Multiarchitecture Programming for Accelerated Compute
Freedom of Choice for Hardware

# oneAPI Industry Initiative & Intel® oneAPI Tools

蔺杰 (Auber Lin)

SATG/DSE/DTCE/APCAP/PRC Customer Engineering Team

Nov. 2023

# Agenda

- oneAPI Goal
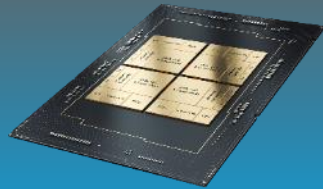- Intel oneAPI product – Toolkits

intel.

# oneAPI Goal

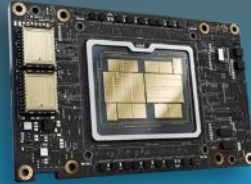# Modern Applications Demand Increased Processing

**Diverse accelerators needed to meet today's performance requirements:**

48% of developers target heterogeneous systems
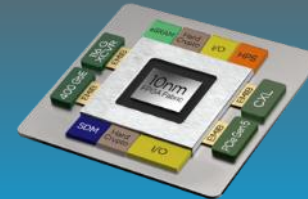that use more than one kind of processor or core[1]

**CPU**

**GPU**

**FPGA**

**Other Accelerators**

**Developer Challenges: Multiple Architectures, Vendors, and Programming Models**
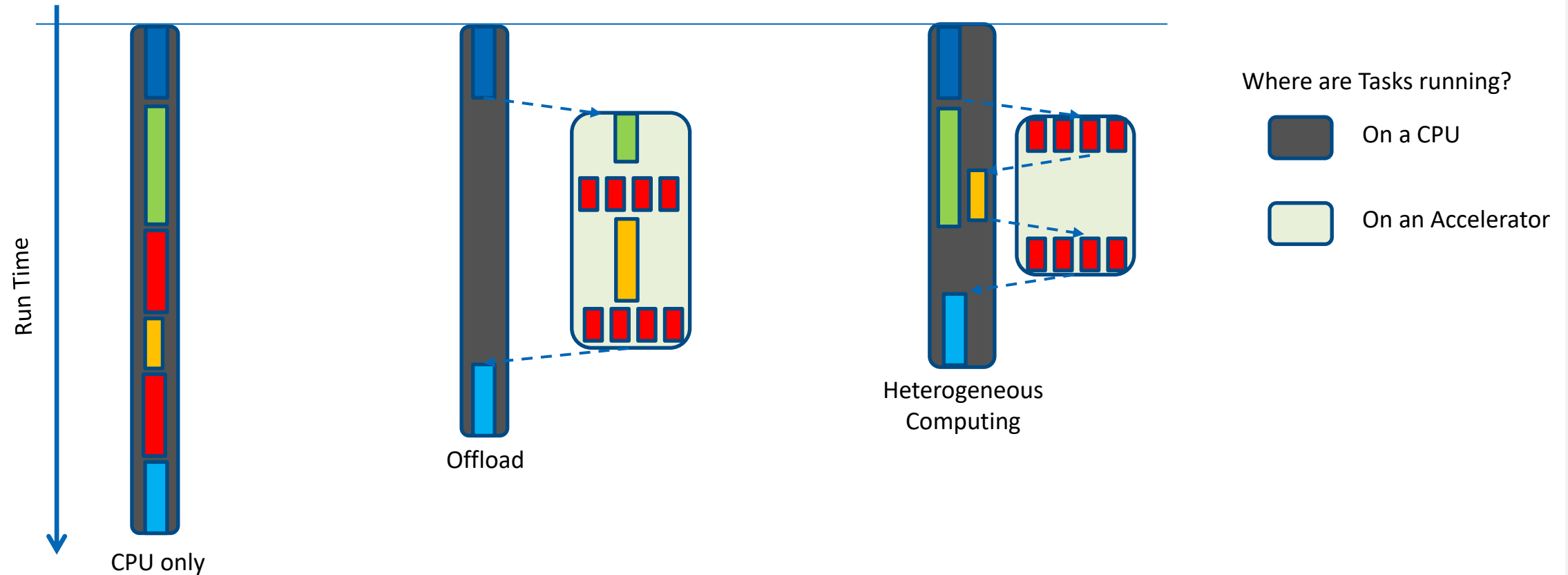
**oneAPI**

**Open, Standards-based, Multiarchitecture Programming**

# Offload v.s. Heterogenous Computing

- Offload

The CPU moves work to an accelerator and waits for the answer



Run Time

CPU only

Offload

Heterogeneous Computing

Where are Tasks running?

On a CPU

On an Accelerator

# 异构计算软件生态现状

| Domain | | Example Workloads | Software Ecosystem |
|---|---|---|---|
| AI | DL | • Training<br>•  Inference | High level frameworks targeting *proprietary interfaces* |
| | ML | • SVM<br>•  k-means | No established standards |
| | Big Data | • Data mining<br>• Analytics | Distributed computing on Xeon |
| HPC | | • Simulation<br>• Modeling | *CUDA based GPU support* |
| Visual | | • Video<br>• Transcode | Well-established industry standards |

- Domain Specific Languages
- Modular, pluggable
- One size does not fit all
- Sparse and unstructured
- Distributed

# oneAPI Industry Initiative

## Break the Chains of Proprietary Lock-in

### Freedom to Make Your Best Choice

- C++ programming model for multiple architectures and vendors
- Cross-architecture code reuse for freedom from vendor lock-in
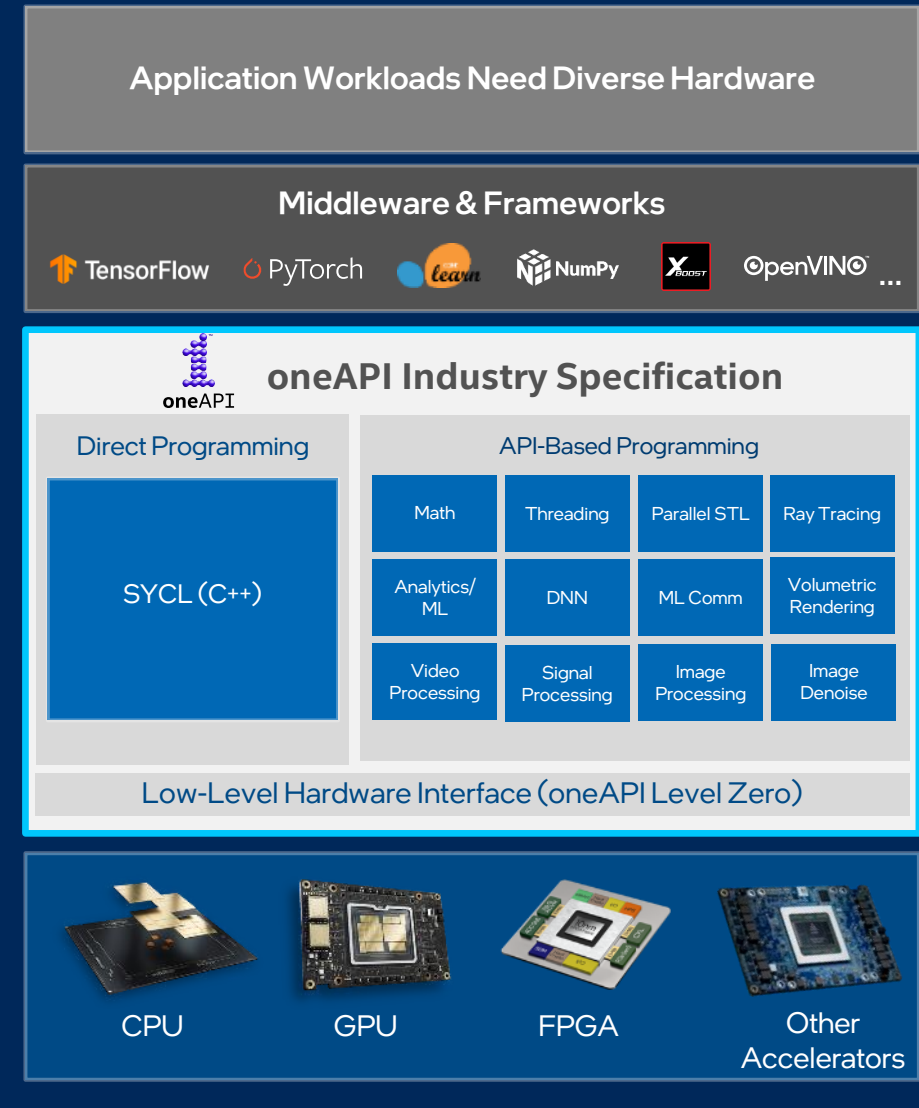
### Realize all the Hardware Value

