# Who Are We?

- Codeplay is a wholly owned subsidiary of Intel

- Focus on advancing and embracing SYCL and oneAPI

# NVIDIA GPUs are Ubiquitous

- CUDA is proprietary

- Defined by NVIDIA for NVIDIA

- Locked to NVIDIA hardware

- Limited input into direction of CUDA

- Protected by NVIDIA legal terms

**DISCRETE GPU MARKET SHARE (Q1 2022)**

|  | Q1'21 | Q4'21 | Q1'22 |
|---|---|---|---|
| AMD | 19% | 18% | 17% |
| INTEL | n/a | 5% | 4% |
| NVIDIA | 81% | 78% | 78% |

https://wccftech.com/nvidia-amd-gain-gpu-market-share-while-overall-shipments-decrease-by-19-in-q1-2022/
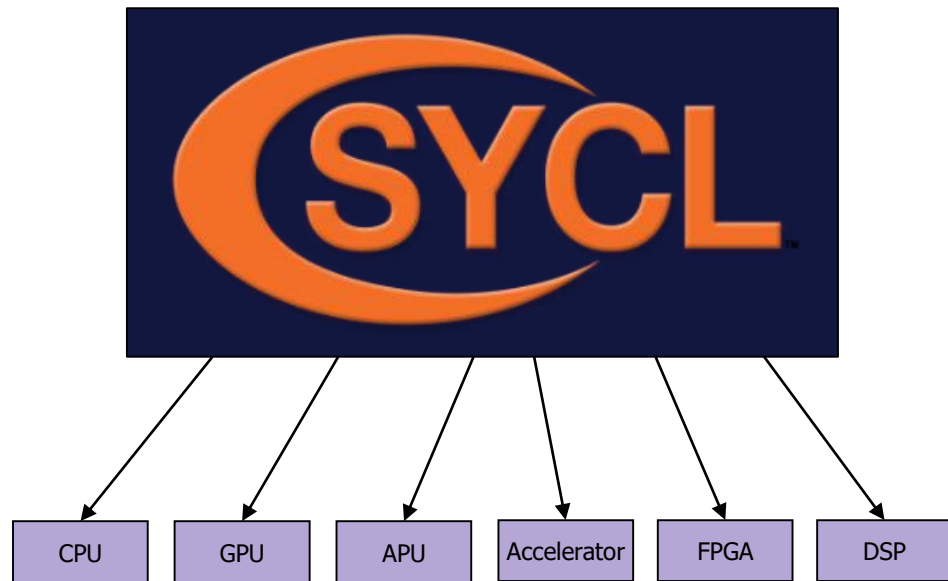
**codeplay**®

# The oneAPI and SYCL Approach

# Open Standards Programming

- SYCL is an open, cross-platform standard programming model based on C++ 17 developed by The Khronos Group

- SYCL supports multiple types of hardware including GPUs, CPUs, and FPGAs from all major vendors

