

2023 SIAM Conference on Computational Science and Engineering Intel® Developer Tools for Serious SYCL

Nisha Patel (Director of Developer Tools Customer Engineering EMEA, WW customer engagements) February 2023



SYCL 2020: Khronos standard ready to support accelerators

Based on modern C++ / Decoupled from OpenCL

Unified address space

sycl::malloc_host
sycl::malloc_device
sycl::malloc_shared

Atomic ops on non-atomic objects

sycl::atomic_ref

Subgroups (similar to CUDA warp) Subgroup and work group algorithms

Any_of, all_of, non_of, reduce, exclusive_scan, inclusive_scan, shift_left, shift_right, select, permute

Reductions

sycl::reduction

SYCL 1.2.1

float *sum = (float *)malloc(sizeof(float)); float *data = (float *)malloc(N*sizeof(float));

const size_t tc = N/ITEMS_PER_THREAD; const size_t itb = ITEMS_PER_THREAD*BS; float *output = (float *)malloc(tc*sizeof(float)); buffer<float, 1> buf(data, N); buffer<float, 1> ans(output, tc); q.submit([&](handler &h) { auto buf_acc = buf.get_access<access::mode::read>(h); auto ans_acc = ans.get_access<access::mode::discard_write>(h); h.parallel_for<class SumK>(nd_range<1>(thread_count, BS), [=](nd_item<1> it) { const size_t st = it.get_group(0)*itb+it.get_local_id(0); float 1sum = 0;for (size_t i = st; i <st+itb; i += it.get_local_range(0)) {</pre> lsum += buf_acc[i]; ans_acc[it.get_global_id(0)] = lsum; }); }).wait(); auto ans_acc = ans.get_access<access::mode::read>(); sum = 0:for (size_t i = 0; i < tc; ++i) { *sum += ans_acc[i]; }</pre>

SYCL 2020

float* sum = malloc_shared<float>(1, q);
float* data = malloc_shared<float>(N, q);

q.parallel_for(N, reduction(sum, std::plus<>()),
 [=](size_t i, auto& sum) {
 sum += data[i];
});



oneAPI

© codeplay°

ComputeCpp^{*}



Accelerating Choice with SYCL **Khronos Group Standard**

- Open, standards-based
- Multiarchitecture performance
- Ereedom 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



Testing Date: Performance results are based on testing by Intel as of August 15, 2022 and may not reflect all publicly available updates

Configuration Details and Workload Setup: Intel® Xeon® Platinum 8360Y CPU @ 2.4GHz, 2 socket, Hyper Thread On, Turbo On, 256GB Hynix DDR4-3200, ucode 0x000363. GPU: Nvidia A100 PCIe 80GB GPU memory. Software: SYCL open source/CLANG 15.0.0, CUDA SDK 11.7 with NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, Ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, Ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, Ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, Ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, Ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, Ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, Ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, Ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64-nvidia-cuda, NVDIA-NVCC 11.7.64, cuMath 11.7, cuDNN 11.7, ubuntu 22.04.1. SYCL open source/CLANG compiler switches: -Iscycl-targets-nvptx64, cuMath 11.7, cuDNN 11.7, ubuntu 22.04.1. SYCL open source/ arch=compute_80, code=sm_80, Represented workloads with Intel optimizations

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See configuration disclosure for details. No product or component can be absolutely secure.

Performance varies by use, configuration, and other factors. Learn more at www.intel.com/PerformanceIndex. Your costs and results may vary.

Architectures

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

Offload vs. Heterogeneous Computing

Heterogeneous computing – not just offload, but leveraging CPUs & accelerators in the same application



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

Heterogeneous Computing: Run sub-problems in parallel on the hardware best suited to them.

intel

4

What is oneAPI?



oneAPI enables cross-platform XPU programming

oneAPI open specification provides a standard programming language, libraries, and hardware abstraction layer



intel

6

Programming challenges for multiple architectures

Application	Workloads N	Need Diverse	Hardware		
Middleware & Frameworks					
CPU programming model	GPU programming model	FPGA programming model	Other accel. programming models		
СРО	GPU	FPGA	Other accel.		
	XP	Us			

oneAPI Libraries

Domain	Name	Description	Open Spec	Open Source
Darallel Drogramming	oneDPL	Data Parallel C++ Library including Parallel STL	Yes	Yes
Parallel Programming	oneTBB	Threading Building Blocks	Yes	Yes
	oneDNN	Deep Neural Networks	Yes	Yes
AI & ML	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