- Performance across CPU, GPUs, FPGAs, and other accelerators
- Expose and exploit cutting-edge features of the latest hardware

### Develop & Deploy Software with Peace of Mind

- Open industry standards provide a safe, clear path to the future
- Interoperable with familiar languages and programming models including Fortran, Python, OpenMP, and MPI
- Powerful libraries for acceleration of domain-specific functions

The productive, smart path to freedom for accelerated computing from the economic and technical burdens of proprietary programming models

**Application Workloads Need Diverse Hardware**

**Middleware & Frameworks**

TensorFlow    PyTorch    learn    NumPy    XGBOOST    OpenVINO  ...

**oneAPI Industry Specification**

Direct Programming

SYCL (C++)

API-Based Programming

| Math | Threading | Parallel STL | Ray Tracing |
| Analytics/ML | DNN | ML Comm | Volumetric Rendering |
| Video Processing | Signal Processing | Image Processing | Image Denoise |

Low-Level Hardware Interface (oneAPI Level Zero)

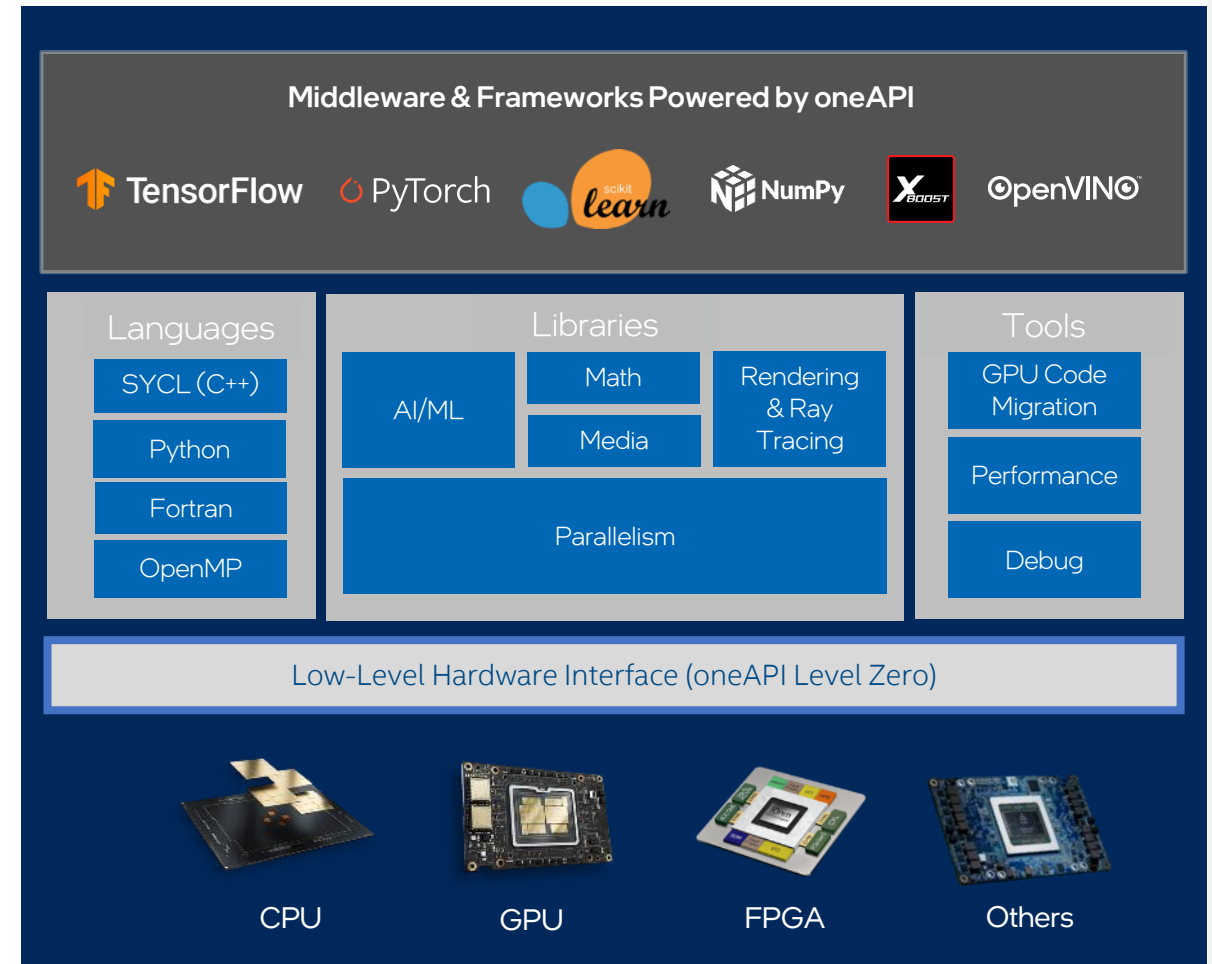CPU    GPU    FPGA    Other Accelerators

# Intel® Developer Tools Supporting oneAPI

**A complete set of proven tools expanded from CPU to accelerators**

- Advanced compilers, libraries, and analysis, debug, and porting tools

- Full support for C, C++ with SYCL, Python, Fortran, MPI, OpenMP

- Intel® Advisor determines device target mix before you write your code

- Intel's compilers optimize code to take full advantage of multiarchitecture workload distribution.