- SYCL is supported by multiple compilers

# SYCL Is a Single-source, High-level, Standard C++ Programming Model

- SYCL can target any device supported by its backend

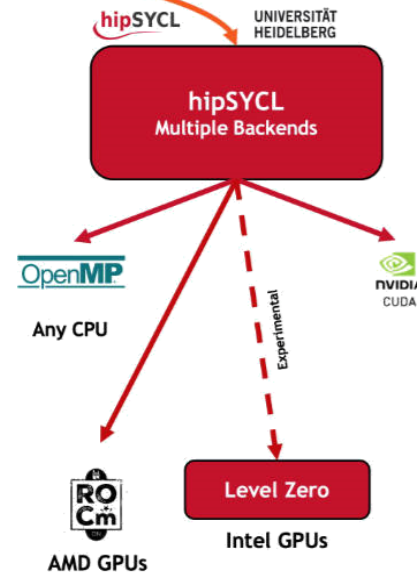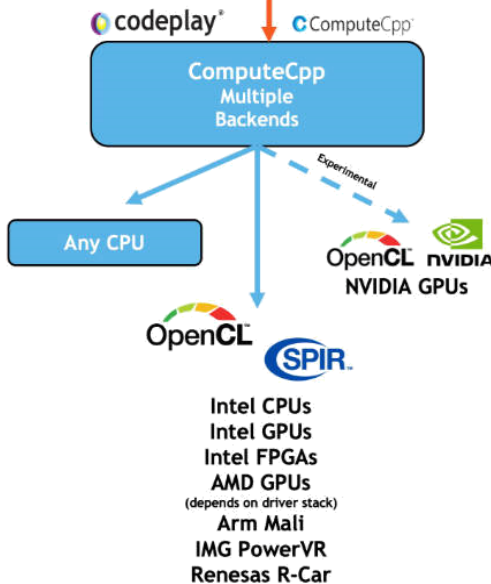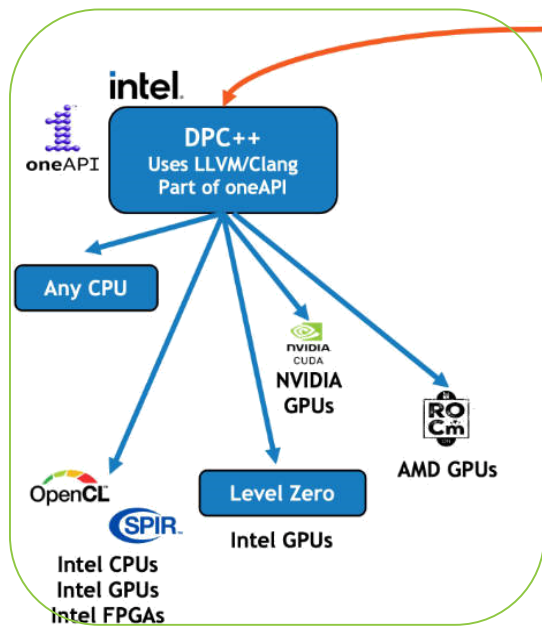- SYCL can target a number of different backends



**SYCL can target a range of heterogeneous platforms**

# SYCL Implementations Under Development

SYCL, OpenCL and SPIR-V, as open industry standards, enable flexible integration and deployment of multiple acceleration technologies

**SYCL Source Code**

SYCL enables Khronos to influence ISO C++ to (eventually) support heterogeneous compute

ISO

**intel oneAPI**

**DPC++**
Uses LLVM/Clang
Part of oneAPI

- Any CPU
- NVIDIA CUDA — NVIDIA GPUs
- ROCm — AMD GPUs
- OpenCL / SPIR — Intel CPUs, Intel GPUs, Intel FPGAs
- Level Zero — Intel GPUs

**codeplay** **ComputeCpp**

**ComputeCpp**
Multiple Backends

- Any CPU
- OpenCL / SPIR — Intel CPUs, Intel GPUs, Intel FPGAs, AMD GPUs (depends on driver stack), Arm Mali, IMG PowerVR, Renesas R-Car
- *Experimental* OpenCL / NVIDIA — NVIDIA GPUs

**hipSYCL** **UNIVERSITÄT HEIDELBERG**

**hipSYCL**
Multiple Backends

- OpenMP — Any CPU
- ROCm — AMD GPUs
- *Experimental* Level Zero — Intel GPUs
- *Experimental* NVIDIA CUDA

