codeplay®

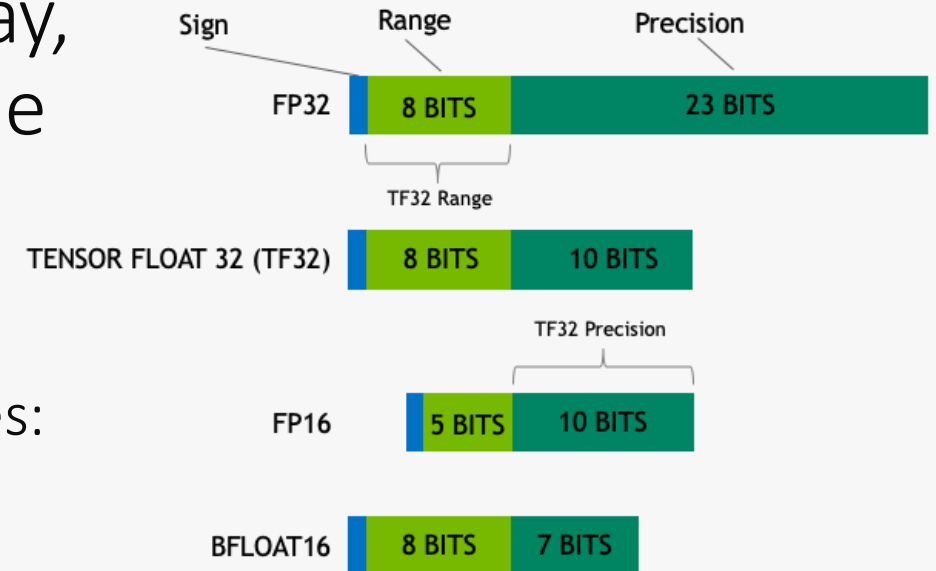Enable AI & HPC to be Open, Safe and Accessible to All

# Intro to SYCL and DPC++

Hugh Delaney – Software Engineer

Intel Bayncore Training – 28th April

- Nvidia is a registered trademark of NVIDIA Corporation

- AMD is a registered trademark of Advanced Micro Devices Corporation

- Intel is a registered trademark of Intel Corporation

- SYCL, SPIR are trademarks of the Khronos Group Inc. OpenCL and the OpenCL logo are trademarks of Apple Inc. used by permission by Khronos

# Introduction

- I am a software engineer at Codeplay, working on the CUDA backend of the Intel SYCL compiler (DPC++).

- Recent work:
  - Extensions to add support for new datatypes: tensorfloat32, bfloat16 to be used in tensor core mixed precision calculations.
  - Designing the SYCL joint_matrix API to allow for different precisions

## 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 with ~80 employees

## Products

### 🐦 Acoran

Integrates all the industry standard technologies needed to support a very wide range of AI and HPC

### ▲ ComputeAorta™

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

### C ComputeCpp™

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

## Partners

intel.
BROADCOM.
SYNOPSYS®
CEVA
Imagination

RENESAS
KMC Kyoto Microcomputer Co., Ltd.
NSI-TEXE

BERKELEY LAB
OAK RIDGE National Laboratory
Argonne NATIONAL LABORATORY

**And many more!**

## ● codeplay®

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

## Markets

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

**Technologies:** Artificial Intelligence
Vision Processing
Machine Learning
Big Data Compute

# Learning Objectives

**1**

Learn why you should use SYCL.

**2**

Understand why industry wants SYCL.

**3**

Learn some neat features of DPC++.

# SYCL First Example

```
1  #include <CL/sycl.hpp>
2  #include <vector>
3
4  constexpr size_t N = 10;
5
6  #define PRINT_ARRAY(a)                                                \
7    for (auto &e : a)                                                   \
8      std::cout << e << " ";                                            \
9    std::cout << std::endl;
10
11 using T = float;
12
13 int main() {
14   auto A = std::vector<T>(N);
15   for (auto i = 0; i < N; i++) {
16     A[i] = (T)i;
17   }
18
19   std::cout << "Before:\n";
20   PRINT_ARRAY(A);
21   auto q = sycl::queue{};
22
23   {
24     auto buf = sycl::buffer(A);
25     q.submit([&](sycl::handler &cgh) {
26       auto acc = buf.get_access(cgh);
27       cgh.parallel_for(sycl::range<1>(N),
28                        [=](sycl::id<1> idx) { acc[idx] *= 2; });
29     });
30   } // memory is updated once buffers are destructed
31
32   std::cout << "After:\n";
33   PRINT_ARRAY(A);
34 }
```

Kernel Code

# What is SYCL?

- High level, standard C++ programming model enabling hardware acceleration.

