**codeplay** ©

Enabling AI to be open, safe and accessible to all

# Building an open, safe, accessible AI & HPC ecosystem

November 2020

Andrew Richards

# Our Brave New World

Single-core CPUs → Multi-core CPUs → Packages of 'chiplets' with multiple optimized cores for each class of algorithm and datatype

Great for processor architects, but how do we write the software?

# How do we write fast software?

*(Your hardware will be obsolete by the time you have optimized it)*

**Hand-code software specifically for the processors we have?**

**Never use any new or innovative processors**

**Use some magical tool that converts any code into fast software for your hardware?**

*(Only works if your magical tool has been pre-programmed to understand the software you just invented)*

*(Those days are over)*

**codeplay**®

# If only there was a proven, practical, solution…

(There is, it's called C++, and it's very widely used)

# C++ has 3 key concepts that enable it to support development of very large, very high performance software

1. Zero-cost abstractions
2. Separation of concerns
3. Composability

codeplay®

# Starting simple: writing a parallel loop

1. We could write a serial loop & hope the compiler parallelizes it

```
void serial_f (float *out,
               const float *in,
               int size) {
  for (int i=0; i<n; i++) {
      out [i] = f (in [i]);
  }
}
```

2. We could write a serial loop & *tell* the compiler to parallelize it

```
void parallel_f (float *out,
                 const float *in,
                 int size) {
#pragma this_loop_is_parallel
    for (int i=0; i<n; i++) {
        out [i] = f (in [i]);
    }
}
```

```
void explicit_parallel_f (float *out,
                          const float *in,
                          int size) {
    parallel_for (0, n, [=] (int i) {
        out [i] = f (in [i]);
    });
}
```
**This is a C++ *zero-cost-abstraction***

3. We could write a parallel loop in C++

***Why would we do it like this?***
⇒ we told the compiler what we want
⇒ now we have complete control
⇒ we can now parallelize very complex software
⇒ now, when we debug the software, it behaves exactly the way we told it to behave

# Writing a parallel loop by hand

```
void parallel_part_f (float *out,
                      const float *in,
                      int start,
                      int end) {
   for (int i=start; i<end; i++) {
      out [i] = f (in [i]);
   }
}
```

4.  Or, we could write the whole thing by hand

***Why would we do it like this?***
⇒ we don't want to maintain this software on
   multiple platforms
⇒ we want to learn how multi-threading works

```
void parallel_threads_f (float *out,
                         const float *in,
                         int size) {
   int part = size / num_cores;
   for (int i=0; i<size; i+=part) {
      create_thread (parallel_part_f,
                        out, in, i,
                        min (size, i + part);
   }
   wait_for_threads_to_complete ();
}
```

codeplay®

# Which of those 4 methods is faster?

(Answer #1: the serial loop is fastest, because I didn't tell you that *n* is 3)

# Lesson: The fastest algorithm varies
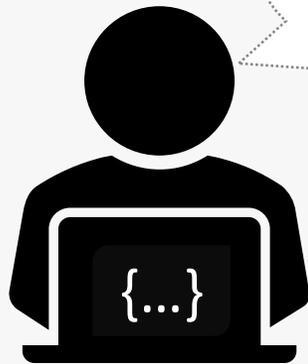
*Your compiler can't know any of this*

1. Performance varies by the size of the data
2. Performance varies by the underlying hardware
3. Performance varies by where the data is
4. Performance varies by what else is running (or could be running) at the same time in the same system

*Your optimized libraries can't know any of this*

# Who knows the answers?

The only person who knows: the size of the data; the hardware it's running on; where the data is, and what else is running on the system is:

The user!

I've written a ten million line program:
Parallelize it yourself

codeplay®

# How do we provide: programs, libraries & tools that can be parallelized and optimized on different systems?

We *separate the concerns*: This is a key modern C++ concept

| | | | |
|---|---|---|---|
| What we want the software to do | The algorithm we will use to calculate the answer | Processor-specific optimizations | How the data is stored |

We can then *independently choose*: the algorithm, the optimizations, which processor each task runs on, how we store the data

# How do we optimize a program where we have *Separated the Concerns?*

Easy: we run it with all the different options and see which runs fastest!