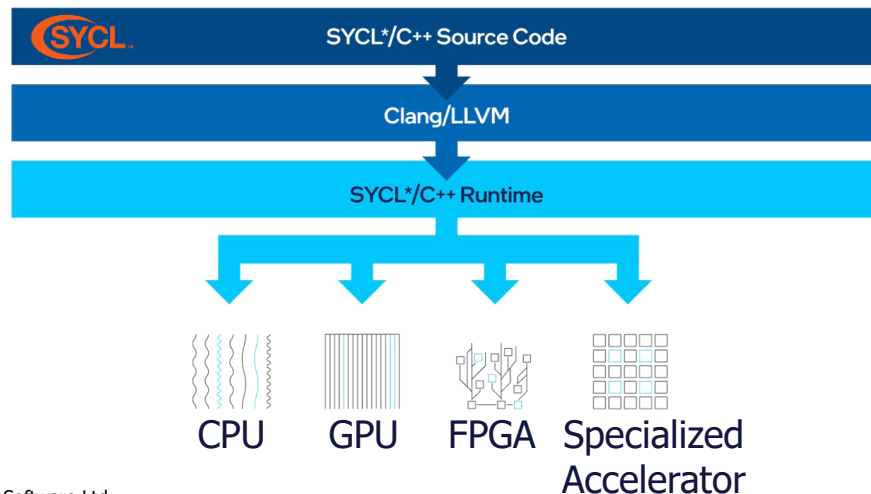Source: https://www.khronos.org/sycl/
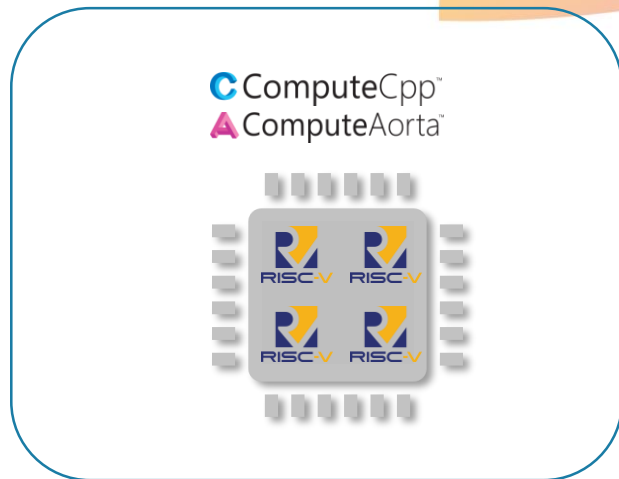
**codeplay** ®

# oneAPI and SYCL

- SYCL sits at the heart of oneAPI

  - Compare with CUDA at the heart of NVIDIA software

- The DPC++ SYCL compiler is open source and based on the LLVM Compiler Infrastructure project



**oneAPI**

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

SYCL*/C++ Source Code

Clang/LLVM

SYCL*/C++ Runtime

CPU    GPU    FPGA    Specialized Accelerator

# SYCL for RISC-V

- Bring industry leading AI and HPC software to the growing range of RISC-V solutions

- Industry-standard compilers & libraries

- Open-source libs and frameworks supported

- Fast migration path of scientific and AI software from NVIDIA GPUs

*"By applying Codeplay's ComputeAorta and ComputeCpp technology, we expect that we can bring state-of-art technology to RISC-V community with our research results."*
*Hideki Sugimoto, CTO NSITEXE Inc, Oct 30th, 2020*

# Migrating from CUDA to Open Standards

# Achieving Multi-Platform Support

Today many programming platforms supported

Future only one software platform needed



CUDA

HIP     SYCL   Proprietary

NVIDIA

AMD    NVIDA   Others
AMD
RISC-V
etc.

CUDA Migration (SYCLomatic)

oneAPI    SYCL

SYCL/CUDA co-exist

NVIDIA

AMD    Intel    Others
RISC-V

codeplay®

# CUDA to SYCL Code Migration Workflow

SYCLomatic / Intel® DPC++ Compatibility Tool assists the migration of code written in CUDA to SYCL once, generating **human readable** code wherever possible

| Nvidia CUDA | Migrate | C++ with SYCL | Build | Deploy |
|---|---|---|---|---|
| CUDA Source Code | SYCLomatic tool | Human Readable **C++ with SYCL** Single Source Code with inline comments | Compilers, Libraries, Analyzers, Debuggers | Run on Multiple Devices (Architecture/VendorAgnostic) |

```
#include
<cuda_runtime.h>

__global__ void
my_cuda_routine()
{
```

90-95%[†] Code Transformed

github.com/oneapi-src/SYCLomatic

Format & Structure Preserved

Tune per Desired Architecture Performance

CPU

GPU

FPGA

Other accel.

† Intel estimates as of September 2021. Based on measurements on a set of 70 HPC benchmarks and samples, with examples like Rodinia, SHOC, PENNANT. Results may vary.