- Accelerators:
  - GPUs (Nvidia, AMD, Intel)
  - SIMD processors
  - Some FPGAs
  - Other niche hardware

# Why use SYCL?

- Code Portability
  - Target GPUs from Nvidia, Intel, AMD
  - Run on CPU, FPGAs, exotic hardware

- Performance Portability

- "But I already know how to use CUDA!"

# Why use SYCL?

- Code Portability
  - Target GPUs from Nvidia, Intel, AMD
  - Run on CPU, FPGAs, exotic hardware

- Performance Portability

- ~~"But I already know how to use CUDA!"~~

- You basically already know how to use SYCL!

codeplay®

# Motivation

- Hardware becoming more and more diverse, specialized to specific tasks.
    - GPUs specialized to graphics/deep learning.
    - FPGAs specialized to niche tasks.
- The programmer of the future will need to interface with lots and lots of different types of hardware.

# Motivation

- However:
  - Interfacing with specialized hardware is hard.
  - Not all hardware vendors can make software like CUDA.

codeplay ®

# oneAPI

- SYCL enables interface with GPUs, FPGAs, and more through the same programming model.

- Programmer can avoid hardware-specific languages, APIs.
  - Increased productivity.
  - High level SYCL API using existing C++ syntax and language features.

codeplay ®

# SYCL Increases productivity

```
1  // CUDA version
2  __device__ int getGlobalIdx_3D_3D() {
3    int blockId =
4        blockIdx.x + blockIdx.y * gridDim.x + gridDim.x * gridDim.y * blockIdx.z;
5    int threadId = blockId * (blockDim.x * blockDim.y * blockDim.z) +
6                  (threadIdx.z * (blockDim.x * blockDim.y)) +
7                  (threadIdx.y * blockDim.x) + threadIdx.x;
8    return threadId;
9  }
```

```
1  // SYCL version
2  template <size_t N>
3  inline size_t getGlobalIdx(sycl::item<N> my_item) {
4    return my_item.get_global_linear_id();
5  }
6
```

- SYCL uses a high-level modern C++ API.
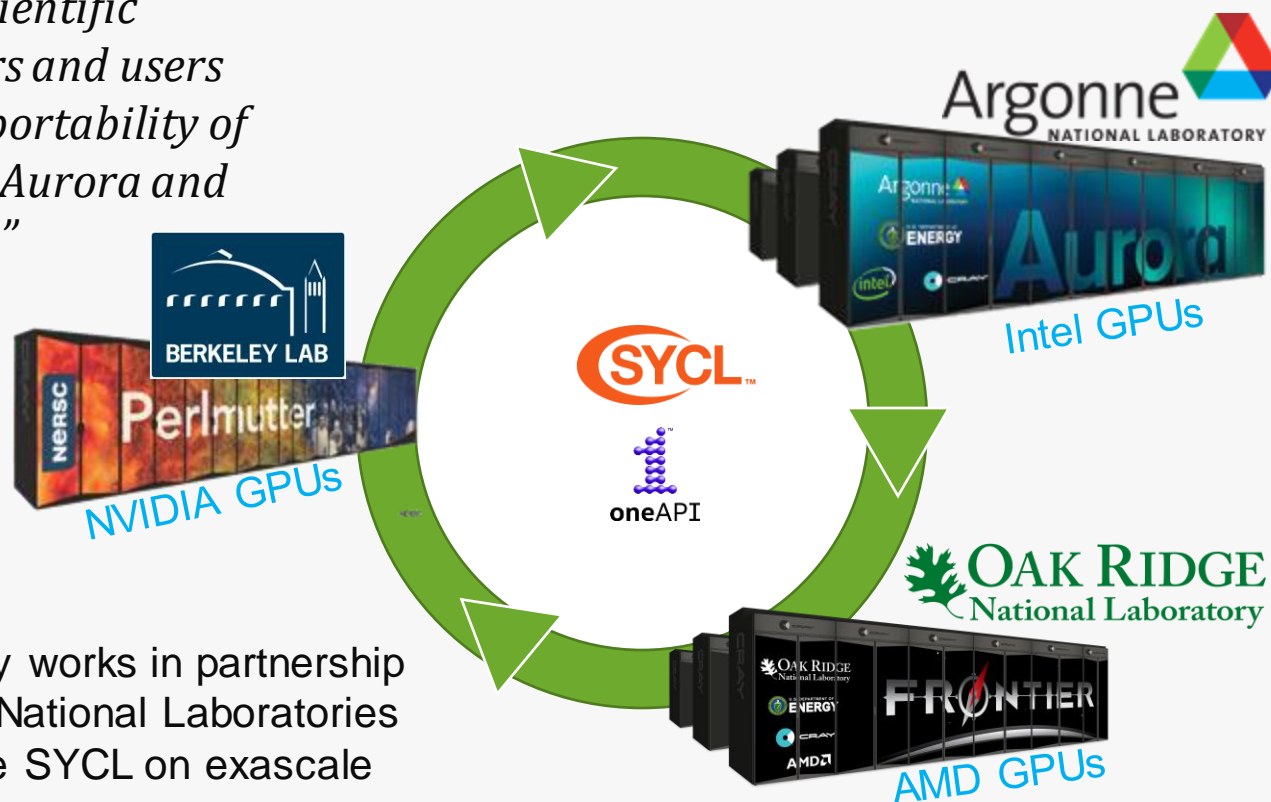  - Likes templates, lambdas and other modern C++ features

