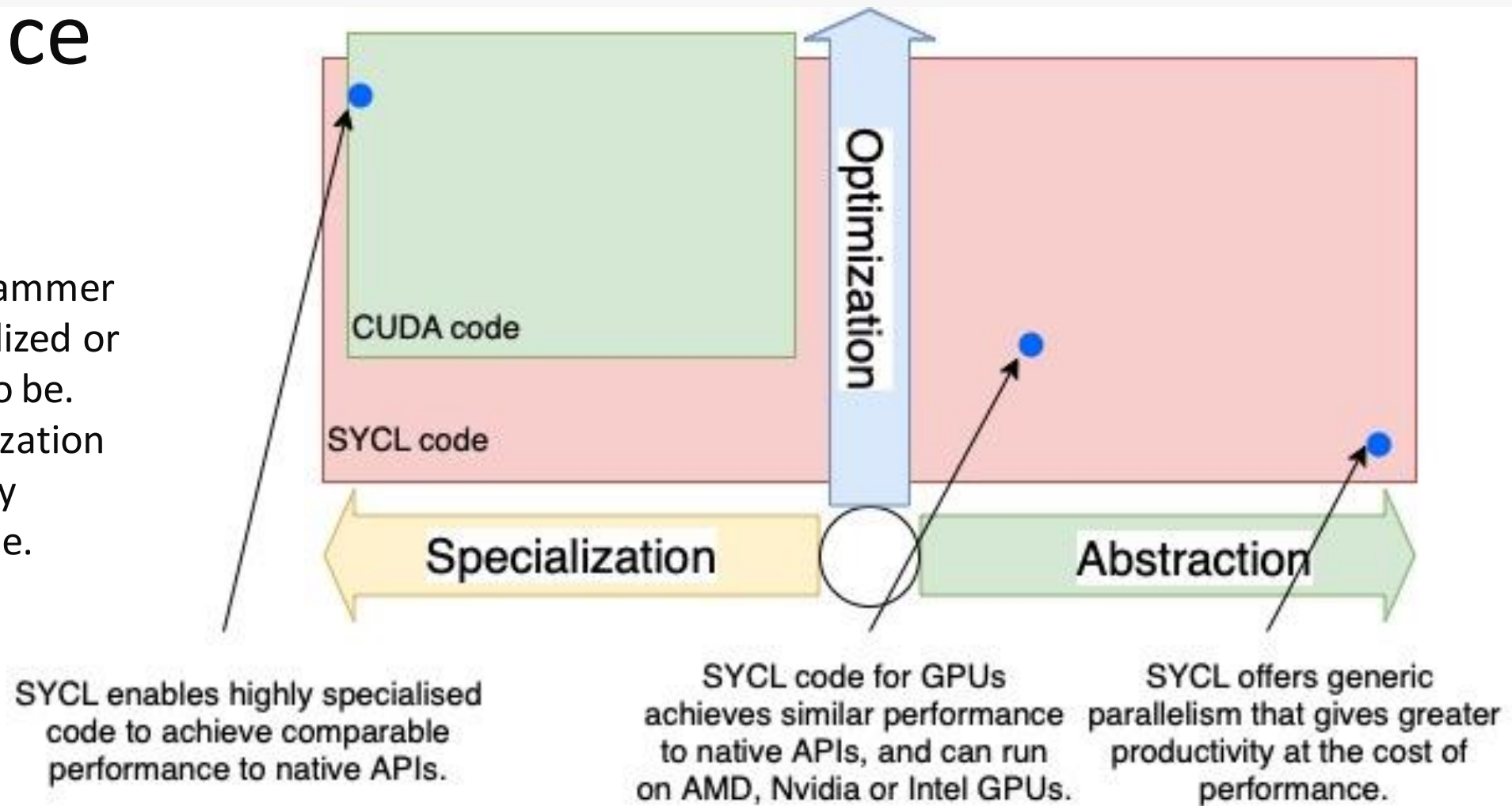
- Kokkos-CUDA vs Kokkos-SYCL for CUDA (same GPUs)