- Intel® VTune™ Profiler analyzes hotspots to optimize code performance

- Intel AI tools support acceleration of major deep learning and machine learning frameworks
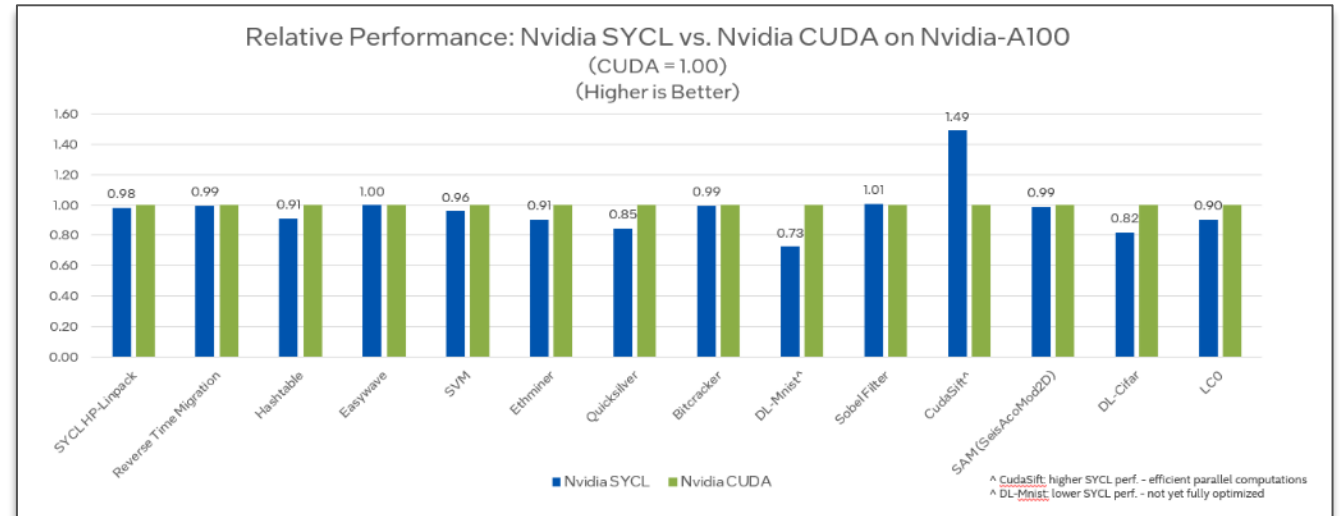
**Middleware & Frameworks Powered by oneAPI**

TensorFlow · PyTorch · scikit learn · NumPy · XBOOST · OpenVINO™

**Languages**
- SYCL (C++)
- Python
- Fortran
- OpenMP

**Libraries**
- AI/ML
- Math
- Media
- Rendering & Ray Tracing
- Parallelism

**Tools**
- GPU Code Migration
- Performance
- Debug

Low-Level Hardware Interface (oneAPI Level Zero)

CPU · GPU · FPGA · Others

# Accelerating Choice with SYCL*

## Khronos Group Standard

- Open, standards-based

- Multiarchitecture performance

- Freedom from vendor lock-in

- Comparable performance to native CUDA on Nvidia GPUs

- Extension of widely used C++ language

- Speed code migration via open source SYCLomatic or Intel® DPC++ Compatibility Tool

Relative Performance: Nvidia SYCL vs. Nvidia CUDA on Nvidia-A100
(CUDA = 1.00)
(Higher is Better)

| | **Architectures** | Intel | Nvidia | AMD CPU/GPU | RISC-V | ARM Mali | PowerVR | Xilinx |

# SYCLomatic: CUDA* to SYCL* Migration Made Easy

**Choose where to run your software, don't let the software choose for you.**

| 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/Vendor Agnostic) |

```
#include
<cuda_runtime.h>

__global__ void
my_cuda_routine()
{
```

90-95%[1] Code Transformed

github.com/oneapi-src/SYCLomatic

Format & Structure Preserved

Tune per Desired Architecture Performance

CPU
GPU
FPGA
Other accel.

Open source SYCLomatic tool assists developers migrating code written in CUDA to C++ with SYCL, generating **human readable** code wherever possible

~90-95% of code typically migrates automatically[1]

Inline comments are provided to help developers finish porting the application

Intel® DPC++ Compatibility Tool is Intel's implementation, available in the Intel® oneAPI Base Toolkit

[github.com/oneapi-src/SYCLomatic](github.com/oneapi-src/SYCLomatic)

# Codeplay oneAPI Plug-ins for Nvidia* & AMD*

Support for Nvidia & AMD GPUs to Intel® oneAPI Base Toolkit

## oneAPI for NVIDIA & AMD GPUs

- Free download of binary plugins to Intel® oneAPI DPC++/C++ Compiler:
- Nvidia GPU
- AMD beta GPU
- No need to build from source!
- Plug-ins updated quarterly in-sync with SYCL 2020 conformance & performance

## Priority Support

- Available through Intel, Codeplay & our channel
- Requires Intel Priority Support for Intel® oneAPI DPC++/C++ Compiler
- Intel takes first call, Codeplay delivers backend support
- Codeplay provides access to older plug-in versions
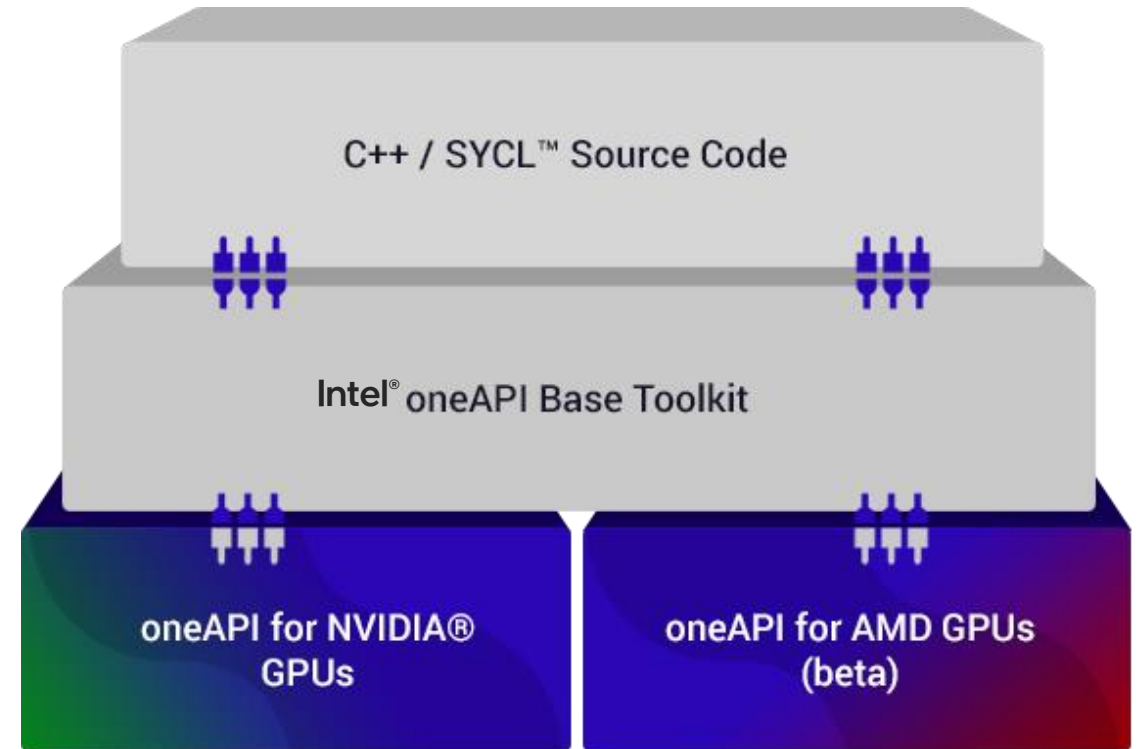
C++ / SYCL™ Source Code

Intel® oneAPI Base Toolkit

oneAPI for NVIDIA® GPUs

oneAPI for AMD GPUs (beta)

Image courtesy of Codeplay Software Ltd.