codeplay®

# oneAPI

- Code written in SYCL can easily migrate between hardware from different vendors.

- Hardware purchasers not locked-in to long term hardware/software commitments.

  - Especially relevant in age of supply-chain shocks.

# SYCL and DPC++ 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

Intel GPUs

NVIDIA GPUs

AMD GPUs

Enables a broad range of software frameworks and applications

kokkos RAJA
Celerity
alpaka BabelStream
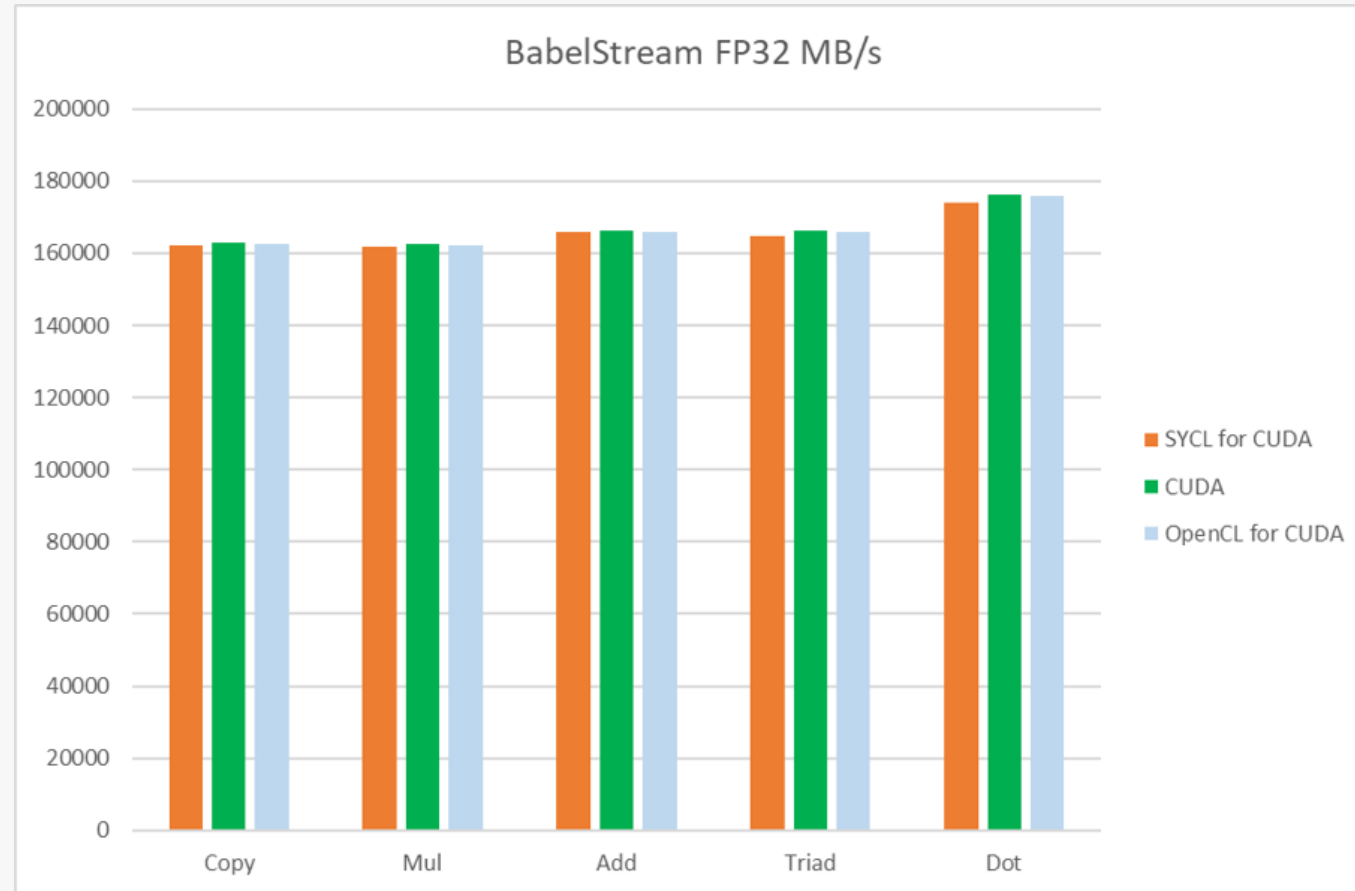LAMMPS
GROMACS
fast, flexible & free