**We break down the optimization problem into three stages:**
1. Writing optimized algorithms, data structures, kernels, schedules
2. Writing our software in a way where we can switch between the different algorithms, data structures, kernels, schedules
3. Choosing the best options for each for the problem we want to solve

# But how do we integrate all the components?

C++ has an answer for this, too: *composability*

- If we write C++ libraries carefully, we can combine them together with user-written code
- If we want to compose across: data formats, different algorithms, different processors, user customization, scheduling, then:
  - ***We need to have the C++ in the same compilation unit, even for different processor cores***
  - ***We call this C++ single-source and it's crucial for making this work on today's heterogeneous multi-core processors***

# This seems like an impossibly big task

But we've already done a lot of the work!

**There are already several C++ libraries that enable this:**
- Kokkos
- Raja
- Eigen
- SYCL-BLAS, SYCL-DNN

**There are already C++ single-source compilers & standards to do this:**
- ISO C++ Parallel STL
- CUDA, HIP
- SYCL – standard for heterogeneous devices
- C++ with OpenMP/OpenAcc
- ComputeCpp, DPC++, triSYCL: implementations of SYCL

**There are already applications doing this:**
- TensorFlow
- A lot of videogame engines

**There are already accelerators supporting this:**
- Most CPUs – out of the box C++
- NVIDIA GPUs – CUDA (& SYCL)
- AMD GPUs – HIP (& SYCL)
- Intel GPUs – DPC++/SYCL
- Renesas R-Car - SYCL
- Imagination Technologies GPUs - SYCL
- ARM Mali GPUs – SYCL
- Intel FPGAs – DPC++/SYCL

# What is SYCL?

- SYCL is a royalty-free vendor-neutral industry standard C++ for parallel software and accelerator processors

- **SYCL takes proven C++ performance ideas & super-charges them for a heterogeneous processing world**

- Now we can:
  - Build our own C++ SYCL compilers for a variety of new processors
  - We can design our own optimizations
  - We can build C++ libraries that can adapt to the performance requirements of lots of different systems
  - We can integrate native compilation for different processors in one source file

codeplay®

# How SYCL handles parallelism

```
cgh.parallel_for<class parallel_demo> (
              cl::sycl::range<1>(n),
              [=](cl::sycl::item<1> i)
{
    out [i] = f (in [i]);
});
```

For more complex parallelism where there are scheduling dependencies, there are a range of options: SYCL requires you to specify where your code *isn't parallel*

- By default, a SYCL `parallel_for` can run entirely parallel

- We define a range to execute in parallel over

- We use a C++ lambda to define the loop body as that's standard now

- It is the job of the programmer to ensure 'f' is safe to run in parallel

- The loop is enqueued and run asynchronously to the CPU thread

- The parallel loop can execute on any SYCL supported core: CPU, GPU, FPGA, DSP, anything programmable

codeplay®

# How SYCL handles data access

```
auto in = inp.get_access<cl::sycl::access::mode::read>(cgh);
auto out = outp.get_access<cl::sycl::access::mode::read_write>(cgh);
cgh.parallel_for<class parallel_demo> (
                  cl::sycl::range<1>(n),
                  [=](cl::sycl::item<1> i)
{
    out [i] = f (in [i]);
});
```

*Access mode specified*

Performance on accelerators is more about data access than compute:

- GPUs have on-board HBM memory and a small amount of fast on-chip SRAM

- DSPs use DMA to transfer data rapidly to a larger amount of on-chip SRAM

- AI accelerators usually have a lot of fast on-chip SRAM

SYCL requires developer specify how to access data: enables maximum performance

codeplay®

# How SYCL handles multiple, different, processors

```
gpu_queue.submit([&](cl::sycl::handler &cgh) {
    auto in = inp.get_access<cl::sycl::access::mode::read>(cgh);
    auto out = outp.get_access<cl::sycl::access::mode::read_write>(cgh);
    cgh.parallel_for<class parallel_demo> (
            cl::sycl::range<1>(n),
            [=](cl::sycl::item<1> i)
    {
        out [i] = f (in [i]);
    });
});
```

*Compiled for CPU by any normal CPU C++ compiler & runs asynchronously on host CPU to enqueue kernels to accelerator*

*This kernel 'name' allows multiple C++ compilers to be stitched together*

*SYCL Device Compiler extracts this kernel and compiles it natively for accelerator processors*