Nvidia GPU plug-in

AMD GPU plug-in

Codeplay blog

Codeplay press release

# oneAPI Libraries

| Domain | Name | Description | Open Spec | Open Source |
|---|---|---|---|---|
| Parallel Programming | oneDPL | Data Parallel C++ Library including Parallel STL | Yes | Yes |
| | oneTBB | Threading Building Blocks | Yes | Yes |
| AI & ML | oneDNN | Deep Neural Networks | Yes | Yes |
| | oneCCL | Collective Communications | Yes | Yes |
| | oneDAL | Data Analytics and Machine learning | Yes | Yes |
| Math | oneMKL | Math Kernels: linear algebra, FFT, random numger generation | Yes | Partial |
| Video | oneVPL | Video Processing: encode, decode, transcode | Yes | Yes |
| Ray Tracing | Embree, VKL, OID, OSPRay | Geometric & Volumetric Ray Tracing, Image Denoise, Scalable Rendering | Yes | Yes |

intel

# oneAPI Industry Momentum

## End Users

@rchanan, FOUNDRY MODO, BioDataAnalysis GmbH, GeoEast, BENTLEY, CINESITE, LAIKA, Accrad Accelerated Radiology, GE, kt, AUTODESK ARNOLD, 中国石油集团东方地球物理勘探有限责任公司 BGP INC.,CHINA NATIONAL PETROLEUM CORPORATION, iOmniscient Autonomous AI, EURECOM Sophia Antipolis, Brightskies, PHILIPS, WeBank, 大势智慧 DASPATIAL, SAMSUNG MEDISON, TANGENTANIMATION, allegro.ai, ILLUMINATION MACGUFF, Verizon, MediaKind

## National Labs

CERN openlab, Argonne NATIONAL LABORATORY, UT-BATTELLE Oak Ridge National Laboratory, CINECA, Peraton Labs, Laboratório Nacional de Computação Científica, lrz Leibniz Supercomputing Centre of the Bavarian Academy of Sciences and Humanities, SANKHYA SUTRA

## ISVs & OSVs

codeplay, CANONICAL, DITI CERTIFACE, ANACONDA, Ansys, MPCDF MAX PLANCK COMPUTING & DATA FACILITY, AI SINGAPORE, SAS, ES, Red Hat OpenShift Data Science, KATANA GRAPH, CGG, AIBLE, CHAOSGROUP, MAXON, spirent, FOUNDRY, mercenaries engineering, Guerilla, Hisense Medical 海信医疗, SENAI CIMATEC SISTEMA FIEB, SAP, E4 COMPUTER ENGINEERING, YUAN, GIGASPACES innovate with confidence, MEGH COMPUTING, KFBIO, Tech Mahindra, AsiaInfo 亚信科技, SUSE, vmware, UNITED IMAGING

## OEMs & SIs

BittWare a molex company, Hewlett Packard Enterprise, DELL Technologies, Atos, Lenovo, HCL, rENIAC, MEGWARE SUPERCOMPUTING • TECHNOLOGY

## Universities & Research Institutes

LOBACHEVSKY UNIVERSITY, TECHNION Israel Institute of Technology, UNIVERSITY OF CAMBRIDGE, ICT 中国科学院计算技术研究所 INSTITUTE OF COMPUTING TECHNOLOGY, CHINESE ACADEMY OF SCIENCES, Ben-Gurion University of the Negev, Berkeley, UC DAVIS UNIVERSITY OF CALIFORNIA, THE UNIVERSITY OF TENNESSEE KNOXVILLE, UNIVERSIDAD COMPLUTENSE MADRID, CTG, OLD DOMINION UNIVERSITY, 北京大学软件与微电子学院 School of Software & Microelectronics, ZIB, ILLINOIS, THE UNIVERSITY OF UTAH, FACULTY OF MATHEMATICS AND PHYSICS Charles University, PURDUE UNIVERSITY Elmore Family School of Electrical and Computer Engineering, UNIVERSITY OF OREGON, TÉCNICO LISBOA, inesc id lisboa, University College London, Durham University, University of Stuttgart Germany, Indian Institutes of Technology Delhi / Kharagpur / Roorkee, SDSC SAN DIEGO SUPERCOMPUTER CENTER, CDAC, UKRI Science and Technology Facilities Council Scientific Computing, TACC, UNIVERSIDAD DE MÁLAGA, URZ HEIDELBERG UNIVERSITY COMPUTING CENTRE, Northern Illinois University NIU, University of BRISTOL, Stockholm University, KTH, Indian Institute of Science Bangalore, IISER PUNE Indian Institute of Science Education & Research Pune

## CSPs & Frameworks

Google Cloud, Microsoft Azure, Alibaba Cloud, TensorFlow, Taboola, Tencent 腾讯, DataRobot, Baidu 百度, PaddlePaddle 飞桨, NAVER CLOVA, PyTorch

# oneAPI Commercial & Community Support Available
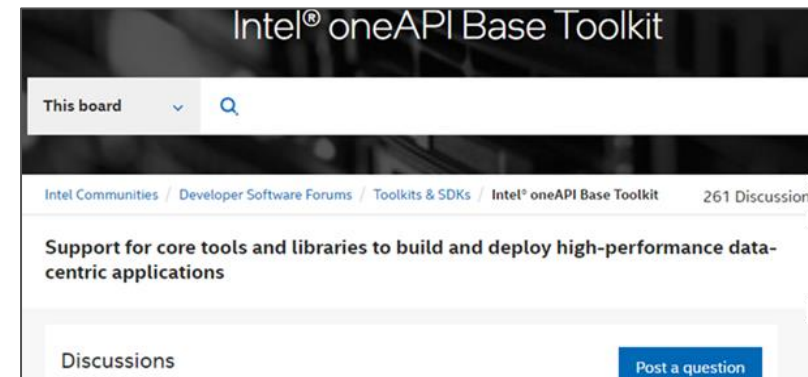
## Priority Support for Intel® oneAPI Toolkits

Every paid version of Intel® oneAPI Developer Toolkits includes Priority Support for that toolkit (Intel oneAPI Base, HPC, IOT, & Rendering Toolkits)

- Direct, private interaction with Intel software support engineers
- Accelerated response time
- Access to—and support for—previous Intel products such as Fortran compiler versions, previous toolkit versions, and more
- Intel Technical Consulting Engineers for on-site or online training and consultation at a reduced cost

**INTEL® PRIORITY SUPPORT**
With All Paid Licenses

## Free Community Support

- Support via the Intel public Community Forum
- Access to only the latest versions of oneAPI Toolkits
- Access to online tutorials and self-help forums

Intel® oneAPI Base Toolkit

This board

Intel Communities / Developer Software Forums / Toolkits & SDKs / Intel® oneAPI Base Toolkit     261 Discussions

**Support for core tools and libraries to build and deploy high-performance data-centric applications**

Discussions     Post a question

# Maximize Your Performance

## With Intel Developer Tools & Hardware Platforms

| HPC & Data Center | AI & Visualization | Embedded Systems & IoT |
|---|---|---|



| **Performance** | **Productivity** | **Freedom** |
|---|---|---|
| Optimize compute performance on the latest Intel CPUs, GPUs and FPGAs | Familiar languages and standards | Open alternative to proprietary lock-in |
| Maximize built-in accelerators | Easily integrate w/ legacy code | Enables easy architecture retargeting |
| Accelerate across AI frameworks | Easily migrate CUDA to SYCL | Code longevity for future hardware |
| | Minimize code re-writes | |

# Details about Intel® oneAPI Toolkits

**Intel® oneAPI Base Toolkit**

**Intel® oneAPI HPC Toolkit**

# Intel® oneAPI Base Toolkit

## Accelerate Data-centric Workloads

A core set of core tools and libraries for developing high-performance applications on Intel® CPUs, GPUs, and FPGAs.