# Performance

- But doesn't 'one size fits all' code equate to poor performance?
  - SYCL code written for GPUs will not be performant on FPGAs.
- Code written for GPUs usually gives good performance across vendors.
- We still need to think about the hardware we are targeting.

# SYCL performance on Nvidia GPUs

- DPC++ calls the underlying CUDA runtime API/native proprietary API.

- Performance is comparable



BabelStream FP32 MB/s

- SYCL for CUDA
- CUDA
- OpenCL for CUDA

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

Run on GeForce GTX 980 with CUDA 10.1
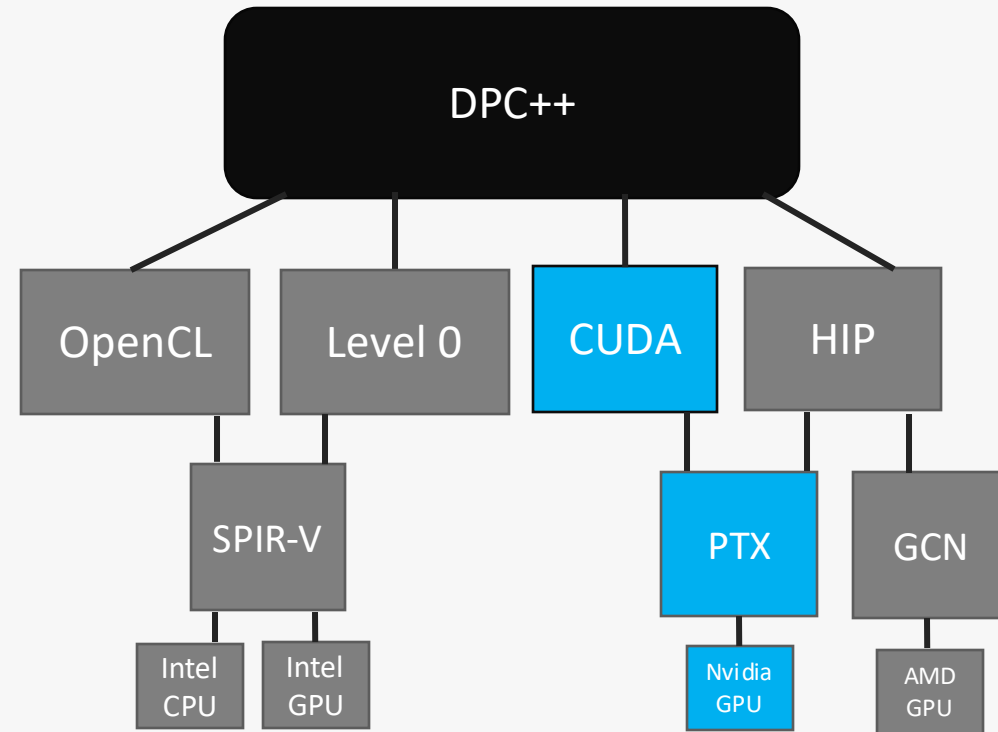
codeplay®

