# Annual Meeting

oneAPI Community Forum

# Agenda

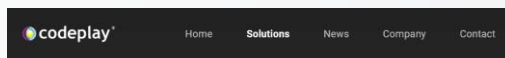| | | |
|---|---|---|
| oneAPI Community Forum Update | Get a roundup of some of the highlights from 2022 as we close off the year and find out about some changes to the organization | Rod Burns, Codeplay Software |
| Breakout Room Discussions | We will break out into rooms on a range of topics for discussion. | •oneAPI 2023 and beyond<br>•AI<br>•Hardware Abstraction<br>•Language |
| Stories from 2022 | Hear from some of our members on what they have contributed to oneAPI this year. | Kevin Harms, Gordon Brown, Kentaro Kawakami and Robert Cohn |
| The future - a cross vendor, industry standard programming | Find out what the future can look like for heterogeneous computing | Andrew Richards, Codeplay CEO |

# 2022 oneAPI Roundup

# oneAPI Initiative to Community Forum

- Intel managed oneAPI initiative → Moved to an open forum
- Codeplay is running the open community forum
- New Spec 1.2 released

- New Steering Committee Members:
  Penporn Konantakool (Google)
  Kevin Harms (Argonne National Lab),
  Antonio Pena (Barcelona Supercomputer Centre)

oneAPI Community Forum announces new members to the steering committee and a new oneAPI Specification 1.2 Release

NOVEMBER 10, 2022

oneAPI Specification 1.2 Release

## oneAPI Specification 1.2 released Nov. 10, 2022

**oneAPI Deep Neural Network Library (oneDNN)**: graph API has been added, which compiles and executes a deep learning computation graph, identifying opportunities for fusing operators and other target-specific optimizations, working closely with industry partners who develop the major frameworks.

**DPC++** (oneAPI's open source SYCL implementation): better management of contexts, queues, and memory management, and enhanced support for images.

**oneMKL**: enhancements for the BLAS libraries with new routines, including support for half/bfloat16, including dense matrix copy and transpose routines as well as updates for BLAS GEMM and GEMM batch.

**oneVPL** added a new API for processing camera RAW data and support for more video color formats.

**Level Zero** added a fabric topology discovery API, support for sRGB, support for image copy with pitch as well as clarifications on existing API.

# SC22 oneAPI Community Forum Meetup

# oneAPI Mission

- Industry defined, open standard-based APIs for accelerated devices

- Bring open source implementations for wide adoption

- Gather industry leaders and contributors to support the mission

- Enable diverse processor designs



oneAPI Community Forum

Define a standards-based open specification

Foster open-source implementations of the specification

oneAPI Projects
Open Source, Open Standard

Libraries

Languages

Hardware Interface

CPU

GPU

FPGA

Other Accel.

Evolving TABs

oneAPI

# 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

oneAPI

# oneAPI Community Forum Steering Committee Formation

- Individuals are being appointed to roles on the Steering Committee
- This will bring on board members of the hardware and software community who have been making contributions to oneAPI

Steering Committee Responsibilities

- Agree and track annual goals
- Vote on proposals for creation of Working Groups
- Vote on proposals for creation of SIGs
- Ratify new versions of the specification
- Approve the Marketing plans

Rod Burns (Codeplay) – Chairperson

Members
Penporn Koanantakool (Google)
Kevin Harms (Argonne National Lab)
Antonio Pena (BSC)
Robert Cohn (Intel)

Others TBA

# Technical Advisory Boards

- Highly engaged
- Great technical discussion
- Providing vital feedback on spec
- Constructive discussions on definitions of the spec

oneMKL

Level Zero

AI

oneIPL

Language

# Governance Goals

- Open membership for groups

- Members vote on specification changes

- Members can propose new groups

- Feedback loop for implementations

These require broader scope and formal specification work alongside more general implementation discussion

oneAPI

# New Technical Group Structure

## Special Interest Groups

- Form from existing TABs
- Facilitate technical discussions
- Can define their own scope
- May or may not feed into a Working Group

## Working Groups

- Deal with specification proposals
- Vote on changes to the specification
- Proposals must be fully formed and draft

oneAPI

# oneAPI Community Forum Organization

**Steering Committee**

| | |
|---|---|
| Rod Burns<br>Chair<br>Codeplay | Robert Cohn<br>Spec Editor<br>Intel |
| Penporn Koanantakool<br>Google | Kevin Harms<br>Argonne National Lab |
| Antonio Pena<br>Barcelona Supercomputer Centre | Others to be announced |