- Both host & device code are compiled via C++ native compilers
- When SYCL goes through OpenCL, it can (optionally) use SPIR-V as the compiler IR
  - ➢ But it's still C++ source compiled to native device ISA
- SYCL device compilers can have per-device extensions
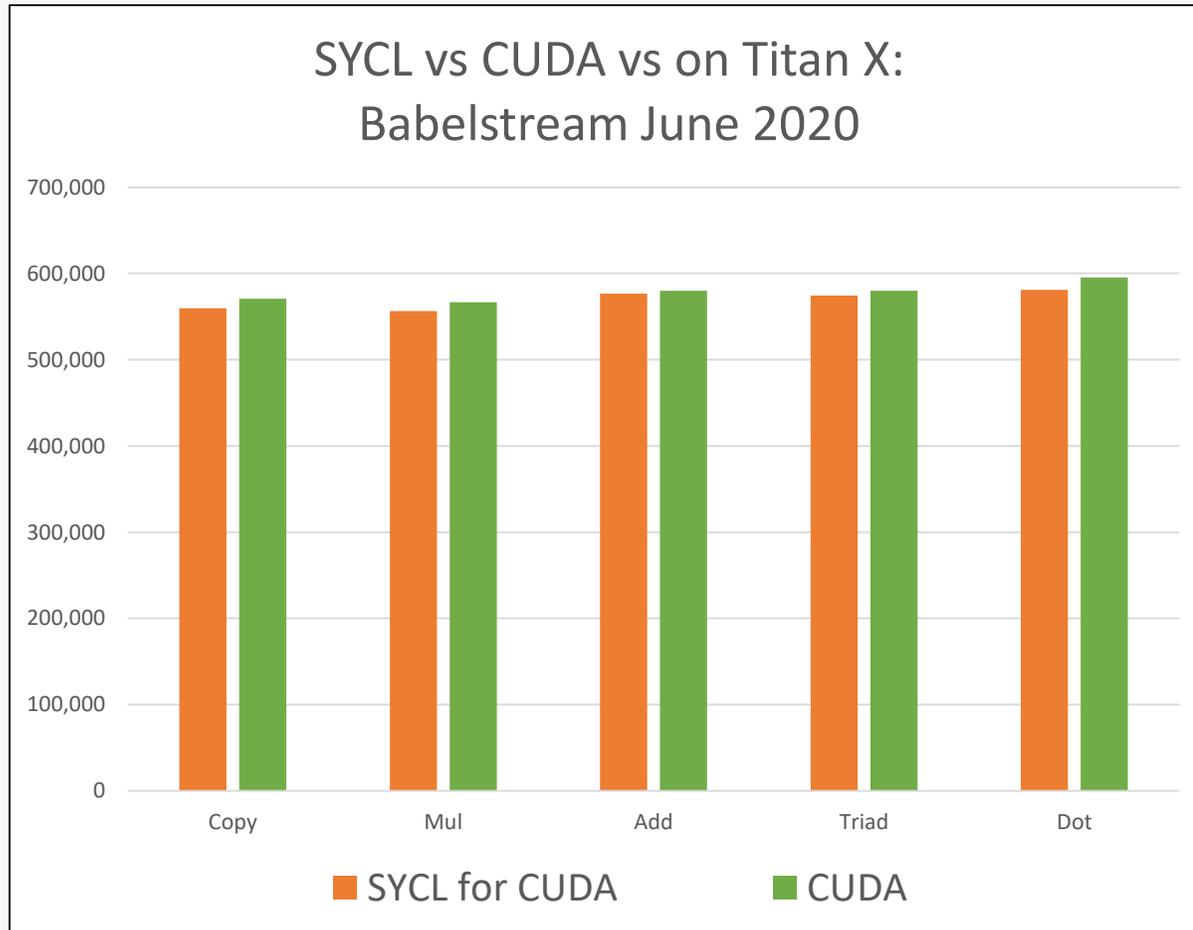- More than one device compiler can compile a single source file

**Combines the benefits of chosen CPU compiler *and* chosen device compiler**

# How SYCL handles processor-specific optimizations

- Most vector instructions and memory models map to SYCL 1.2.1 today

- New instructions or memory systems can be mapped to SYCL extensions – there's a clear mechanism for this