Architectures supported



Intel[®] oneAPI Toolkits



Intel [®] oneAPI Base Toolkit	intel BASE TOOLKIT	A core set of high-performanc building C++, SYCL, C/OpenM	rmance libraries and tools for)penMP, and Python applications			
Add-on Domain-specific Toolkits	Intel® oneAPI Tools for HPC Deliver fast Fortran, OpenMP & MPI applications that scale	intel Image: For Edge & developers Intel® oneAPI Tools for IoT Build efficient, reliable solutions that run at network's edge	IOT Intel RENDERING ONCAPI RENDERING Intel® ONEAPI Rendering Toolkit Create performant, high-fidelity visualization applications			
Toolkits powered by oneAPI	intel. AI ANALYTICS TOOLKIT FOR AI develor Intel [®] AI Analy Accelerate mad science pipelin optimized DL & performing Pyt	opers and data scientists Image: Comparison of the scientists ytics Toolkit Intelearning & data chine learning & data Intelearning & data es end-to-end with Depers and begin of the scientists ML frameworks & high- edge hon libraries Image: Comparison of the scientists	PenVINO™ For deep learning inference developers I® OpenVINO™ toolkit Ioy high performance inference & applications from e to cloud			
	Download at <mark>intel.com/oneAPI</mark> Or run tools in the <u>Intel® Developer Cloud</u>					

Intel[®] oneAPI and AI Tools 2023 – Highlights Optimized, Standards-based Support for Powerful New Architectures

Optimized support for Intel's upcoming portfolio of CPU and GPU architectures	Compilers & SYCL Support	Performance Libraries
intel4th Gen Intel® Xeon® Scalable Processors with Intel® Advanced Matrix Extensions, Quick assist Technology, Intel® AVX-512, bfloat16, and more built-in acceleratorsintelMatrix Extensions, Quick assist Technology, Intel® AVX-512, bfloat16, and more built-in acceleratorsintelIntel® Xeon® Max Series CPUs with high- bandwidth memoryintelIntel® Data Center GPUs, including Flex Series with hardware AV1 encode and Max Series with datatype flexibility, Intel® Xe Matrix Extensions, vector engine, XE-Link, and other features	 Intel® oneAPI DPC++/C++ Compiler improves CPU and GPU offload performance and broadens SYCL language support for improved code portability and productivity Intel® oneAPI DPC++ Library (oneDPL) expands support of the C++ standard library in SYCL kernels with additional heap and sorting algorithms and adds the ability to use OpenMP for thread-level parallelism. Intel® DPC++ Compatibility Tool (based on the open source SYCLomatic project) improves migration of CUDA library APIs, including those for runtime and drivers, cuBLAS, and cuDNN. Intel® Fortran Compiler implements coarrays, eliminating the need for external APIs such as MPI or OpenMP, expands OpenMP 5.x offloading features, adds DO CONCURRENT GPU offload, and improves optimizations for source-level debugging. 	Intel® oneAPI Math Kernel Library increases CUDA library function API compatibility coverage for BLAS and FFT; for 4 th Gen Xeon, leverages Intel® XMX to optimize matrix multiply computations for TF32, FP16, BF16, and INT8 data types; and provides interfaces for SYCL and C/Fortran OpenMP offload programming. Intel® oneAPI Threading Building Blocks improves support and use of the latest C++ standard for parallel_sort, offers an improved synchronization mechanism to reduce contention when multiple task_arena calls are used concurrently, and adds support for Microsoft Visual Studio 2022 and Windows Server 2022. Intel® oneAPI Video Processing Library supports the industry's only hardware AVI codec in the Intel Data Center GPU Flex Series and Intel® Arc™ processors; expands OS support for RHEL9, CentOS, Stream 9, SLESI5Sp4, and Rocky 9 Linux; and adds parallel encoding feature to sample_multi_transcode.

Intel[®] oneAPI and AI Tools 2023 – Highlights Optimized, Standards-based Support for Powerful New Architectures

Analysis & Debug

Intel[®] VTune[™] Profiler enables ability to identify MPI imbalance issues via its Application Performance Snapshot feature; delivers visibility into Xe Link cross-card traffic for utilization, bandwidth consumption, and other issues; and adds support for 4th Gen Intel[®] Xeon[®] Scalable Processors, Intel[®] Data Center GPU Max Series, and 13th Gen Intel[®] Core[™] processors.