**Marketing Committee**

Alison Richards

**Special Interest Groups (SIGs)**

- Language
- Hardware Abstraction
- AI
- Math

**Working Groups**

To be defined

**oneAPI Specification**

# Next Steps

- TABs will become SIGs
- oneMKL TAB becomes **Math**
- Level Zero TAB becomes **Hardware Abstraction**
- Scope of SIGs will be agreed

- The GitHub project has been updated with governance and other information
  - **https://github.com/oneapi-src/oneAPI-tab**
- Steering Committee will meet to set goals and agree new group processes

| AI | oneAPI Future |
| Language | Hardware Abstraction |
| You | Decide… |

# Contribute to the oneAPI Forum

- Join and lead SIGs and Working Groups
- Submit proposals for features and changes
- Vote on proposals

https://www.oneapi.io/community

Talk to me about the changes and give me your feedback rod@codeplay.com

# Breakout Rooms

oneAPI Future

Rod Burns (Codeplay)

AI
Penporn Koanantakool
(Google)

Language

Robert Cohn (Intel)

Hardware Abstraction

Kevin Harms (ANL)

oneAPI

# oneAPI Future

- What areas are you most interested in?
  - **Introduce yourself and your interest in oneAPI**
- What is your perception of "oneAPI"?
- How "open" do you think oneAPI is?
- What can bring wide adoption to oneAPI

# Stories from 2022

oneAPI

# DOE oneAPI Enablement

**Kevin Harms**
**Argonne Leadership Computing Facility**

# oneAPI Enablement – Aurora



- Support ALCF's Aurora supercomputer
- Collaboration between ALCF and Intel
  - ⬚ Funding via Non-Recurring Engineering (NRE) to support HPC focused improvements
- DPC++
  - ⬚ Support for Intel PVC GPU
- Level Zero
  - ⬚ Reviews and comments to the specification
  - ⬚ Identificaton for multiple GPU nodes
- oneMKL
  - ⬚ Batched interfaces
- oneDNN
  - ⬚ Support for TF and pyTorch + optimizatoins
- oneDAL
  - ⬚ Prioritize list of algorithms to be ported/optimized to GPU

**Compute Node**
2 Xeon Intel® Xeon® CPU Max processors
6 Intel® Data Center GPU Max
Node Unified Memory Architecture
8 fabric endpoints

**GPU Architecture**
Intel XeHPC architecture
High Bandwidth Memory Stacks

**Node Performance**
>130 TF

**System Size**
>9,000 nodes

# oneAPI Enablement – Perlmutter / Polaris

- Support for NERSC's Perlmutter supercomputer
- Collaboration between NERSC, ALCF and Codeplay
  - Focus on support of Nvidia A100 and SM_80 architecture
- Tensor Core support
- Atomics support
- Support for std::complex
  - https://github.com/argonne-lcf/SyclCPLX
- Interoperability of SYCL and OpenMP
- Support for multi-device contexts and peer-to-peer copy
- All work done against intel-llvm (Data Parallel C++)