## Who Uses It?

- A broad range of developers across industries
- Add-on toolkit users since this is the base for all toolkits

## Top Features/Benefits

- Data Parallel C++ compiler, library and analysis tools
- SYCLomatic / DPC++ Compatibility tool helps migrate CUDA code to C++ with SYCL
- Python distribution includes accelerated scikit-learn, NumPy, SciPy libraries
- Optimized performance libraries for threading, math, data analytics, deep learning, and video/image/signal processing

### Learn More & Download

---

## Intel® oneAPI Base Toolkit

### Direct Programming

- Intel® oneAPI DPC++/C++ Compiler
- Intel® DPC++ Compatibility Tool
- Intel® Distribution for Python
- Intel® FPGA Add-on for oneAPI Base Toolkit

### API-Based Programming

- Intel® oneAPI DPC++ Library oneDPL
- Intel® oneAPI Math Kernel Library - oneMKL
- Intel® oneAPI Data Analytics Library - oneDAL
- Intel® oneAPI Threading Building Blocks - oneTBB
- Intel® oneAPI Collective Communications Library oneCCL
- Intel® oneAPI Deep Neural Network Library - oneDNN
- Intel® Integrated Performance Primitives - Intel® IPP

### Analysis & debug Tools

- Intel® VTune™ Profiler
- Intel® Advisor
- Intel® Distribution for GDB

intel
1
oneAPI

BASE
TOOLKIT

# Productive and Performant SYCL Compiler

**Intel® oneAPI DPC++/C++ Compiler**

## Uncompromised parallel programming productivity and performance across CPUs and accelerators

- Allows code reuse across hardware targets, while permitting custom tuning for a specific accelerator
- Open, cross-industry alternative to single architecture proprietary language

## Khronos SYCL Standard

- Delivers C++ productivity benefits, using common and familiar C and C++ constructs
- Created by Khronos Group to support data parallelism and heterogeneous programming

## Builds upon Intel's decades of experience in architecture and high-performance compilers

**Learn More & Download**

oneAPI DPC++/C++ Compiler and Runtime

C++ with SYCL Source Code

Clang/LLVM

SYCL Runtime

CPU    GPU    FPGA

# Analysis & Debug Tools

**Get More from Diverse Hardware**

## Design

### Intel® Advisor

- Efficiently offload code to GPUs

- Optimize your CPU/GPU code for memory and compute

- Enable more vector parallelism and improve efficiency

- Add effective threading to unthreaded applications

## Debug

### Intel® Distribution for GDB

- Multiple accelerator support with CPU, GPU and FPGA

- Enables deep, system-wide debug of SYCL, C, C++, OpenMP and Fortran cross-architecture applications

- IDE Integration into Microsoft Visual Studio, VS Code and Eclipse

## Tune

### Intel® VTune™ Profiler

- Tune for GPU, CPU, and FPGA

- Optimize offload performance

- Supports SYCL, C, C++, Fortran, Python, Go, Java or a mix of languages

# Intel® oneAPI Tools for HPC
# Intel® oneAPI HPC Toolkit

**Deliver Fast Applications that Scale**

## What is it?

A toolkit that adds to the Intel® oneAPI Base Toolkit for building high-performance, scalable parallel code on C++, Fortran, SYCL, OpenMP & MPI from enterprise to cloud, and HPC to AI applications.

## Who needs this product?

- OEMs/ISVs
- C++, Fortran, OpenMP, MPI Developers

## Why is this important?

- Accelerate performance on Intel® Xeon® & Core™ processors & Intel accelerators
- Deliver fast, scalable, reliable parallel code with less effort built on industry standards

**Learn More & Download**

## Intel® oneAPI Base & HPC Toolkits

### Direct Programming

- Intel® C++ Compiler Classic
- Intel® Fortran Compiler Classic
- Intel® Fortran Compiler
- Intel® oneAPI DPC++/C++ Compiler?
- Intel® DPC++ Compatibility Tool
- Intel® Distribution for Python
- Intel® FPGA Add-on for oneAPI Base Toolkit

### API-Based Programming

- Intel® MPI Library
- Intel® oneAPI DPC++ Library oneDPL
- Intel® oneAPI Math Kernel Library - oneMKL
- Intel® oneAPI Data Analytics Library - oneDAL
- Intel® oneAPI Threading Building Blocks - oneTBB
- Intel® oneAPI Collective Communications Library oneCCL
- Intel® oneAPI Deep Neural Network Library - oneDNN
- Intel® Integrated Performance Primitives – Intel® IPP

### Analysis & debug Tools

- Intel® Inspector
- Intel® Trace Analyzer & Collector
- Intel® VTune™ Profiler
- Intel® Advisor
- Intel® Distribution for GDB

■ Intel® oneAPI **HPC** Toolkit +
■ Intel® oneAPI **Base** Toolkit

intel 1 oneAPI | HPC TOOLKIT

intel®

# oneAPI Resources
**software.intel.com/oneapi**

## Get Started

- software.intel.com/oneapi
- Documentation + dev guides
- Code Samples
- Intel® Developer Cloud

## Industry Initiative

- oneAPI.io
- oneAPI open Industry Specification
- Open-source Implementations

## Learn

- Training: Webinars & courses
- Intel® DevMesh Innovator Projects
- Summits & Workshops: Live & on-demand virtual workshops, community-led sessions
- Training by certified oneAPI experts worldwide for HPC & AI

## Ecosystem

- Community Forums
- Academic Programs: oneAPI Centers of Excellence: research, enabling code, curriculum, teaching

# SYCL Basics

# SYCL View of Heterogenous Computing Platform



Host (CPU)

Host Memory

Executed on...

SYCL Application

Host code   Device code

submits...

Device

Compute Unit (CU)

Private Memory

Local Memory

Global Memory

Executed on...

Command Group

Command Queue

# Anatomy of a SYCL Application

```cpp
#include <CL/sycl.hpp>
using namespace sycl;

int main() {
std::vector<float> A(1024), B(1024), C(1024);
// some data initialization
    {
      buffer bufA {A}, bufB {B}, bufC {C};
      queue q;
      q.submit([&](handler &h) {
          auto A = bufA.get_access(h, read_only);
          auto B = bufB.get_access(h, read_only);
          auto C = bufC.get_access(h, write_only);
          h.parallel_for(1024, [=](auto i){
              C[i] = A[i] + B[i];
          });
      });
    }
for (int i = 0; i < 1024; i++)
      std::cout << "C[" << i << "] = " << C[i] << std::endl;

}
```

Host code

Accelerator device code

Host code

intel.

# Anatomy of a SYCL Application

```cpp
#include <CL/sycl.hpp>
using namespace sycl;

int main() {
std::vector<float> A(1024), B(1024), C(1024);
// some data initialization
    {
      buffer bufA {A}, bufB {B}, bufC {C};
      queue q;
      q.submit([&](handler &h) {
          auto A = bufA.get_access(h, read_only);
          auto B = bufB.get_access(h, read_only);
          auto C = bufC.get_access(h, write_only);
          h.parallel_for(1024, [=](auto i){
             C[i] = A[i] + B[i];
          });
      });
    }
for (int i = 0; i < 1024; i++)
     std::cout << "C[" << i << "] = " << C[i] << std::endl;

}
```

Application scope

Command group scope

Device scope

Application scope

# SYCL Basics

```cpp
std::vector<float> A(1024), B(1024), C(1024);
    {
        buffer bufA {A}, bufB {B}, bufC {C};
        queue q;
        q.submit([&](handler &h) {
            auto A = bufA.get_access(h, read_only);
            auto B = bufB.get_access(h, read_only);
            auto C = bufC.get_access(h, write_only);
            h.parallel_for(1024, [=](auto i){
                C[i] = A[i] + B[i];
            });
        });
    }
  for (int i = 0; i < 1024; i++)

      std::cout << "C[" << i << "] = " << C[i] << std::endl;

}
```

Buffers creation via host vectors/pointers

Buffers encapsulate data in a SYCL application

- Across both devices and host!