- Then, these processor-specific performance features are *integrated into the template libraries* in an appropriate place
  - ➢ The aim is to enable processor-specific optimizations in as least a disruptive way as possible
  - ➢ Enables us to run the same software with high performance on lots of different processors

codeplay®

# SYCL Performance Equivalent to CUDA

## Bandwidth-Bound Benchmarks

### SYCL vs CUDA vs on Titan X: Babelstream June 2020



Legend: ■ SYCL for CUDA  ■ CUDA

## Compute-Bound Benchmarks

### SYCL vs CUDA on Titan RTX: GEMM October 2020



GEMM Size: 4096*4096, 2048*2048, 1024*1024, 512*512

Time/µs (lower is better)

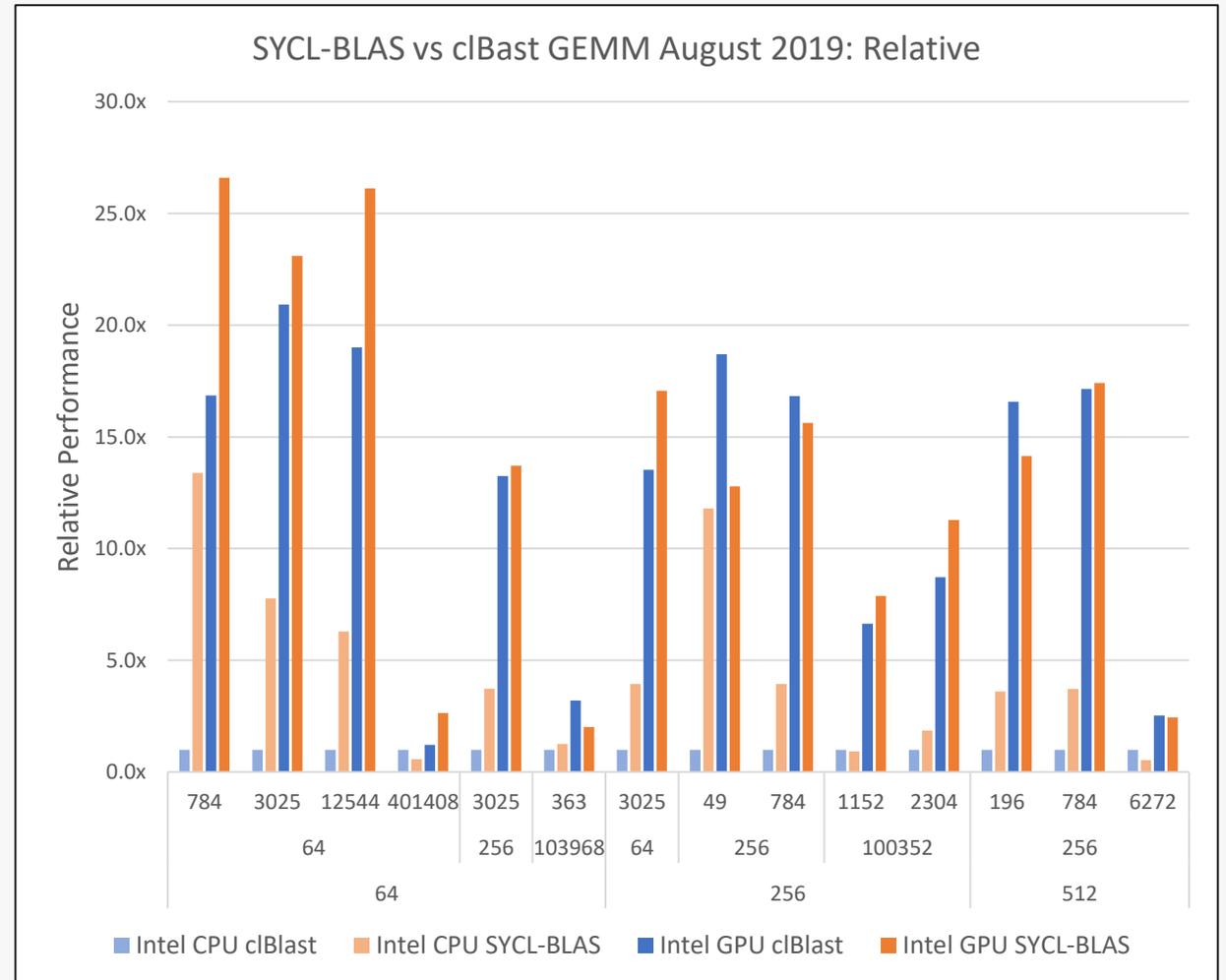Legend: ■ CUDA (cuBLAS)  ■ SYCL (MKL)

codeplay®

# Getting High Performance with Accelerators

1. Move the computation to the data
   - Data movement is slower than computation
   - Compute as much as you can on the data when you have it locally