# oneAPI Enablement – Frontier

- Initial support for OLCF's Frontier system
- Collaboration between OLCF, ALCF and Codeplay
  - Focus on initial implementation supporting AMD MI-50/MI-100 GPUs
  - Frontier uses AMD MI-250x
- Port SYCL CUDA backend to new HIP backend (PI_HIP)
- Support for four specific benchmarks
  - LULESH
  - BabelStream
  - SYCLDslash
  - RSbench
- Approximately 98% of performance for HIP version
- Performance comparisons across Nvidia and AMD hardware
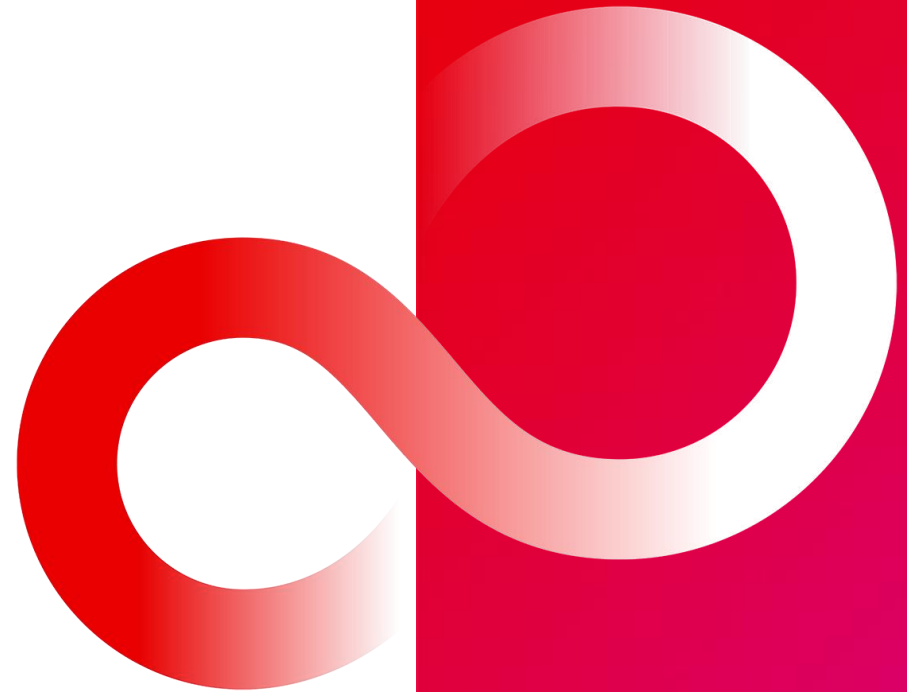
# Acknowledgements

Argonne
NATIONAL LABORATORY

# oneAPI TAB/ annual community forum meeting

Kenaro Kawakami

- kawakami.k@fujitsu.com
- https://github.com/kawakami-k

Computing Laboratory,

Fujitsu Research, Fujitsu Ltd.

# oneDNN for Arm

- **oneDNN has established an essential position in DL S/W stack.**
  - The oneDNN allows users to run DL applications without being aware of device differences.
  - Implementations for Arm have already been merged into oneDNN and freely used by everyone.

| Applications | Image recognition | Object detection | Speech recognition | NLP etc. |
|---|---|---|---|---|

| Framework | TensorFlow | | PyTorch | |
|---|---|---|---|---|

| Library | oneDNN for Arm | oneDNN | cuDNN |
|---|---|---|---|
| | Xbyak_aarch64 | Xbyak | |

| Device | Fujitsu A64FX | Intel CPU | NVIDIA GPU |
|---|---|---|---|

Arm v8.2a+SVE compliant

26

© 2022 Fujitsu Limited

# Recent activity of oneDNN for Arm

- Fujitsu developers, including myself, has been working on SVE 512(512-bit SIMD) of Armv8/9 instruction set support for oneDNN.
  - JIT-based implementation with Xbyak_aarch64
    - This allows efficient use of CPU resources to achieve high performance, just like the implementation for x64 with Xbyak.
  - The computational kernel required for CNN-based DL is almost ready.
    - Kernel type: convolution, batch-norm, relu and its variant, pooling, sum, reorder.
    - Data type: fp32, int8, uint8, int32.
  - Support for v3.0 has also been completed.
    - Release note
      https://github.com/oneapi-src/oneDNN/releases
    - API changes for v3.0
      https://github.com/oneapi-src/oneDNN/blob/rfcs/rfcs/20220815-v3.0-API-cleanup/README.md

# Kernel support status

**FUJITSU**

| Convolution | Batch Norm. | Eltwise | Pooling | ... | Sum | Convolution | Batch Norm. | Eltwise | Pooling | ... | Sum | Convolution | Batch Norm. | Eltwise | Pooling | ... | Sum | Convolution | Batch Norm. | Eltwise | Pooling | ... | Sum | Convolution | Batch Norm. | Eltwise | Pooling | ... | Sum | Convolution | Batch Norm. | Eltwise | Pooling | ... | Sum |

| Calc. kernel generation for AVX512 (512-bit SIMD) | Calc. kernel generation for AVX2 (256-bit SIMD) | Calc. kernel generation for SSE4.1 (128-bit SIMD) | Calc. kernel generation for SVE512 (512-bit SIMD) | Calc. kernel generation for SVE256 (256-bit SIMD) | Calc. kernel generation for SVE128/ASIMD (128-bit SIMD) |