**codeplay®** **Migrate and Deploy Code in 5 Easy Steps**

# Migration Approaches

| Semi-Automatic | Incremental Porting |
|---|---|
| Use conversion tools | Port your kernels alongside existing CUDA code |
| Some engineering work to complete migration | Run CUDA and SYCL code together |

# Semi-Automatic

- **DPCT**
  - Intel released tool

- **SYCLomatic**
  - Open source


- Migrates CUDA code to SYCL
- ~90% of code is migrated

Intel® DPC++ Compatibility Tool Usage Flow



80-90% Transformed

Complete Coding & Tune to Desired Performance

Human Readable DPC++ with Inline Comments

Developer's CUDA* Source

Compatibility Tool

DPC++ Source Code

# Incremental Porting

- Migrating large codebases is a major effort
- It is possible to incrementally migrate CUDA kernels to SYCL
- Run SYCL and CUDA co-existing in same application on NVIDIA GPU

( CUDA + SYCL ) ➜ NVIDIA GPU

<u>Evaluate</u> and transition application code to SYCL and oneAPI

# oneAPI for NVIDIA GPUs and AMD GPUs

- Codeplay contributes plugins

  - Application developers can continue to execute SYCL and oneAPI software on NVIDIA and AMD GPUs

  - Adds support for NVIDIA and AMD GPUs to the oneAPI Base Toolkit



C++ / SYCL™ Source Code

oneAPI Base Toolkit

oneAPI for NVIDIA® GPUs

oneAPI for AMD GPUs (beta)

Uses existing NVIDIA and AMD tools and libraries

*Download from developer.codeplay.com*

codeplay®

# Use Familiar NVIDIA GPU Tools

- Developers can profile code on NVIDIA GPUs with nsys and ncu

- Developers can debug on NVIDIA GPUs with CUDA-gdb

- All of these tools are used with oneAPI in the same way as an application written in CUDA
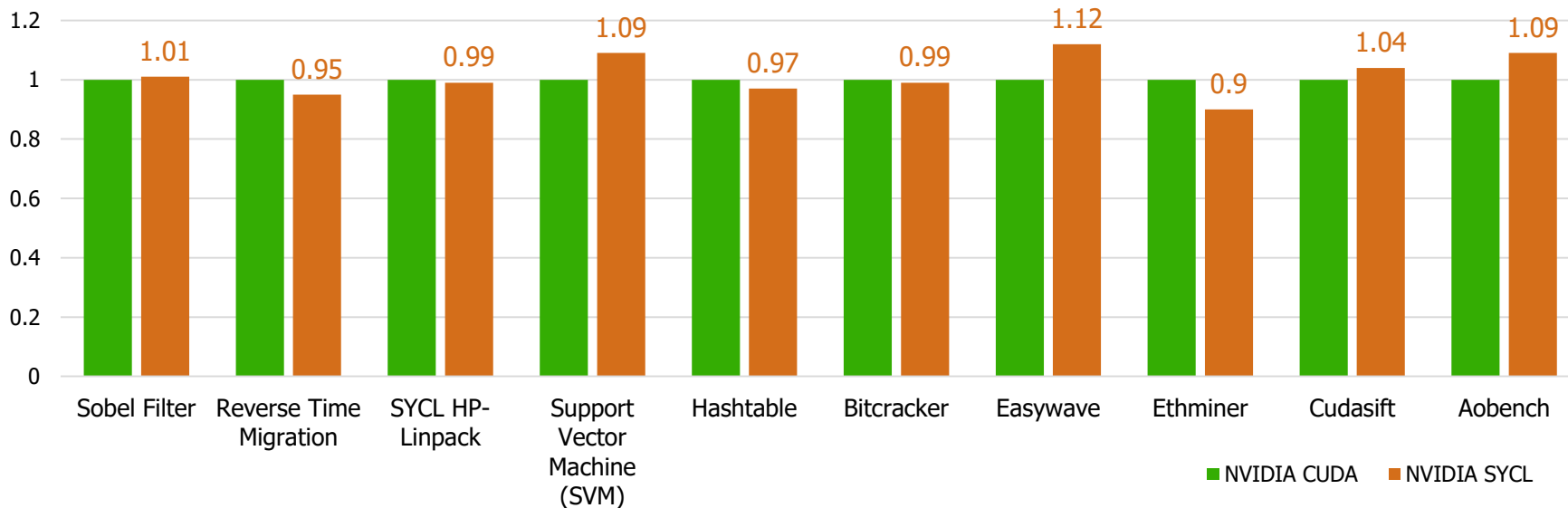
# How to Get the Plugins