# What is DPC++

- DPC++ is an Intel-backed open-source implementation of SYCL.
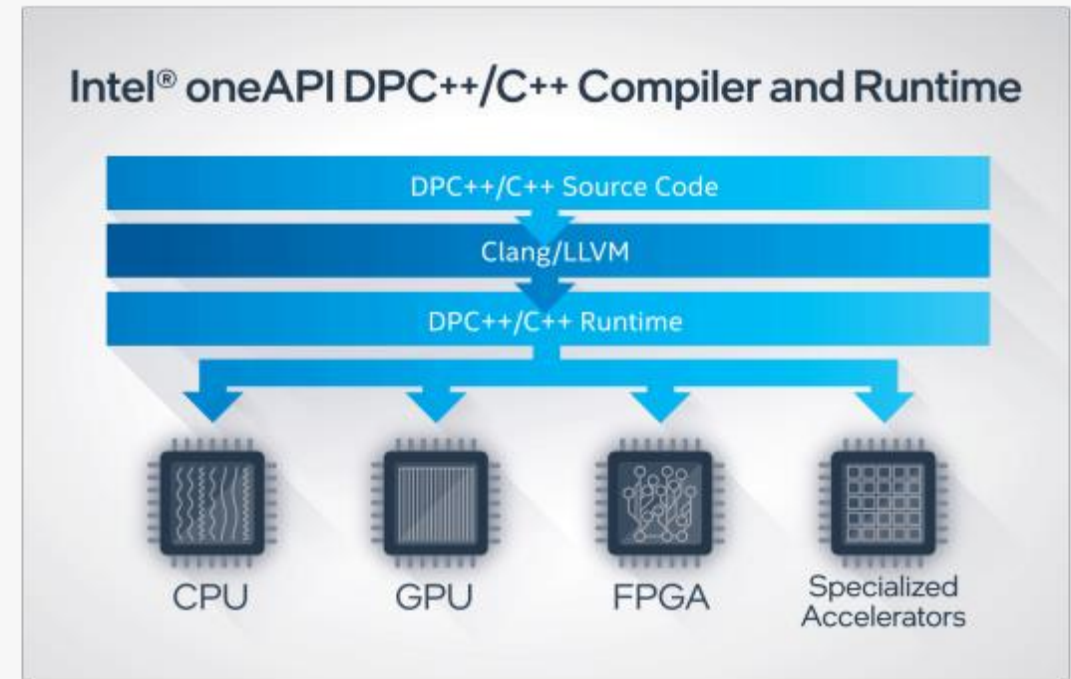- Built and designed by users.

# What is DPC++

- DPC++ currently implements 4 GPU backends:
  - OpenCL
  - Level Zero
  - CUDA
  - HIP

# What is DPC++

- Includes extensions for specific hardware features (tensor cores etc).

- Supports oneAPI libraries oneMKL, oneDNN.



Intel® oneAPI DPC++/C++ Compiler and Runtime

DPC++/C++ Source Code

Clang/LLVM

DPC++/C++ Runtime

CPU  GPU  FPGA  Specialized Accelerators

codeplay®

# Porting your brain to SYCL

- Most CUDA concepts map to a SYCL equivalent

- SYCL is an abstraction layer on top of CUDA & others

| CUDA | SYCL |
|------|------|
| CUstream | sycl::queue |
| CUcontext | sycl::context |
| CUdevice | sycl::device |
| CUevent | sycl::event |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |
|  |  |

# Where to get started with DPC++

- Build from source:
  - https://github.com/intel/llvm
  - https://intel.github.io/llvm-docs/GetStartedGuide.html

- Intel Devcloud:
  - https://devcloud.intel.com/oneapi/get_started/baseToolkitSamples/

# Using DPC++ Compiler

- Check for available devices using sycl-ls

```
→ cuda01 ~ /opt/llvm/build/bin/sycl-ls
[ext_oneapi_cuda:gpu:0] NVIDIA CUDA BACKEND, NVIDIA GeForce RTX 2080 SUPER 0.0 [CUDA 11.4]
[ext_oneapi_cuda:gpu:1] NVIDIA CUDA BACKEND, Tesla K40c 0.0 [CUDA 11.4]
[host:host:0] SYCL host platform, SYCL host device 1.2 [1.2]
→ cuda01 ~ []
```