| Xbyak (JIT assembler for x64) | Xbyak_aarch64 (JIT assembler for Arm) |

intel Xeon processor

FUJITSU A64FX™

AWS Graviton3

M2

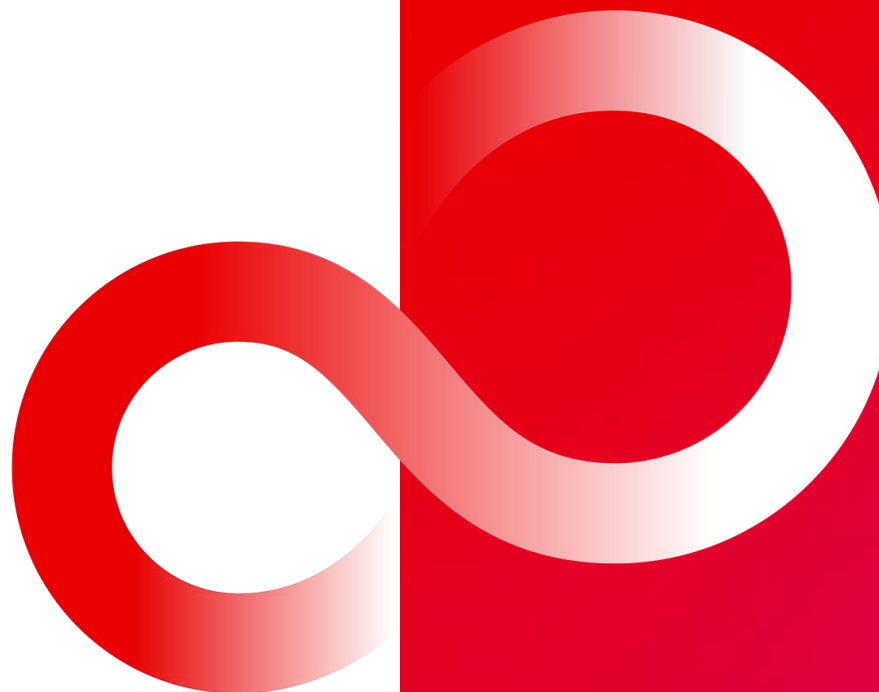**Fully supported**     **Fully supported †**     **Partially supported ‡**

## Let's work together to extend the implementation for Arm devices!

† Data type of fp32/int8/uint8/int32 are supported.
‡ The figure shows the support status of JIT-based implementations.
oneDNN also includes Arm-compute-library-based implementation for Arm devices.

# Thank you

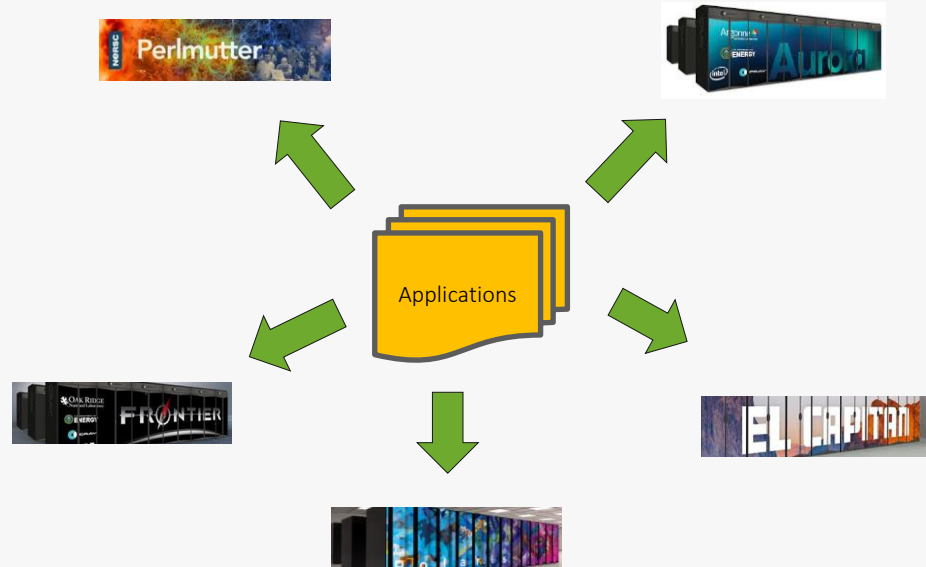# Codeplay's Journey to a Common Hardware Interface