# SYCL Basics

```cpp
std::vector<float> A(1024), B(1024), C(1024);
    {
      buffer bufA {A}, bufB {B}, bufC {C};
      queue q;
      q.submit([&](handler &h) {
          auto A = bufA.get_access(h, read_only);
          auto B = bufB.get_access(h, read_only);
          auto C = bufC.get_access(h, write_only);
          h.parallel_for(1024, [=](auto i){
              C[i] = A[i] + B[i];
          });
      });
    }
  for (int i = 0; i < 1024; i++)

      std::cout << "C[" << i << "] = " << C[i] << std::endl;

}
```

- A queue submits command groups to be executed by the SYCL runtime

- Queue is a mechanism where work is submitted to a device.

# SYCL Basics

```cpp
std::vector<float> A(1024), B(1024), C(1024);
    {
       buffer bufA {A}, bufB {B}, bufC {C};
       queue q;
       q.submit([&](handler &h) {
           auto A = bufA.get_access(h, read_only);
           auto B = bufB.get_access(h, read_only);
           auto C = bufC.get_access(h, write_only);
          h.parallel_for(1024, [=](auto i){
              C[i] = A[i] + B[i];
          });
       });
    }
 for (int i = 0; i < 1024; i++)

       std::cout << "C[" << i << "] = " << C[i] << std::endl;


}
```

- Accessors creation
- Mechanism to access buffer data
- Create data dependencies in the SYCL graph that order kernel executions

# SYCL Basics

```cpp
std::vector<float> A(1024), B(1024), C(1024);
    {
        buffer bufA {A}, bufB {B}, bufC {C};
        queue q;
        q.submit([&](handler &h) {
            auto A = bufA.get_access(h, read_only);
            auto B = bufB.get_access(h, read_only);
            auto C = bufC.get_access(h, write_only);
            h.parallel_for(1024, [=](auto i){
                C[i] = A[i] + B[i]; }
            );
        });
    }
for (int i = 0; i < 1024; i++)

        std::cout << "C[" << i << "] = " << C[i] << std::endl;

}
```

`range<1>{1024}`   `id<1>`

- Vector addition kernel enqueues a parallel_for task.
- Pass a function object/lambda to be executed by each work-item

# Host-side Memory Model – Buffer & Accessor

**Buffers:** Encapsulate data in a SYCL application

- Across both <mark>devices</mark> and host!

**Accessors:** Mechanism to access buffer data

- Create data dependencies in the SYCL graph that order kernel executions

```
int main() {
  auto R = range<1>{ num };
  buffer<int> A{ R }, B{ R };
  queue Q;

  Q.submit([&](handler& h) {
    accessor out(A, h, write_only);

    h.parallel_for(R, [=](auto idx) {
      out[idx] = idx[0]; }); });

  Q.submit([&](handler& h) {
        accessor out(A, h, write_only);
    h.parallel_for(R, [=](auto idx) {
      out[idx] = idx[0]; }); });
...
```

# Host-side Memory Model – Unified Shared Memory



```cpp
using namespace sycl;
int main() {
  queue q;
  float *a = malloc_shared<float>(N, q);
  float *b = malloc_shared<float>(N, q);
  float *c = malloc_shared<float>(N, q);
  q.parallel_for(1024, [=](auto i) {
      c[i] = a[i] + b[i];
  }).wait();
  for (int i=0; i<1024; i++) std::cout << c[i] << "\n";
  free(a, q); free(b, q); free(c, q);
  return 0;
}
```

Host code

Accelerator device code

Host code

# Notices & Disclaimers

**oneAPI Essentials**

# SYCL Program Structure

Learn about SYCL Program Structure, important SYCL Classes and Buffer Memory Model in SYCL

intel.

# Learning Objectives

Explain the SYCL fundamental classes

Use device selection to offload kernel workloads

Decide when to use basic parallel kernels and ND-Range kernels

Understand various ways to synchronize data between host and device with using buffer memory model

Write a complete SYCL program that offload computation to accelerator device

# oneAPIs implementation of SYCL

oneAPIs implementation of SYCL (DPC++)

   =  C++ and SYCL* standard and extensions

Based on modern C++

- C++ productivity benefits and familiar constructs

Standards-based, cross-architecture

- Incorporates the SYCL standard for data parallelism and heterogeneous programming

# Extends SYCL* standard

## Enhance Productivity

- Simple things should be simple to express

- Reduce verbosity and programmer burden

## Enhance Performance

- Give programmers control over program execution

- Enable hardware-specific features


## Fast-moving open collaboration feeding into the SYCL* standard

- Open source implementation with goal of upstream LLVM

- extensions aim to become core SYCL*, or Khronos* extensions

# A Complete SYCL Program

## Single source

- Host code and heterogeneous accelerator kernels can be mixed in same source files

## Familiar C++

- Library constructs add functionality, such as:

| Construct | Purpose |
| --- | --- |
| queue | Work targeting |
| malloc_shared | Data management |
| parallel_for | Parallelism |

Host code

Accelerator device code

Host code

```cpp
#include <CL/sycl.hpp>
constexpr int N=16;
using namespace sycl;
int main() {
  queue q;
  int *data = malloc_shared<int>(N, q);
  q.parallel_for(N, [=](auto i) {
      data[i] = i;
  }).wait();
  for (int i=0; i<N; i++) std::cout << data[i] << "\n";
  free(data, q);
  return 0;
}
```

# SYCL Classes

# Device

- The **device** class represents the capabilities of the accelerators in a oneAPI system.

- The device class contains member functions for querying information about the device, which is useful for DPC++ programs where multiple devices are created.

- The function get_info gives information about the device:

  - Name, vendor, and version of the device

  - The local and global work item IDs

  - Width for built in types, clock frequency, cache width and sizes, online or offline

```cpp
queue q;
device my_device = q.get_device();
std::cout << "Device: " << my_device.get_info<info::device::name>() << std::endl;
```

# Device Selector

- The **device_selector** class enables the runtime selection of a particular device to execute kernels based upon user-provided heuristics.

- The following code sample shows use of the standard device selectors (default_selector, cpu_selector, gpu_selector...) and a derived **device_selector**

```cpp
default_selector selector;
// host_selector selector;
// cpu_selector selector;
// gpu_selector selector;
queue q(selector);
std::cout << "Device: " << q.get_device().get_info<info::device::name>() << std::endl;
```

# Queue

- A queue **submits command groups** to be executed by the SYCL runtime

- Queue is a mechanism where work is submitted to a device.

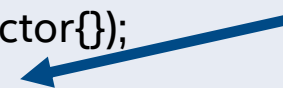- A Queue map to one device and multiple queues can be mapped to the same device.

```
queue q;


q.submit([&](handler& h) {

    // COMMAND GROUP CODE

});
```

# Choosing Where Device Kernels Run

## Work is submitted to queues

- Each queue is associated with exactly one device (e.g. a specific GPU or FPGA)

- You can:

  - Decide which device a queue is associated with (if you want)

  - Have as many queues as desired for dispatching work in heterogeneous systems

| | |
|---|---|
| Create queue targeting any device: | `queue();` |
| Create queue targeting a pre-configured classes of devices: | `queue(cpu_selector{});`<br>`queue(gpu_selector{});`<br>`queue(intel::fpga_selector{});`<br>`queue(accelerator_selector{});`<br>`queue(host_selector{});`    **Always available** |
| Create queue targeting specific device (custom criteria): | `class custom_selector : public device_selector {`<br>  `int operator()(......` **// Any logic you want!**<br>`...`<br>`queue(custom_selector{});` |

# Kernel

- The kernel class encapsulates methods and data for executing code on the device when a command group is instantiated

- Kernel object is not explicitly constructed by the user

-  Kernel object is constructed when a kernel dispatch function, such as parallel_for, is called

```
q.submit([&](handler& h) {
  h.parallel_for(range<1>(N), [=](id<1> i) {
    A[i] = B[i] + C[i]);
  });
});
```

# DPC++ language and runtime

- DPC++ language and runtime consists of a set of C++ classes, templates, and libraries

- Application scope and command group scope :

  - Code that executes on the host

  - The full capabilities of C++ are available at application and command group scope

- Kernel scope:

  - Code that executes on the device.