Intel® Advisor adds automated roofline analysis for Data Center GPU MAX Series to identify and prioritize memory, cache, or compute bottlenecks and understand their causes, and delivers actionable recommendations for optimizing data-transfer reuse costs of CPU-to-GPU offloading.

Al & Analytics

Intel® AI Analytics Toolkit can now be run natively on Windows with full parity to Linux except for distributed training (GPU support is coming in Q1 2023).

Intel® oneAPI Deep Neural Network Library further supports delivery of superior CNN performance by enabling advanced features in 4th Gen Xeon CPUs including Intel AMX, AVX-512, VNNI, and bfloat16.

Intel[®] Distribution of Modin integrates with new heterogeneous data kernels (HDK) solution in the back end, enabling AI solution scale from low-compute resources to large- or distributed-computed resources.

Beta additions for <u>Intel® Distribution for Python</u> include compute-follows-data model extension to GPU, data exchange between libraries and frameworks, and dataparallel extensions for NumPy and Numba packages.

Rendering & Visual Computing

Intel® oneAPI Rendering Toolkit includes the Intel® Implicit SPMD Program Compiler runtime library for fast SIMD performance on CPUs.

Intel® Open Volume Kernel Library increases memorylayout efficiency for VDB volumes and adds an AVX-512 8-wide CPU device mode for increased workload performance.

Intel® OSPRay and Intel® OSPRay Studio add features for multi-segment deformation motion blur for mesh geometry, primitive, and objects; face-varying attributes for mesh and subdivision geometry; new light capabilities such as photometric light types; and instance ID buffers to create segmentation images for AI training.

oneAPI: A Bridge to Our Heterogeneous/Distributed Future

My vision for how we bring oneAPI into a future dominated by power-optimized heterogenous chips organized into distributed systems:



Conclusion

- The present and future of computing is heterogeneous
- oneAPI proposes an open specification for programming heterogenous systems
- Results demonstrate the promise of multi-architecture, multi-vendor oneAPI
- Future opportunities : continued ease of development, multi-architecture, distributed, disaggregated

Call to action: Interested? Join the community at oneapi.io

Developer Tools for Intel[®] Data Center GPU Flex Series Multiarchitecture Acceleration for Media, Visual Inferencing & Cloud Workloads

Intel[®] oneAPI Video Processing Library for GPU, and GStreamer + FFmpeg for CPU speed media processing, delivery and cloud gaming streaming.

Intel® oneAPI DPC++/C++ Compiler & oneVPL, both part of the Intel® oneAPI Base Toolkit, scale AI visual inference workloads up to 150 TOPS INT8 between CPU and GPU with the same code.

Intel® VTune™ Profiler analyzes application performance and accelerates compute-intensive tasks by identifying the most time-consuming parts of GPU code and optimizing GPU offload schema and data transfers for SYCL, OpenCL code, Microsoft DirectX*, or OpenMP* offload code.



Intel[®] Graphics Performance Analyzers

optimize visual compute and graphics workloads on CPUs and GPUs .

Intel[®] OpenVINO[™] toolkit, powered by oneAPI, supports popular deep learning frameworks such as TensorFlow and PyTorch to streamline AI visual inferencing and deploy quickly.

Intel[®] Al Analytics Toolkit, powered by oneAPI, accelerates machine learning & data science pipelines end-to-end with optimized DL & ML frameworks & high-performing Python libraries

Developer Tools for Intel[®] Data Center GPU Max Series Build Multiarchitecture Applications with Breakthrough Performance for HPC & AI

Accelerate HPC

oneAPI DPC++/C++ and Fortran compilers and performance libraries (oneMKL, oneDNN, Intel® MPI Library) activate X^e Matrix Extensions (Intel[®] XMX) for acceleration and Intel® Xe Link for direct GPU-to-GPU communications to speed up applications running on multiple GPUs.

Analysis and debug tools help developers analyze and accelerate compute-intensive tasks by identifying the most time-consuming parts of GPU code, ensure correctness and efficiently debug applications on discrete GPUs. Get actionable advice to design code that runs optimally on Intel GPUs.

Boost AI & Deep Learning Inference

oneDNN in the Intel® oneAPI Base Toolkit utilizes Intel XMX to accelerate AI training/inference in deep learning frameworks.

Intel® OpenVINO™ toolkit powered by oneAPI. streamlines AI visual inferencing and accelerates deployment.

Both tools support popular deep learning frameworks such as TensorFlow and PyTorch. GPU support for tools in the Intel[®] AI Analytics Toolkit is coming soon.