@codeplaysoft  /codeplaysoft  codeplay.com

Gordon Brown, Principal Product Owner, oneAPI
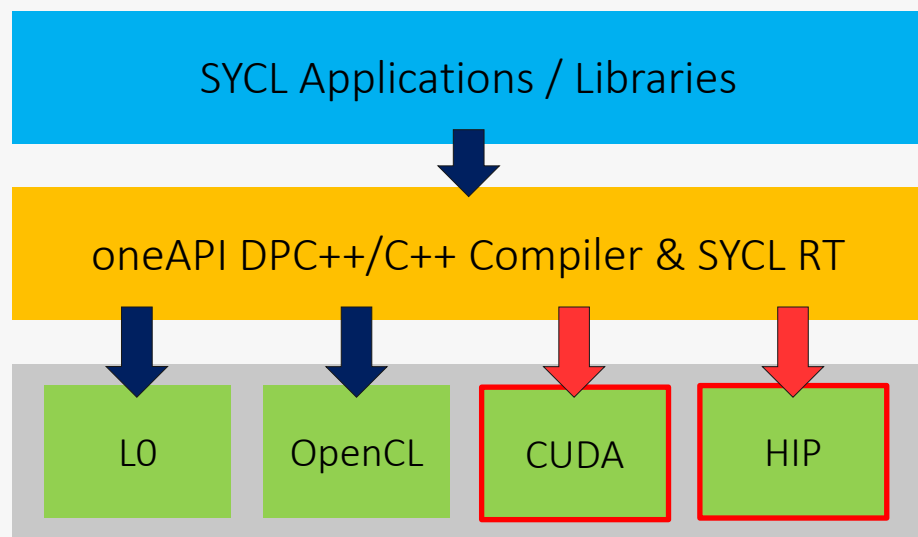
# Extend oneAPI to Support Nvidia & AMD



- Performance portability & long term stability in programming models are more important than ever
- Codeplay have been working with Intel and the US DoE to extend oneAPI to support Nvidia & AMD GPUs

# Current Status

| oneAPI Component | Nvidia GPUs | AMD GPUs |
|---|---|---|
| DPC++/C++ Compiler | ~95% supported | ~50% supported |
| oneDNN | Yes | In progress |
| oneMKL | Yes | In progress |

First binary release of oneAPI for Nvidia/AMD GPUs coming soon!

SYCL Applications / Libraries

oneAPI DPC++/C++ Compiler & SYCL RT

L0

OpenCL

CUDA

HIP

- DPC++ has full support for Nvidia GPUs and partial support AMD GPUs

- oneMKL & oneDNN supports Nvidia GPUs and support for AMD GPUs is in progress

- Extensions have been introduced to support CUDA capabilities: tensor cores, cooperative groups, extended atomics, etc

codeplay®

# What's to Come in 2023

Proposition of DPC++ extensions to SYCL Next

Continued alignment with SYCL and ISO C++

Continued maintenance of the DPC++ CUDA and HIP backends

Further support for the DPC++ HIP backend

Further performance optimizations for the CUDA and HIP backends

Support for additional oneAPI libraries such as oneDPL and oneCCL

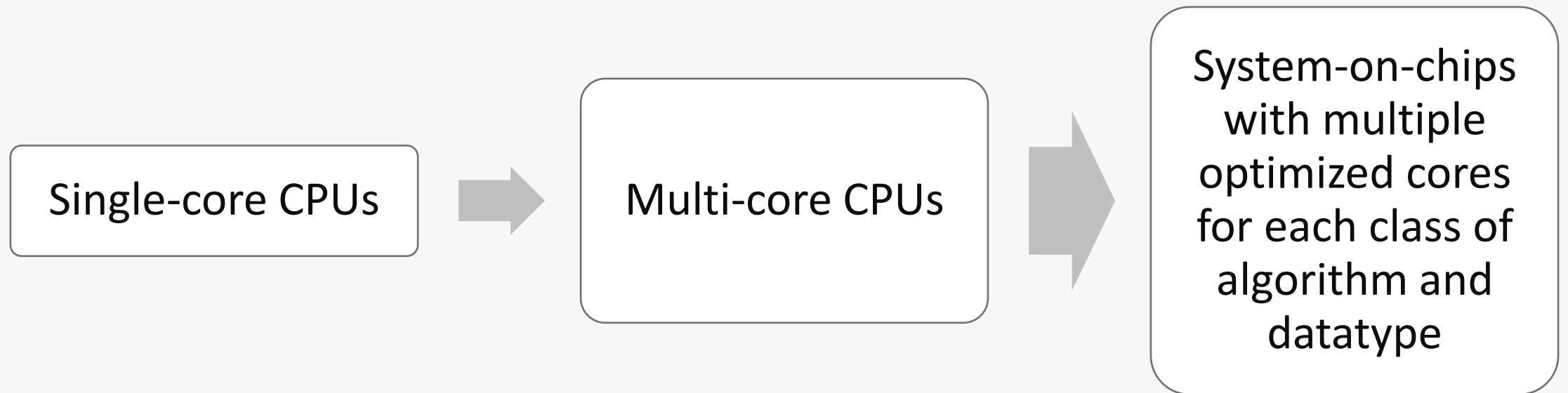# Presentation from Andrew Richards, Codeplay CEO