  - At kernel scope there are limitations in accepted C++

# Parallel Kernels

- Parallel Kernel allows multiple instances of an operation to execute in parallel.

- Useful to offload parallel execution of a basic for-loop in which each iteration is completely independent and in any order.

- Parallel kernels are expressed using the parallel_for function

**for**-loop in CPU application

```
for(int i=0; i < 1024; i++){
    a[i] = b[i] + c[i];
});
```

Offload to accelerator using **parallel_for**

```
h.parallel_for(range<1>(1024), [=](id<1> i){
    A[i] =  B[i] + C[i];
});
```

# Basic Parallel Kernels

The functionality of basic parallel kernels is exposed via range, id and item classes

- range class is used to describe the iteration space of parallel execution

- id class is used to index an individual instance of a kernel in a parallel execution

- item class represents an individual instance of a kernel function, exposes additional functions to query properties of the execution range

```cpp
h.parallel_for(range<1>(1024), [=](id<1> idx){

        // CODE THAT RUNS ON DEVICE

});
```
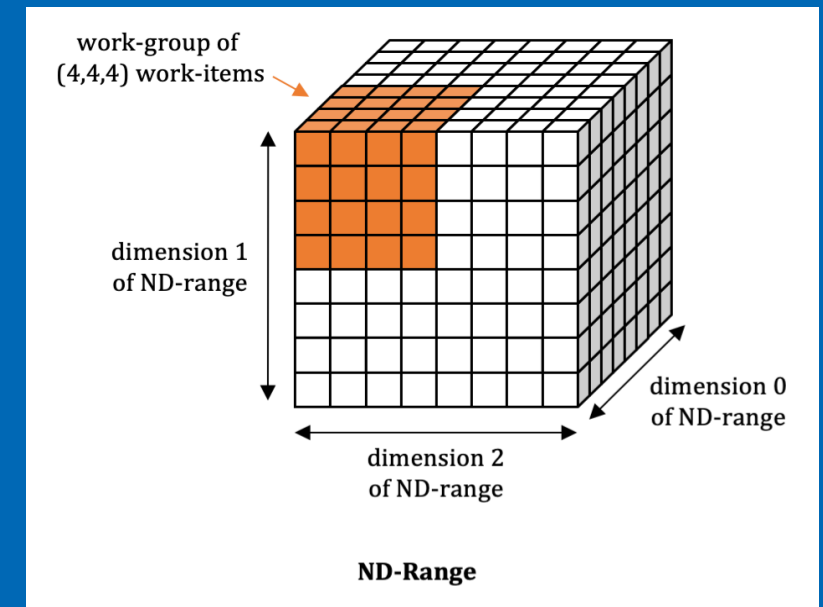
```cpp
h.parallel_for(range<1>(1024), [=](item<1> item){

        auto idx = item.get_id();

        auto R = item.get_range();

        // CODE THAT RUNS ON DEVICE

});
```

# ND-Range Kernels

Basic Parallel Kernels are easy way to parallelize a for-loop but does not allow performance optimization at hardware level.

ND-Range kernel is another way to expresses parallelism which enable low level performance tuning by providing access to local memory and mapping executions to compute units on hardware.

- The entire iteration space is divided into smaller groups called work-groups, work-items within a work-group are scheduled on a single compute unit on hardware.

- The grouping of kernel executions into work-groups will allow control of resource usage and load balance work distribution.



work-group of (4,4,4) work-items

dimension 1 of ND-range

dimension 0 of ND-range

dimension 2 of ND-range

**ND-Range**

# ND-Range Kernels

The functionality of nd_range kernels is exposed via nd_range and nd_item classes

```
h.parallel_for(nd_range<1>(range<1>(1024),range<1>(64)), [=](nd_item<1> item){

    auto idx = item.get_global_id();

    auto local_id = item.get_local_id();

    // CODE THAT RUNS ON DEVICE

});
```

global size          work-group size

- nd_range class represents a grouped execution range using global execution range and the local execution range of each work-group.
- nd_item class represents an individual instance of a kernel function and allows to query for work-group range and index.

# Memory Models

SYCL programs can either use a pointer-based memory model called Unified Shared Memory or can use Buffer-Accessor memory model

- Unified Shared Memory – pointer-based memory model to access data on host and device

- Buffer Memory Model – defines shared array of one, two or three dimensions that can be used by the SYCL kernel and has to be accessed using accessor classes

# Unified Shared Memory

**Unified Shared Memory** enables accessing memory on the host and device with same pointer reference

```
queue q;

auto data =  malloc_shared<int>(N, q);

for(int i=0;i<N;i++) data[i] = 10;

q.parallel_for(N, [=](auto i){

        data[i] += 1;

}).wait();

for(int i=0;i<N;i++) std::cout << data[i] << " ";

free(data, q);
```

Setup Unified Shared Memory →

Host can initialize →

Device can modify →

Host has output →

# Buffer Memory Model

**Buffers:** Encapsulate data in a SYCL application

- Across both devices and host!

**Accessors:** Mechanism to access buffer data

- Create data dependencies in the SYCL graph that order kernel executions

```cpp
queue q;
std::vector<int> v(N, 10);
{
  buffer buf(v);
  q.submit([&](handler& h) {
    accessor a(buf, h , write_only);
    h.parallel_for(N, [=](auto i) { a[i] = i; });
  });
}
for (int i = 0; i < N; i++) std::cout << v[i] << " ";
```

# SYCL Code Anatomy

- SYCL programs require the include of CL/sycl.hpp

- It is recommended to employ the namespace statement to save typing repeated references into the sycl namespace

```
#include <CL/sycl.hpp>

using namespace sycl;
```

# SYCL Code Anatomy

```cpp
void dpcpp_code(int* a, int* b, int* c) {
  // Setting up a DPC++ device queue
  queue q;
  // Setup buffers for input and output vectors
  buffer buf_a(a, range<1>(N));
  buffer buf_b(b, range<1>(N));
  buffer buf_c(c, range<1>(N));
  //Submit Command group function object to the queue
  q.submit([&](handler &h){
    //Create device accessors to buffers allocated in global memory
    accessor A(buf_a, h, read_only);
    accessor B(buf_b, h, read_only);
    accessor C(buf_c, h, write_only);
    //Specify the device kernel body as a lambda function
    h.parallel_for(range<1>(N), [=](auto i){
      C[i] = A[i] + B[i];
    });
  });
}
```

**Step 1:** create a device queue (developer can specify a device type via device selector or use default selector)

**Step 2:** create buffers (represent both host and device memory)

**Step 3:** submit a command for (asynchronous) execution

**Step 4:** create buffer accessors to access buffer data on the device

**Step 5:** send a kernel (lambda) for execution

**Step 6:** write a kernel

Kernel invocations are executed in parallel

Kernel is invoked for each element of the range

Kernel invocation has access to the invocation id

Done!
The results are copied to vector `c` at `buf_c` buffer destruction

# Custom Device Selector

The following code shows derived device_selector that employs a device selector heuristic. The selected device prioritizes a GPU device because the integer rating returned is higher than for CPU or other accelerator.

```cpp
#include <CL/sycl.hpp>
using namespace cl::sycl;
class my_device_selector : public device_selector {
public:
  int operator()(const device& dev) const override {
    int rating = 0;
    if (dev.is_gpu() & (dev.get_info<info::device::name>().find("Intel") != std::string::npos))
      rating = 3;
    else if (dev.is_gpu()) rating = 2;
    else if (dev.is_cpu()) rating = 1;
    return rating;
  };
};
int main() {
  my_device_selector selector;
  queue q(selector);
  std::cout << "Device: " << q.get_device().get_info<info::device::name>() << std::endl;
  return 0;
}
```