Download for free from **developer.codeplay.com**

# Performance

# Relative Performance
# Nvidia SYCL vs Nvidia CUDA on Nvidia GPU

### Relative Performance: NVIDIA CUDA vs NVIDIA SYCL on NVIDIA-A100
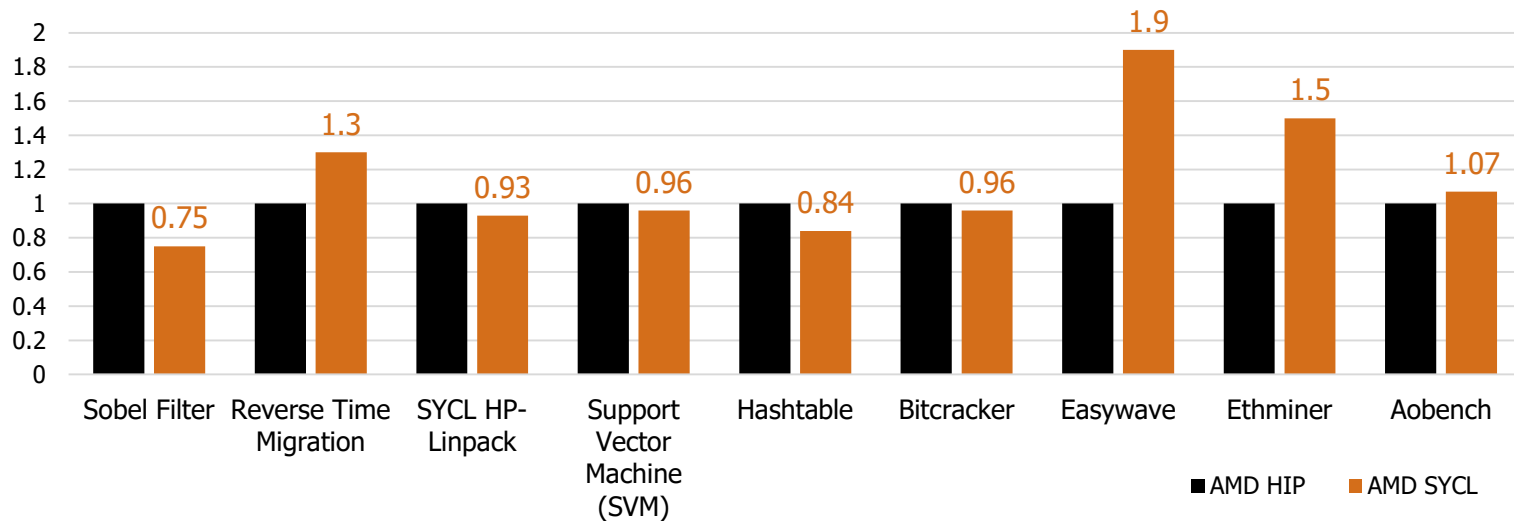### (CUDA=1.00, Higher is Better)



Bar chart values (NVIDIA CUDA = green, NVIDIA SYCL = orange):

| Workload | NVIDIA CUDA | NVIDIA SYCL |
| --- | --- | --- |
| Sobel Filter | 1.00 | 1.01 |
| Reverse Time Migration | 1.00 | 0.95 |
| SYCL HP-Linpack | 1.00 | 0.99 |
| Support Vector Machine (SVM) | 1.00 | 1.09 |
| Hashtable | 1.00 | 0.97 |
| Bitcracker | 1.00 | 0.99 |
| Easywave | 1.00 | 1.12 |
| Ethminer | 1.00 | 0.9 |
| Cudasift | 1.00 | 1.04 |
| Aobench | 1.00 | 1.09 |