# First Compilation

```
→ cuda01 ~ /opt/llvm/build/bin/clang++ -fsycl -fsycl-targets=nvptx64-nvidia-cuda hello.cc
warning: linking module '/opt/llvm/build/lib/clang/15.0.0/../../clc/remangled-l64-signed_char.libspirv-nvptx64
--nvidiacl.bc': Linking two modules of different target triples: '/opt/llvm/build/lib/clang/15.0.0/../../clc/r
emangled-l64-signed_char.libspirv-nvptx64--nvidiacl.bc' is 'nvptx64-unknown-nvidiacl' whereas 'hello.cc' is 'n
vptx64-nvidia-cuda'
 [-Wlinker-warnings]
1 warning generated.
→ cuda01 ~ LD_LIBRARY_PATH=/opt/llvm/build/lib ./a.out
Running on device: NVIDIA GeForce RTX 2080 SUPER
→ cuda01 ~ ▯
```
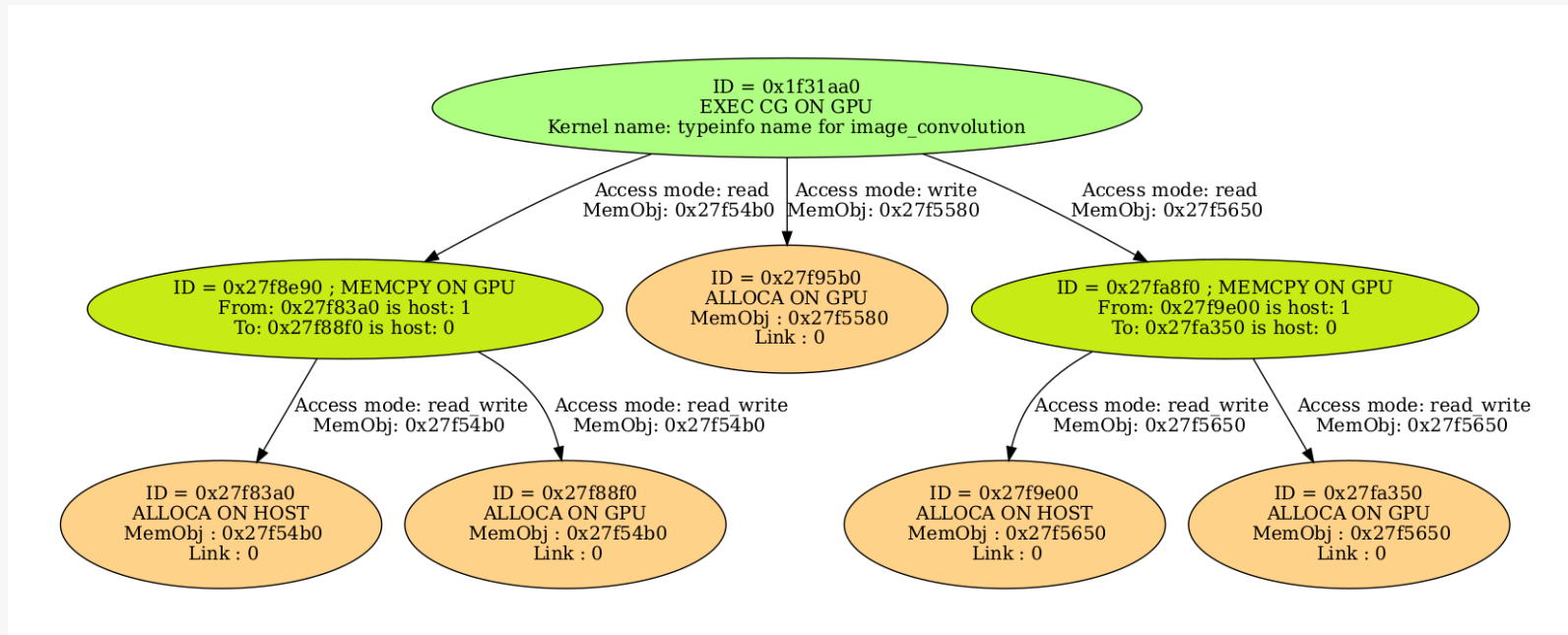
# SYCL Device Filter

- Use SYCL_DEVICE_FILTER to specify one of the devices detected by sycl-ls.

```
→ cuda01 ~ SYCL_DEVICE_FILTER=cuda LD_LIBRARY_PATH=/opt/llvm/build/lib ./a.out
Running on device: NVIDIA GeForce RTX 2080 SUPER
→ cuda01 ~ SYCL_DEVICE_FILTER=host LD_LIBRARY_PATH=/opt/llvm/build/lib ./a.out
Running on device: SYCL host device
→ cuda01 ~ 
```

# Neat DPC++ Environment Variables

- SYCL_PI_TRACE
- SYCL_PRINT_EXECUTION_GRAPH

**codeplay** ©

Enable AI & HPC to be Open, Safe and Accessible to All

# Questions

@codeplaysoft            info@codeplay.com            codeplay.com