# Asynchronous Execution

Think of a SYCL application as two parts:

1. Host code

2. The graph of kernel executions

These execute independently, except at synchronizing operations

- The host code submits work to build the graph (and can do compute work itself)

- The graph of kernel executions and data movements executes asynchronously from host code, managed by the SYCL runtime
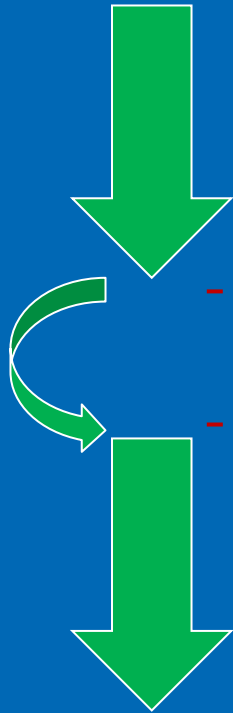
# Asynchronous Execution

**Host code execution**

**Graph executes asynchronously to host program**
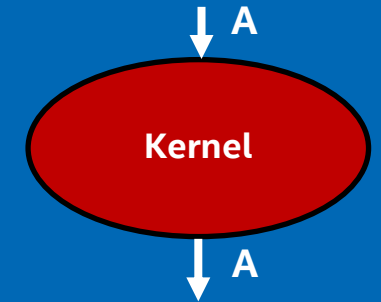
**Enqueues kernel to graph, and keeps going**

```cpp
#include <CL/sycl.hpp>
constexpr int N=16;
using namespace sycl;
int main() {
  std::vector<int> data(N);
  {
    buffer A(data);
    queue q;
    q.submit([&](handler& h) {
      accessor out(A, h, write_only);
      h.parallel_for(N, [=](auto i) {
        out[i] = i;
      });
    });
  }
  for (int i=0; i<N; ++i) std::cout << data[i];
}
```

**A**

**Kernel**

**A**

# Implicit dependency between kernels

```cpp
int main() {
  auto R = range<1>{ num };
  buffer<int> A{ R }, B{ R };
  queue q;

  q.submit([&](handler& h) {
    accessor out(A, h, write_only);
    h.parallel_for(R, [=](id<1> i) {
      out[i] = i; }); });              }  Kernel 1

  q.submit([&](handler& h) {
    accessor out(A, h, write_only);
    h.parallel_for(R, [=](id<1> i) {
      out[i] = i; }); });              }  Kernel 2

  q.submit([&](handler& h) {
    accessor out(B, h, write_only);
    h.parallel_for(R, [=](id<1> i) {
      out[i] = i; }); });              }  Kernel 3

  q.submit([&](handler& h) {
    accessor in(A, h, read_only);
    accessor inout(B, h);
    h.parallel_for(R, [=](id<1> i) {
      inout[i] *= in[i]; }); });       }  Kernel 4
}
```
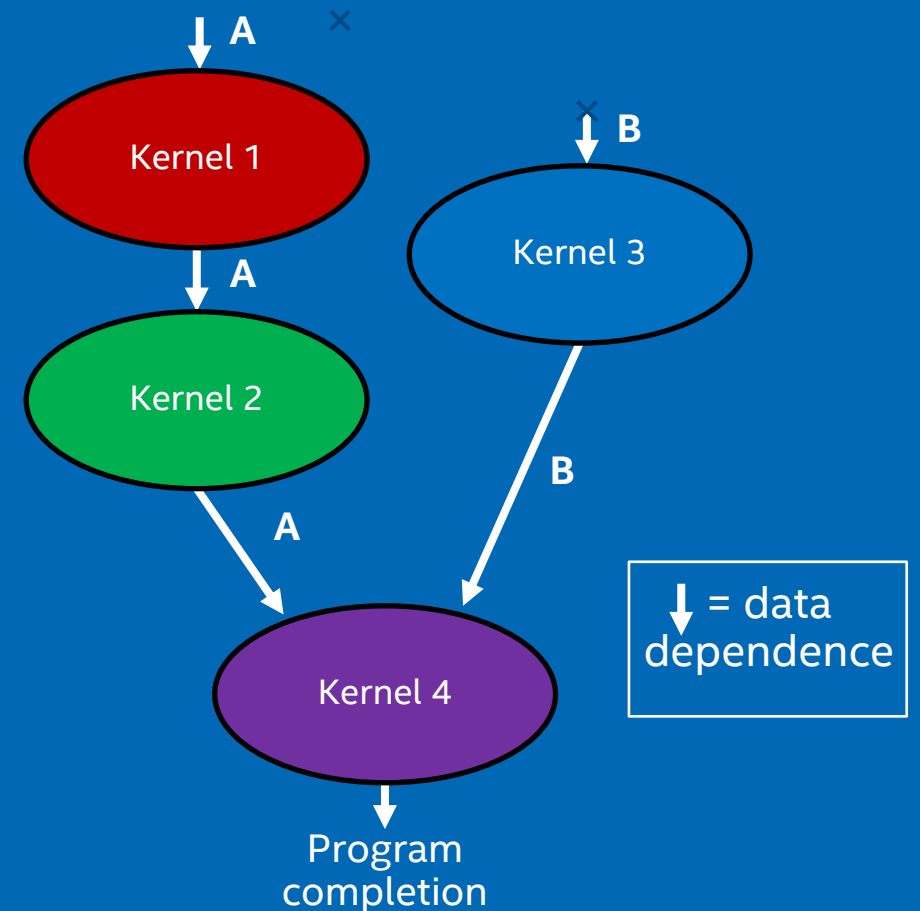
Automatic data and control dependence resolution!



⬇ = data dependence

# Host Accessors

- An accessor which uses host buffer access target

- Created outside of command group scope

- The data that this gives access to will be available on the host

- Used to <span style="color:gold">synchronize the data back to the host</span> by constructing the host accessor objects

intel.

```cpp
#include <CL/sycl.hpp>
using namespace sycl;
constexpr int N = 16;

int main() {
  std::vector<double> v(N, 10);
  queue q;

  buffer buf(v);
  q.submit([&](handler& h) {
    accessor a(buf, h)
    h.parallel_for(N, [=](auto i) {
      a[i] -= 2;
    });
  });

  host_accessor b(buf, read_only);
  for (int i = 0; i < N; i++)
    std::cout << b[i] << "\n";
  return 0;
}
```

Buffer takes ownership of the data stored in vector.

Creating host accessor is a blocking call and will only return after all enqueued kernels that modify the same buffer in any queue completes execution and the data is available to the host via this host accessor.

# Synchronization – Buffer Destruction

```cpp
#include <CL/sycl.hpp>
using namespace sycl;
constexpr int N=16;

void dpcpp_code(std::vector<double> &v, queue &q){
  buffer buf(v);
  q.submit([&](handler& h) {
    accessor a(buf, h);
    h.parallel_for(N, [=](auto i) {
      a[i] -= 2;
    });
  });
}

int main() {
  std::vector<double> v(N, 10);
  queue q;
  dpcpp_code(v,q);
  for (int i = 0; i < N; i++)
      std::cout << v[i] << "\n";
  return 0;
}
```

Buffer creation happens within a

separate function scope.


When execution advances

beyond this function scope,

buffer destructor is invoked

which relinquishes the ownership

of data and copies back the data

to the host memory.