Large-scale Atomic/Molecular
Massively Parallel Simulator

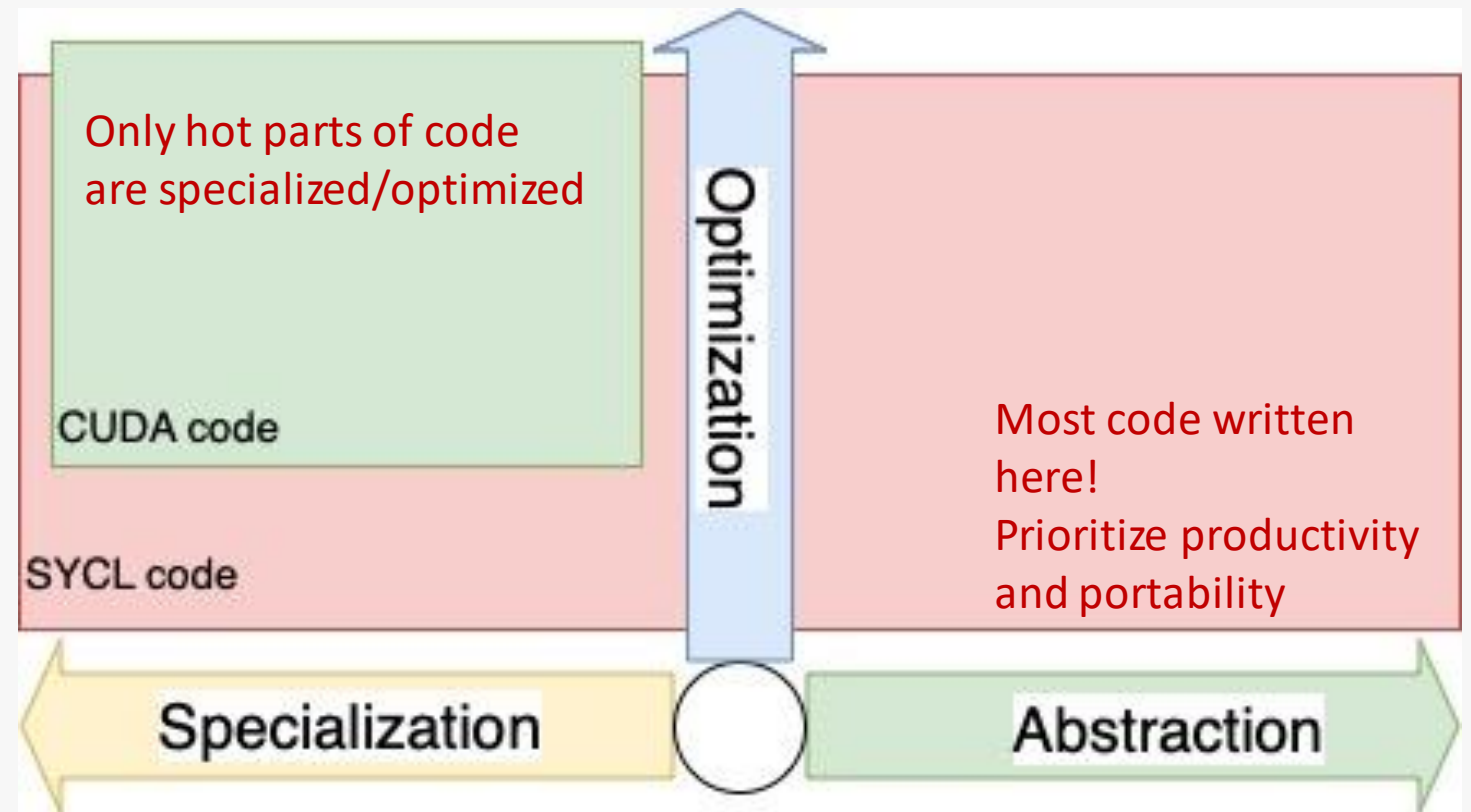
Performance Portability

- SYCL allows the programmer to choose how specialized or abstract code ought to be.
- The degree of specialization or abstraction can vary throughout a codebase.