# oneAPI: Building the Future

December 22

Andrew Richards

codeplay®

Enabling AI & HPC to be open, safe and accessible to all

# Our Brave New World

Single-core CPUs → Multi-core CPUs → System-on-chips with multiple optimized cores for each class of algorithm and datatype

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

codeplay ®

# 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?**
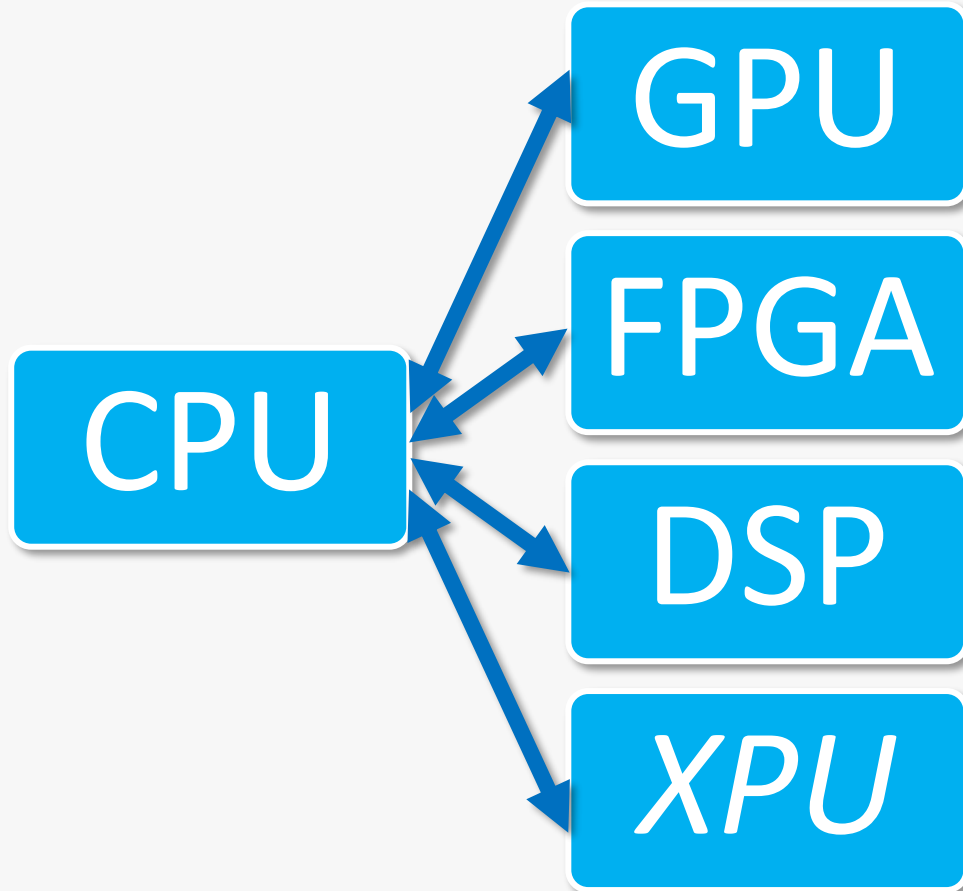
**Never use any new or innovative processors**

*(Those days are over)*

**Use some magical tool that converts any code into fast software for your hardware?**
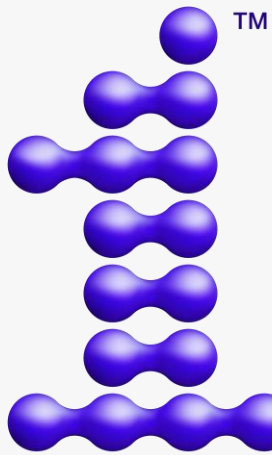
*(Only works if your magical tool has been pre-programmed to understand the software you just invented)*