Create Multiarchitecture Code Efficiently with Code Migration Tools

Migrate CUDA code to SYCL to create a single source code base for easy portability across multiple vendors' architectures - including Intel Max Series GPU.

The Intel® DPC++ Compatibility Tool, based on open source SYCLomatic, automates most of the process. Save significant time on ongoing code maintenance.

Coming Soon – High-Performance, High-Fidelity Ray Tracing & Rendering

Intel[®] Embree, fully ported to C++ with SYCL will take advantage of hardware ray tracing units and expanded memory capacity to deliver increased rendering performance across larger datasets.

Intel[®] Open Image Denoise will leverage Intel XMX delivering high fidelity, high-performance denoising capabilities and Al optimization.

oneAPI for NVIDIA and AMD GPUs

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

Develop code using SYCL and run on NVIDIA and AMD GPUs.



Download free plugins from developer.codeplay.com



Free and Open Source



Paid Priority Support



- Download free plugin from Codeplay website
- Code is open source
- Compatibility

CUDA SDK >= 11.7

GPUs with at least sm_50

- Raise support issues direct to Codeplay engineers
- Escalate defect reports







oneAPI Industry Momentum



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.

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 core¹



Developer Challenges: Multiple Architectures, Vendors, and Programming Models



Open, Standards-based, Multiarchitecture Programming

oneAPI Industry Initiative Break the Chains of Proprietary Lock-in

Freedom to Make Your Best Choice

- One 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
- Compatible with existing languages and programming models including C, C++ with SYCL, Python, OpenMP, Fortran, 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



faafa.

oneAPI

CUDA to SYCL Migration Made Easy Open Source SYCLomatic Tool Reduces Code Migration Time



Assists developers migrating code written in CUDA to C++ with SYCL, generating **human readable** code wherever possible

~90-95% of code typically migrates automatically¹

Inline comments are provided to help developers finish porting the application

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

Codeplay Compiler Plug-ins for Nvidia and AMD GPUs Adding support for NVIDIA and AMD GPUs to the Intel® oneAPI Base Toolkit

oneAPI for NVIDIA & AMD GPUs

- Free Codeplay download of latest binary plugins to the Intel DPC++/C++ compiler:
 - Nvidia GPU
 - AMD Beta GPU
- Availability at the same time as the Intel oneAPI Base Toolkit
- Plug-ins updated quarterly in-sync with oneAPI

Priority Support

- Sold by Intel and Codeplay and our channel
- Requires Intel Priority support for Intel DPC++ /C++ compiler
- Intel takes first call and Codeplay delivers backend support
- Codeplay access to older versions of plugins

Nvidia GPU plug-in

AMD GPU plug-in

Codeplay blog

Codeplay press release

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







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



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



Maximize Your Performance With Intel Developer Tools & Hardware Platforms

HPC & Data Center

AI & Visualization

Embedded Systems & IoT

intel. intel intel. intel intel. AI ANALYTICS TOOLKIT oneAPI oneAPI oneAPI oneAPI RENDERING TOOLKI HPC TOOLKIT INT TOOLKIT BASE TOOLKIT intel intel intel intel intel intel intel intel. intel intel intel. intel GPU GPU ATOM CORE Xeon AGILEX Xeon Xeon AGILEX CYCLONE MAX SERIES MAX SERIES FLEX SERIES FLEX SERIES **Productivity** Performance Freedom Familiar languages and standards Open alternative to proprietary lock-in Optimize compute performance on the latest Intel CPUs, GPUs and FPGAs Easily integrate w/legacy code Enables easy architecture retargeting Maximize built-in accelerators Easily migrate CUDA to SYCL Code longevity for future hardware Accelerate across Al frameworks Minimize code re-writes

Details about Intel[®] oneAPI Toolkits Intel[®] oneAPI Base 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

	t	
Direct Programming	API-Based Programming	Analysis & debug Tools
Intel® oneAPI DPC++/C++ Compiler	Intel [®] oneAPI DPC++ Library oneDPL	Intel® VTune [™] Profiler
Intel [®] DPC++ Compatibility Tool	Intel [®] oneAPI Math Kernel Library - oneMKL	Intel [®] Advisor
Intel [®] Distribution for Python	Intel® oneAPI Data Analytics Library - oneDAL	Intel [®] Distribution for GDB
Intel® FPGA Add-on for oneAPI Base Toolkit	Intel® oneAPI Threading Building Blocks - oneTBB	
	Intel® oneAPI Video Processing Library - oneVPL	
	Intel [®] oneAPI Collective Communications Library oneCCL	
	Intel® oneAPI Deep Neural Network Library - oneDNN	intel
	Intel® Integrated Performance Primitives - Intel® IPP	ONEAPI