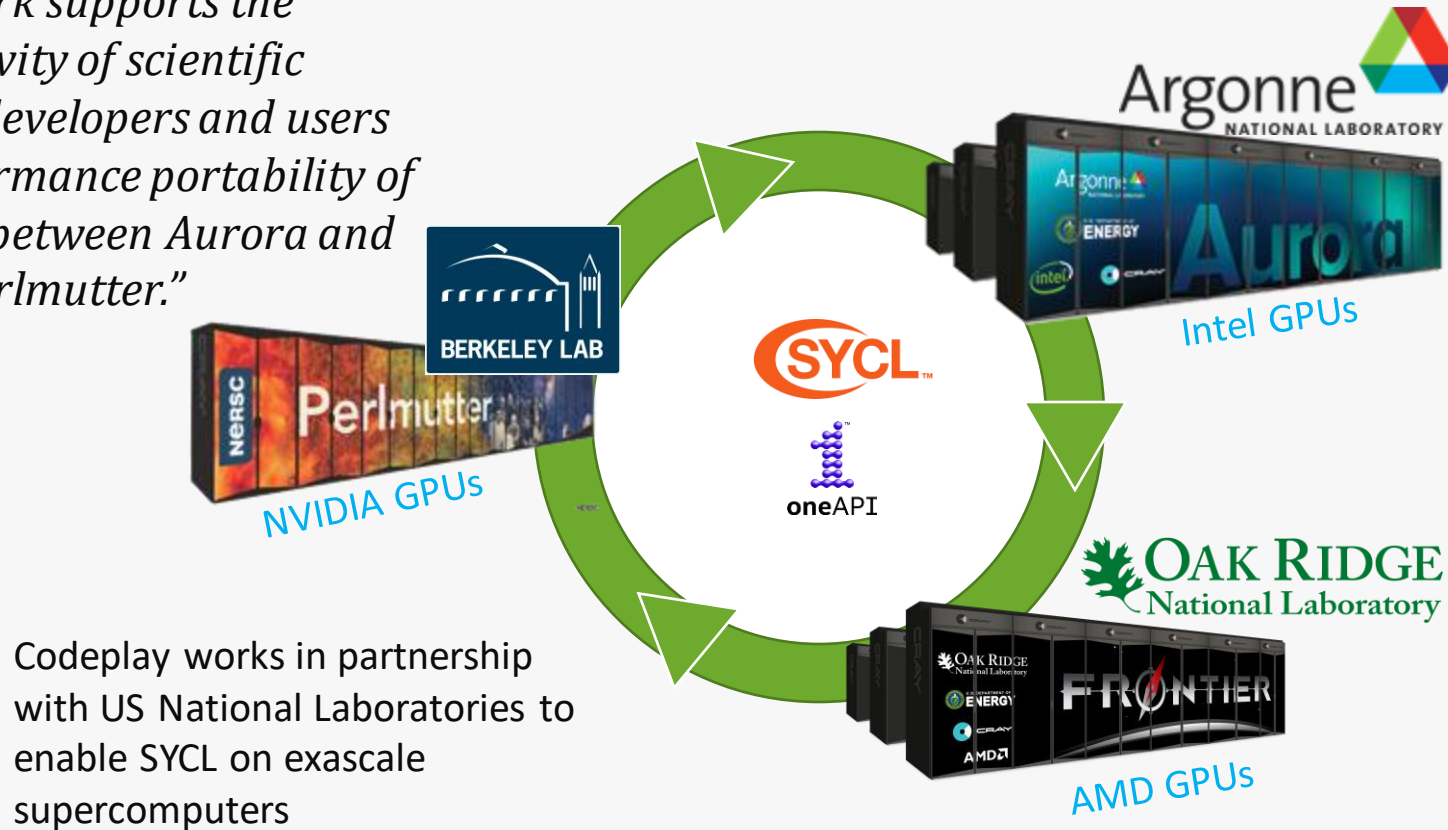
Performance Portability

- In a given codebase, only the "hot" parts need be highly specialized.
- Non "hot" parts are written in generically parallel SYCL.
 - Improves code portability.
 - Improves programmer productivity.



SYCL Enables Supercomputers

“this work supports the productivity of scientific application developers and users through performance portability of applications between Aurora and Perlmutter.”



Codeplay works in partnership with US National Laboratories to enable SYCL on exascale supercomputers

Enables a broad range of software frameworks and applications



Takeaway



With SYCL you can:

- Use modern C++.
- Write code that can target multiple architectures/platforms.
- Attain comparable performance to native APIs.
- Choose how generic or specialized your code ought to be.

Questions



COMPILER EXPLORER

Support Compiler Explorer on [Patreon!](#)

Sponsors: Backtrace, intel, Solid Sands

```
1 #include <CL/sycl.hpp>
2
3 using T = float;
4 constexpr size_t n = 1638400;
5
6 int main() {
7     std::vector<T> A(n);
8     std::vector<T> B(n);
9
10    for (auto i = 0; i < n; ++i) {
11        A[i] = i;
12    }
13
14    sycl::queue q{};
15    std::cout << "Running on device: "
16              << q.get_device().get_info<sycl::info::device::name>() << std::endl;
17
18    try {
19        sycl::buffer bufA{A};
20        sycl::buffer bufB{B};
21
22        q.submit([&](sycl::handler &cgh) {
23            sycl::accessor accA{bufA, cgh};
24            sycl::accessor accB{bufB, cgh};
25            cgh.parallel_for(sycl::range{n},
26                            [=](sycl::id<1> i) { accB[i] = accA[i] * accA[i]; });
27        });
28    } catch (sycl::exception &e) {
29        std::cerr << "Caught exception: " << e.what();
30    }
31
32    std::cout << "Results: " << std::endl;
33    std::for_each(B.begin(), B.begin() + 10, [](T &c) { std::cout << c << " "; });
34    std::cout << "\n\n";
35 }
36
```

Device code

Executor x86-64 icx 2022.1.0 (C++, Editor #1)

x86-64 icx 2022.1.0 -fsycl

<https://godbolt.org/z/E33oTG8vE>