2. Write computation in a way that can be "moved"
   - Separation of data access from data storage in SYCL enables this
   - Write parameterizable data-movement and composable data + computation

3. Try out a range of different data-movement strategies
   - *"Autotuning"* or *"Empirical Optimization"*
   - For many ML or linear-algebra problems, we find large combinatorial explosion

# SYCL-BLAS Empirical Optimization

- Some processors have hand-optimized libraries integrated into SYCL-BLAS for specific operations

- Can tune the algorithms per-processor:
  - With and without on-chip local memory, work-group size, double buffering, avoid bank conflicts in on-chip memory, cache-line-size , top-level-tile-size, block-level-tile-size

- We continuously run huge number of benchmarks with different optimization parameters for different platforms and different matrix shapes

https://github.com/codeplaysoftware/sycl-blas



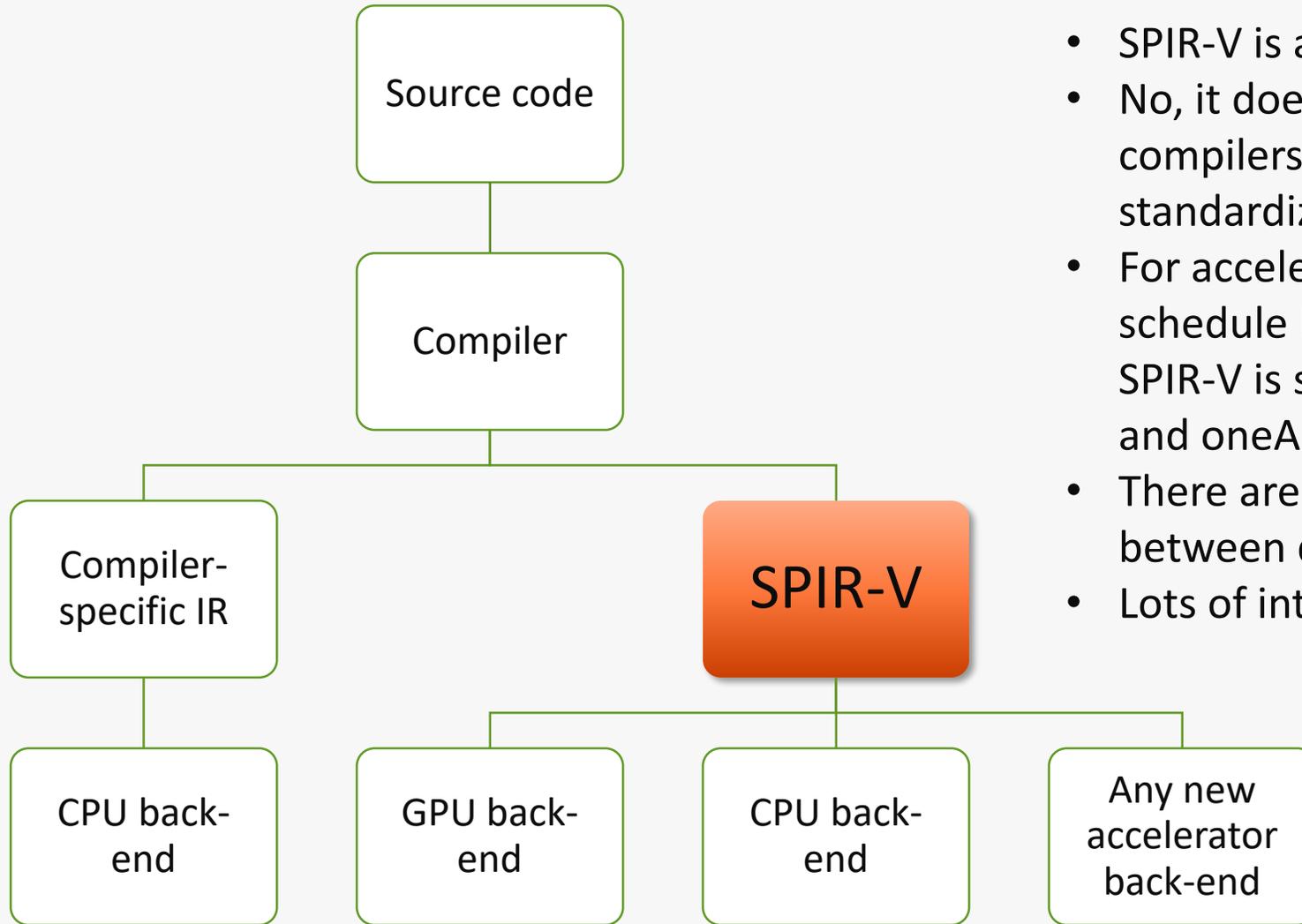SYCL-BLAS vs clBast GEMM August 2019: Relative