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



Intel® DPC++ Compatibility Tool Minimizes Code Migration Time

Assists developers migrating code written in CUDA to C++ with SYCL once, generating human readable code wherever possible

~90-95% of code typically migrates automatically¹

Inline comments are provided to help developers finish porting the application

SYCLomatic tool is the open source version

Intel DPC ++ Compatibility Tool Usage Flow



Free Yourself from Vendor Lock-in – SYCL* Adoption Made Easy Learn More & Download

¹Intel 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.

Analysis & Debug Tools Get More from Diverse Hardware

	Debug	E Tune
Intel [®] Advisor	Intel [®] Distribution for GDB	Intel [®] VTune [™] Profiler
 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 	 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 for GPU, CPU, and FPGA Optimize offload performance Supports SYCL, C, C++, Fortran, Python, Go, Java or a mix of languages

Powerful Performance Libraries

Rich Functionality

Optimized performance libraries for every use case:

Threading, offload, math, data analytics, data processing, rendering, ray tracing, DNN, comms, crypto, and more.

Realize all the Hardware Value

Designed for acceleration of key domain-specific functions

Freedom of Choice

Pre-optimized for each target platform for maximum performance

Intel® oneAPI Video Processing Library oneVPL Intel® oneAPI Threading Building Blocks oneTBB Intel® oneAPI DPC++ Library oneDPL	ntel® oneAPI Data Analytics Library oneDAL Intel® oneAPI Collective Communications Library oneCCL				
Intel® oneAPI Threading Building Blocks oneTBB Intel® oneAPI DPC++ Library oneDPL	Intel® oneAPI Collective Communications Library oneCCL				
Intel® oneAPI DPC++ Library oneDPL					
	Intel® oneAPI Rendering Toolkit Rendering and Ray-Tracing Libraries				
Intel® Integrated Performance Primitives Intel® IPP	Intel [®] MPI Library				
The One-Stop Shop for All your					

Intel[®] oneAPI DPC++ Library (oneDPL) Accelerate SYCL C++ Kernels on Intel CPUs, GPUs & FPGAs

Optimized C++ Standard Algorithms

Contains 75 parallelized C++17 algorithms and utilities for efficient application development and deployment on a variety of hardware.

Based on parallel libraries that C++ developers are already familiar with

Incorporates popular libraries Parallel STL and Boost. Compute for easier developer adoption.

Integrated with Intel® DPC++ Compatibility Tool

Complements all oneAPI DPC++ components to simplify migration of developers' CUDA* code to DPC++ code.

Bringing Multi-Architecture Compute to C++ Learn More & Download

Intel® oneAPI Deep Neural Network Library (oneDNN) Deliver High Performance Deep Learning

Helps developers create high performance deep learning frameworks

Abstracts out instruction set & other complexities of performance optimizations

Same API for both Intel CPUs and GPUs, use the best technology for the job

Supports Linux, Windows

Open sourced for community contributions



Intel® oneAPI Video Processing Library (oneVPL) Accelerated Video Processing with a Unified Programming API

Deliver fast, high-quality video transcoding from camera to cloud

Boost media and video application performance with hardware-accelerated codecs and programmable graphics on Intel CPUs and GPUs

Simple API that works the same on CPU and GPU

Using the API, developers have full control over codec visual quality and performance



Intel[®] oneAPI Collective Communications Library (oneCCL) Optimize Communication Patterns

Provides optimized communication patterns for high performance on Intel CPUs & GPUs to distribute model training across multiple nodes

Transparently supports many interconnects, such as Intel® Omni-Path Architecture, InfiniBand, & Ethernet

Built on top of lower-level communication middleware-MPI & libfabrics

Enables efficient implementations of collectives used for deep learning training-all-gather, all-reduce, & reduce-scatter



Intel[®] VTune[™] Profiler Tune for CPU, GPU & FPGA

Analyze SYCL code

See the lines of SYCL that consume the most time

Tune for Intel CPUs, GPUs & FPGAs

Optimize for any supported hardware accelerator

Optimize Offload

Tune OpenMP offload performance

Wide Range of Performance Profiles

CPU, GPU, FPGA, threading, memory, cache, storage... Flame graph display improves visualization of hot spots