codeplay®

# Our high performance future



GPU

FPGA

CPU

DSP

*XPU*

**Moving software from the CPU to the whole system**

# What are we doing about it?

1. **Everything** you need to build software from CPUs to XPUs

2. **Open**: mix of open-source, open-standard, open-governance

3. **High performance**: optimized libraries for different processors

**one**API ™

codeplay ®

# Components of oneAPI

**DPC++**
- SYCL compiler
- For C++ programmers who want performance across CPU, GPU, FPGA, *XPUs*

**oneDPL**
- ISO C++ Standard Parallelism library
- For C++ programmers who want performance easily

**oneDNN**
- AI graph compiler
- For people who want high performance deep learning

**oneDAL**
- Data analytics library
- For C++ programmers doing data analytics

**oneTBB**
- CPU-only parallelism library
- For C++ developers who want high performance CPU code

**oneCCL**
- Library for distributed processing across multiple hardware devices
- For big systems

**SPIR-V**
- Virtual instruction set for accelerators
- Enables different compilers & languages

**Level Zero**
- Low-level hardware interface
- Enables more languages to be accelerated with SPIR-V

**oneVPL**
- Accelerated video codec
- For people processing video: encodes, decodes & processes video

**oneMKL**
- Optimized math library
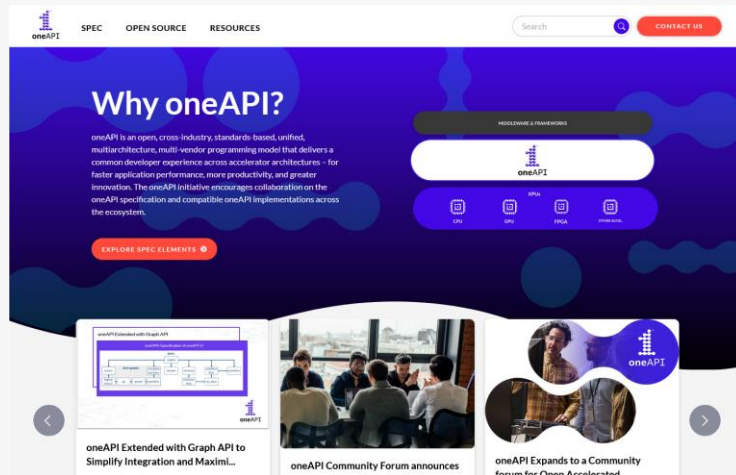- Linear algebra (sparse+ dense), FFT/DFT, random numbers, LAPACK

**Ray Tracing**
- Accelerated ray tracing
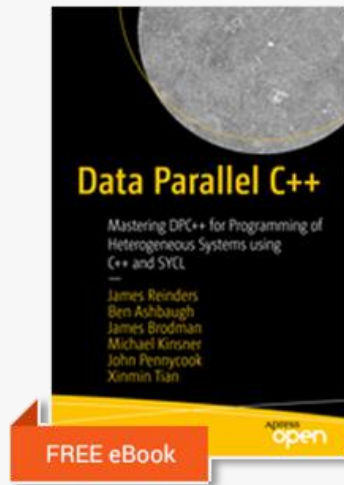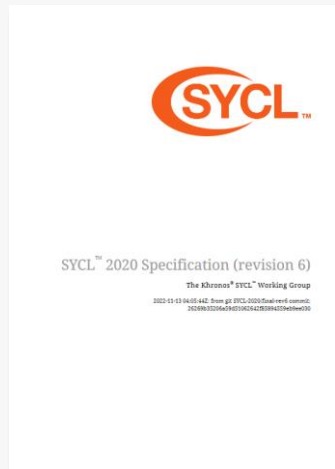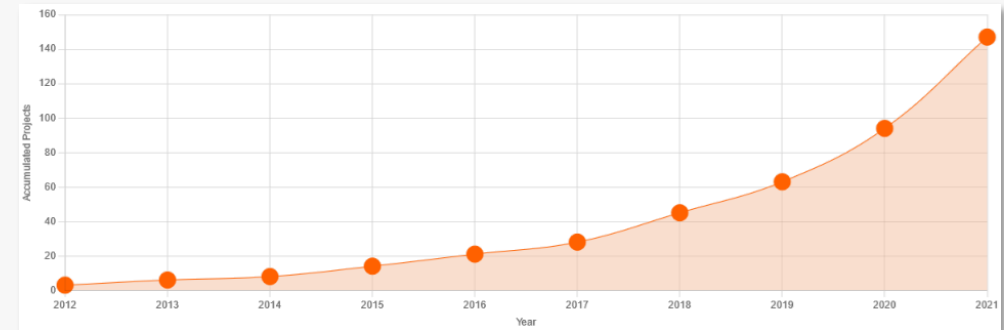- For graphics programmers