# Relative Performance
# AMD SYCL vs AMD HIP on AMD GPU

Relative Performance: AMD **HIP** vs AMD **SYCL** on AMD Instinct MI100 Accelerator
(HIP=1.00, Higher is Better)

# Example Conversion

# Example Conversion : N-Body

- Simulates gravitational interaction in a fictional galaxy

$$\vec{F}_i = -\sum_{i \neq j} G \frac{(\vec{r}_i - \vec{r}_j)}{|\vec{r}_i - \vec{r}_j|^3}$$

- Intentionally simple kernel
- OpenGL for graphics

```cpp
for (int i = 0; i < params.numParticles; i++) {
    vec3 other_pos{pPos.x[i], pPos.y[i], pPos.z[i]};
    vec3 r = other_pos - pos;
    // Fast computation of 1/(|r|^3)
    coords_t dist_sqr = dot(r, r) + params.distEps;
    coords_t inv_dist_cube = rsqrt(dist_sqr * dist_sqr * dist_sqr);

    // assume uniform unit mass
    force += r * inv_dist_cube * (i != id);
}
```



N-body demo running with DPC++ on device: NVIDIA GeForce RTX 3060

# Try It Out for Yourself

- https://github.com/codeplaysoftware/cuda-to-sycl-nbody

- Run it on your own hardware
- Raise issues
- Contribute

- <mark>Visit the demo at the Codeplay booth</mark>



N-body demo running with DPC++ on device: NVIDIA GeForce RTX 3060

oneAPI Community Forum

# What is the oneAPI Community Forum?

**1**

A cross industry group of hardware and software experts

**2**

Defines standard interfaces for accelerator computing

**3**

Multiple specialist technical working groups

**4**

Drives the future of open-standard accelerator computing

# Benefits

## For Software Developers

- Develop with open standards for accelerator computing

- Single code base for multiple processors targets

- Standards and industry defined libraries

- Future proof your software

## For Processor Developers

- Adopt an open standard with existing open-source implementations

- Enable an existing ecosystem of software and educational resources

- Leverage an existing tested and optimized toolchain

## Free and based on open standards

These organizations support the oneAPI initiative for a single, unified programming model for cross-architecture development.
It does not indicate any agreement to purchase or use of Intel's products. *Other names and brands may be claimed as the property of others.

# Conclusions

- NVIDIA with CUDA is dominant and starting place for most AI applications, but locks into one supplier

- SYCL is the best alternative and provides platform independence for heterogeneous processor programming

- oneAPI, based on SYCL, will provide the ecosystem and tools needed

- Start now with oneAPI

  - Experimenting with existing solutions and evolving your own

  - Join oneAPI Community Forum

# Other Performance Research

- Excellent published papers and presentations

  - "State of SYCL – ECP BOF Showcases Progress and Performance"
    by John Russell, February 28, 2023

    - https://www.hpcwire.com/2023/02/28/state-of-sycl-ecp-bof-showcases-progress-and-performance/

  - "SYCL's impact on algorithms, data structures and implementations"
    by Tom Deakin and Tobias Weinzierl, February 27, 2023

    - https://tobiasweinzierl.webspace.durham.ac.uk/research/workshops/siam-cse-23-sycl/  (SeisSol project)

  - "Evaluation of Intel's DPC++ Compatibility Tool in heterogeneous computing"
    by Germán Castaño a, Youssef Faqir-Rhazoui a, Carlos García a b, Manuel Prieto-Matías
    July, 2022

    - https://www.sciencedirect.com/science/article/pii/S0743731522000727?via%3Dihub

- Intel's list of CUDA to SYCL resources

  - https://www.intel.com/content/www/us/en/developer/tools/oneapi/training/migrate-from-cuda-to-cpp-with-sycl.html
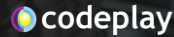
# Codeplay Software

## 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.

## Collaborations

SYNOPSYS

BROADCOM.

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

## Supported Solutions

oneAPI

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

## Markets

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

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

codeplay®