Supports Popular Languages

SYCL, C, C++, Fortran, Python, Go, Java, or a mix

So	ource Assembly 💵 = 😽 👉 🛶	٩
🔺	Source	 ♦ GPU Instructions Executed by Instruction T[≫] ♥ Control Flow ♥ Send & Wait ♥ Int32 & SP Float ♥ Int64 & DP Float ♥ Other
158	dx = ptr[j].pos[0] - ptr[i].pos[0];	75,002,500
159	dy = ptr[j].pos[1] - ptr[i].pos[1]	12,500,000 🦲
160	<pre>dz = ptr[j].pos[2] - ptr[i].pos[2];</pre>	12,500,000 📒
161		
162	distanceSqr = dx*dx + dy*dy + dz*dz	87,500,000
163	distanceInv = 1.0 / sqrt(distanceSo	12,500,000 📒
164		
165	ptr[i].acc[0] += dx * G * ptr[j].ma	162,503,750
166	ptr[i].acc[1] += dy * G * ptr[j].ma	150,000,000
167	ptr[i].acc[2] += dz * G * ptr[j].ma	150,000,000



Images above show analysis of SYCL code and GPU Offload profiling.

Learn More & Download

Intel® Advisor Configure Your Accelerated Computing Solution

Offload Advisor

Estimate performance of offloading to an accelerator

Roofline Analysis

Optimize CPU/GPU code for memory and compute

Vectorization Advisor

Add and optimize vectorization

Threading Advisor

Add effective threading to unthreaded applications

Flow Graph Analyzer

Create and analyze efficient flow graphs





Know Before You Offload

Learn More & Download

Intel[®] Distribution for GDB* Heterogeneous Application Debug

High-level language debug support

Multiple accelerator support: Intel CPU, GPU, FPGA emulation

Auto-detect accelerator architecture during application runtime

Non-proprietary open-source solution based on GDB

eclipse-workspace - Sepia_Filter/src/sepia_dpcpp.cpp - Eclipse IDE 🛛 🗇 \ominus						
File Edit Source Refactor Navigate Search Project Run Intel Window Help						
🔨 体 🔳 体 Debug 🗸 😧 Sepia_Filter Debug 🔨 🖗 🖂 マ 🖓 🖓 😒 🖉 🖉 🖉 🖉 🖉 🖉 🖉 🖉 🖉 🖉 🖉 🖉	- ¢¢ - ¢¢	Quick Access				
🛊 Debug 🕮 🕒 Project Explorer 🛛 🙀 💆 🗢 🗆 🕞 sepia_dpcpp.cpp 🕮 🗠	🕬 Variab 😫 💁	Break 🐔 Expre 🛋 Modul 🔗 🗆				
[Sepia_Filter Debug [C/C++ Application] [49 // able to find the full call graph automatically.		200 C C C C C C C C C C C C C C C C C C				
▼ ② sepia [12598] [cores: 2,6,7] 50 // always_initine as calls are expensive on Gen GPU.	Name	Type Value				
▼ Phread #1 [sepia] 12598 [core: 2] (Suspended: Breakpo 52 // - coeffs can be declared outside of the function, but still must be constant	• • src image	float * 0x7ffff0ba3010				
<pre>sepia_impl() at sepia_dpcpp.cpp:70 0x4254ec</pre> 54 // - SYGL compiler will automatically deduce the address space for the two 54 // - printers: operalization for particular address space	▶ ● dst image	float * 0x7fffecba2010				
main() at sepia_dpcpp.cpp:224 0x4254ec 55 // can used for more control	00+j	int 0				
▶ @ Thread #2 [sepia] 12607 [core: 6] (Suspended : Containe 569attribute_((always inline))	00-k	int 0				
Dread #3 [sepia] 12634 [core: 7] (Suspended : Contain 50 static voia sepia_imp((toat *src_image, toat *ost_image, int i) { sepialize the second sepialize the second se	00-W	float 0				
j opt/intel/inteloneapi/debugger/latest/gdb/intel64/bin/t 59 { 0.2f, 0.3f, 0.3f, 0.0f,	00-j	int 0				
60 0.1f, 0.5f, 0.5f, 0.0f,	▼ Coeffs	const fli 0x7ffffffb8b0				
62 0.0f, 0.0f, 0.0f, 0.0f };	⊷coeffs[0]	const fl 0.20000003				
	⊷coeffs[1]	const fl 0.300000012				
Enter location here 2008 10 0 4 1 = CHANNELS PER PIXEL;	⊷coeffs[2]	const fl= 0.300000012				
0 0000000004254ec: mov -0x06(%rbp),%eax = 66 for (int j = 0; j < 4; ++j) {	↔coeffs[3]	const fle 0				
000000000004254f2: add -0x6(%rbp),%eax 67 float w = 0.0f;	⊷coeffs[4]	const fl+ 0.100000001				
00000000004254f5: movslq %eax,%rcx 60 for (int k = 0; k < 4; ++k) (⊷coeffs[5]	const fli 0.5				
0000000000425476: mov = 0x80(%rb).%rcx = 7, xxm 000000000425476: mov = 0x8(%rb).%rcx = 70 w += coeffs[4 * j + k] * src image[i + k];	49-coeffs[6]	const fli 0.5				