# Using SYCL today

- ComputeCpp - Codeplay
  - Closed-source. Community Edition free. Professional Edition fully-supported
  - Supports OpenCL SPIR-V processors (ARM GPU, Renesas R-Car, PowerVR GPU, Intel GPU, +add your own)
- oneAPI/DPC++ - Intel–led, but Codeplay and others contributing
  - Open-source, very active development
  - Intel GPU, NVIDIA GPU, Intel FPGA support released so far
- triSYCL - Xilinx
- hipSYCL – Heidelberg University
  - Open-source active development
  - AMD/NVIDIA GPUs: doesn't go through OpenCL
  - Open-source, less active development now

Check out the growing SYCL ecosystem at sycl.tech & the growing oneAPI ecosystem

codeplay ®

But, you promised a magic compiler that optimizes everything for me!

# SPIR-V

```
                  ┌──────────────┐
                  │ Source code  │
                  └──────┬───────┘
                         │
                  ┌──────┴───────┐
                  │   Compiler   │
                  └──────┬───────┘
          ┌──────────────┴──────────────┐
   ┌──────┴───────┐              ┌───────┴──────┐
   │  Compiler-   │              │    SPIR-V    │
   │ specific IR  │              └───────┬──────┘
   └──────┬───────┘        ┌─────────────┼─────────────┐
   ┌──────┴───────┐  ┌─────┴────┐  ┌─────┴────┐  ┌──────┴──────┐
   │  CPU back-   │  │ GPU back-│  │ CPU back-│  │  Any new    │
   │     end      │  │    end   │  │    end   │  │ accelerator │
   └──────────────┘  └──────────┘  └──────────┘  │  back-end   │
                                                 └─────────────┘
```

- SPIR-V is a standardized compiler IR
- No, it doesn't slow down your code, all compilers go through an IR, SPIR-V just standardizes it
- For accelerators, you need an API to schedule kernels and manage data, SPIR-V is supported by: OpenCL, Vulkan and oneAPI
- There are lots of conversion tools between different flavours of SPIR-V
- Lots of integration with LLVM

Now you can write your own domain-specific compiler + integrate it with other hardware & software

codeplay®

Are you sure this is safe enough to drive a car?

# Functional Safety

- We want to drive cars with accelerate AI software
- Standards make it much easier for us to assess safety
  - We can have independent verification
  - We can make tools to verify standards-based software
  - We can maintain standards-based software for decades
- MISRA C++
  - Coding guidelines for automotive safety being updated for SYCL & acceleration
- RTOS integration
  - We are working with RTOS vendors to integrate safety and acceleration
- This is an industry-wide challenge & we need to work together to solve it
- Great area of research to work in: lots of individual challenges

# What now?

- We're building out this open ecosystem
- Tools & libraries are vastly more useful if they integrate
- We're getting more & more specialized processors for different tasks, so we need software that supports multiple types of processor
- Industry standards make integration much easier: if everyone invents incompatible solutions, nothing works together
- Adoption in: HPC, automotive (ADAS), machine learning, robotics
- These C++ techniques have been proven time & again: they work!
- SYCL is the best C++ for accelerated software: open, fast, modern

codeplay®

codeplay®

Enabling AI to be open, safe and accessible to all