*Your Idea Here?*
- It's an open project, so you can add your own concepts
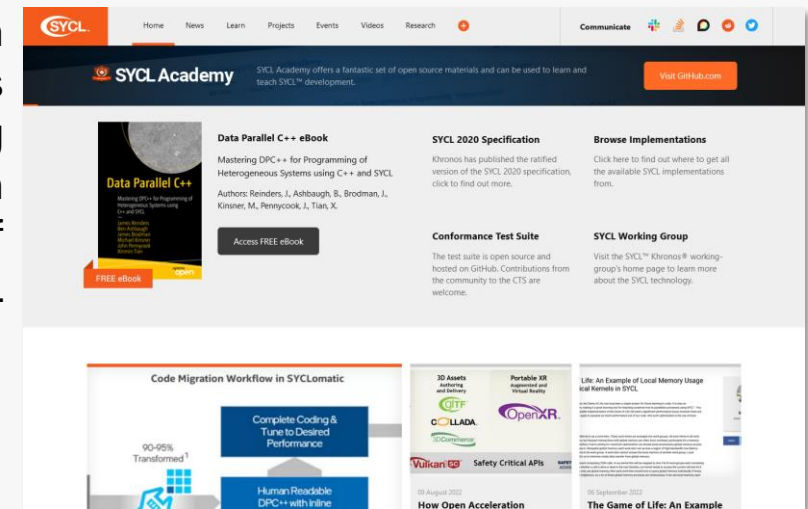
# The oneAPI ecosystem

oneapi.io
Holds all the specifications and tracks everything going on in the ecosystem

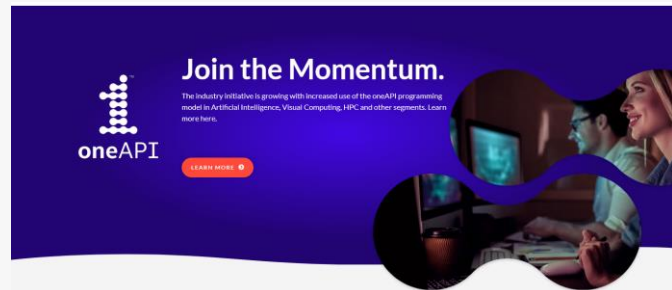Rapidly-growing range of projects using SYCL and oneAPI

The SYCL specification and Data Parallel C++ book document the programming model of oneAPI

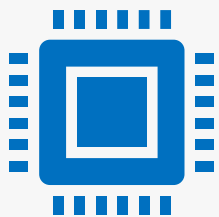sycl.tech tracks everything going on in the world of SYCL

# The future



- **Come and join us!**
  - You can join both the oneAPI Community Forum and add your input to drive the project forwards

- **Build out the performance-portable software ecosystem**
  - There's a huge opportunity to build the performance-portable frameworks of the future using oneAPI & SYCL, including ISO C++



- **Bring your own hardware to developers**
  - You can design your own chip for oneAPI, so you can accelerate all this software with your own hardware innovations

# What interests me

## Safety

- We're increasingly seeing AI being used to do things requiring safety such as driving cars or operating medical systems
- There's huge benefits to using AI in these ways
- But: there are huge dangers
- We're working on solving some of these challenges

## Performance portability

- In 2023, we'll have a range of hardware supporting oneAPI and SYCL
- SYCL doesn't magically give you performance everywhere
- But: you can build performance frameworks using SYCL
- Let's build performance portable frameworks in 2023

## Software-first-hardware

- When we build software with oneAPI we can target a variety of hardware
- It's very hard to design processors for very complex software
- Instead of designing software-for-hardware, let's design hardware-for-software
- We're particularly working with RISC-V to enable this

codeplay®

# codeplay

Enabling AI & HPC to be open, safe and accessible to all

# What Happens Next?

- SIG meetings for 2023 will be scheduled in January
- We invite your proposals for new SIGs during 2023
- Processes and mechanisms will be updated
- The Steering Group will meet to set the goals for 2023

Tell us what you want to see in 2023 for the oneAPI Community Forum
oneapi@codeplay.com rod@codeplay.com