000000000425502: mov -0x14(%rbp),%eax 71 }	Miccoeffe[7]	const fl. A				
000000000425595; add -0x6c(%rbp),%eax 72 Ust_Amage[1 +]] - w,	Details:{0	.200000003. 0.300000012. 0.30000				
000000000425306: mUlss (%rcx,%rdx,4),%xmm0 74 }	Default:0x	711111116860				
000000000425510: addss -0x68(%rbp),%xmm0 75 76 // Few useful acronyms- 'using namespace cl-+svcl+' also helps	Decimal:14	0737488337072				
0000000000425515: movs5 %xmm0, 0x08(%rDp) 69 69 69 69 69 69 60 60 60 60 60 60 60 60 60 60	Binary:111	111111111111111111111111111111111111111				
00000000042551a: mov -0x6c(%rbp),%eax 78 constexpr auto sycl write = cl::sycl::access::mode::write:	Octal:0377	777777734260				
00000000042551d: add \$0x1,%eax // constexpr auto sycl_global_burrer = cl::sycl::access::larget::global_burrer;						
00000000004255220: mov veax, vxxc(srop) 0000000000425523: impo 0x425e2 emain(int, ch. 81⊖// This is alternative (to a lambda) representation of a SYCL kernel.						
72 dst_image[i + j] = w; Console IIII Registers R Problems O Executables I Memory Debugger Console 22 DPC+	Compatibility Tool	📕 🖳 🖾 🖛 📼 🗋				
000000000425528: movs5 -0x66(%rbp),%xm0 openengaaaasts24 Sepia Filter Debug [C/C++ Application] /opt/intel/inteloneapi/debugger/latest/gdb/intel64/bin/gdb-one	api (8.3)					
00000000425331 mov -0x10(stop),stax 00000000425531 mov -0x14(stop),stax						
000000000425534: add -0x64(%rbp),%ecx Thread 1 "sepia" hit Breakpoint 2, sepia impl (src image=0x7ftTtBba3010, dst image=0x7ftTbba3010, dst i						
uuuuuuuuuuuusassii movsti teexiina dicaa tadi ad						
66 for (int j = 0; j < 4; ++j) { 70 w += coeffs[4 * j + k] * src_image[i + k];						
00000000042533f: mov -0x64(%rbp),%ecx [gob]						
	ên :					

oneAPI FPGA Add-On SYCL Coding for Spatial Architecture

For Experienced FPGA Developers

Ease of Use

The same heterogeneous SYCL code. Experienced FPGA users can take advantage of a streamlined programming model using Data Parallel C++

Runtime Analysis Support

Collect profiling data at runtime to analyze CPU and FPGA interaction with Intel® VTune™ Profiler

Real-time Processing

Process data faster with deterministic low latency, lower power and high throughput

Device Specific Optimizations

One-day class provides experienced FPGA developers training to begin optimizing oneAPI code for FPGA



Where to get started? Writing a new program **Optimizing your** Adding accelerators to 3 2 or porting existing source code? vour hardware config? CUDA code? Intel[®] DPC++ Compatibility Tool Intel[®] DPC++ Compatibility Tool Intel[®] DPC++ Compatibility Tool Compilers Intel[®] oneAPI DPC++/C++ Compiler Intel[®] oneAPI DPC++/C++ Compiler & Code Intel[®] Fortran Compiler Intel[®] Fortran Compiler Migration Intel[®] Distribution for Python Intel[®] Distribution for Python Libraries Intel[®] oneAPI Libraries Intel[®] oneAPI Libraries Intel[®] Advisor Intel[®] VTune Profiler Intel[®] VTune Profiler Analyzers Intel[®] Advisor Intel[®] VTune Profiler Intel[®] Advisor Start with VTune to analyze code New Program? Start with our industry Start with Intel Advisor to define your performance in your application. leading compliers & libraries found in optimal device mix before optimizing our toolkits for your best performance. and target the right accelerators. Using an older compiler? Re-compiling Porting? Use DPC++ Compatibility Tool Use VTune to further analyze and tune with Intel's newest compiler can offer real to migrate your CUDA code to SYCL for code performance in your application porting across architectures speed-up. Start with Intel System Bring Up Kit to

speed up system bring-up and validation of new hardware

Bring-Up

4 Wor

Working on new hardware?

Intel® System Bring Up Toolkit

intel⁴⁸

oneAPI Resources software.intel.com/oneapi

Get Started

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



oneAPI

Industry Initiative

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



Learn

- Training: <u>Webinars</u> & 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
- Intel[®] DevMesh Innovator
 Projects



 <u>Academic Programs</u>: oneAPI Centers of Excellence: research, enabling code, curriculum, teaching

Domain-specific Toolkits for Specialized Workloads

- Intel[®] oneAPI HPC Toolkit
- Intel® AI Analytics Toolkit
- Intel[®] oneAPI Rendering Toolkit
- Intel[®] oneAPI loT Toolkit
- Intel® Distribution of OpenVINO[™] toolkit, powered by oneAPI

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

ts

Direct Programming	API-Based Programming	Analysis & debug Tools		
Intel [®] C++ Compiler Classic	Intel [®] MPI Library	Intel [®] Inspector		
Intel® Fortran Compiler Classic	Intel® oneAPI DPC++ Library oneDPL	Intel® Trace Analyzer & Collector		
Intel [®] Fortran Compiler	Intel® oneAPI Math Kernel Library - oneMKL	Intel [®] Cluster Checker		
Intel® oneAPI DPC++/C++ Compiler?	Intel [®] oneAPI Data Analytics Library - oneDAL	Intel® VTune™ Profiler		
Intel® DPC++ Compatibility Tool	Intel® oneAPI Threading Building Blocks - oneTBB	Intel [®] Advisor		
Intel® Distribution for Python	Intel® oneAPI Video Processing Library - oneVPL	Intel [®] Distribution for GDB		
Intel® FPGA Add-on for oneAPI Base Toolkit	Intel [®] oneAPI Collective Communications Library oneCCL			
	Intel® oneAPI Deep Neural Network Library - oneDNN	intel		
Intel® oneAPI HPC Toolkit + Intel® oneAPI Base Toolkit	Intel® Integrated Performance Primitives – Intel® IPP	ONEAPI HPC		

Deliver Fast HPC Applications that Scale Customer Use Cases – Intel[®] one API Base & HPC Toolkits



Intel oneAPI tools help prepare code for Aurora. Aurora, Argonne Leadership Computing Facility's Intel-HPE/Cray supercomputer, will be one of the U.S.'s 1st exascale systems

SAMPLE USE CASES & PROOF POINTS



Zuse Institute Berlin (ZIB) ported the *easyWave* tsunami simulation application from CUDA to Data Parallel C++ (DPC++) delivering performance on Intel CPUs, GPUs, FPGAs, & Nvidia P100





Accelerating Google Cloud for HPC C2 provides great performance for HPC workloads: 40% higher performance/core. Runs on Intel® Xeon® processors + AMD, optimized by Intel® oneAPI Base & HPC Toolkits. Video Video | Podcast



Acceleration for HPC & Al Inferencing

CERN, SURFsara, and Intel are investigating approaches driving **breakthrough performance on simulations** used in scientific, engineering, and financial applications*.



Texas Advanced Computing Center (TACC) Frontera SuperComputer Visualization & Filesystem Use Cases Show Value of Large Memory Fat Nodes on Intel® Xeon® processors & Intel® Optane Persistent Memory*



University of Stockholm/KTH

GROMACS, a simulation application used to design new drugs, was optimized by oneAPI. CUDA code was migrated to oneAPI to create new cross-architecture code targeting Intel CPUs and multiple accelerators.

Learn more: *oneAPI Discussions with HPC Thought Leaders* Video [2.20] *Uses Intel® oneAPI Rendering Toolkit



Intel[®] AI Analytics Toolkit

Accelerate end-to-end AI and data analytics pipelines with libraries optimized for Intel® architectures

Who needs this product?

Data scientists, AI researchers, ML and DL developers, AI application developers

Top Features/Benefits

 Deep learning performance for training and inference with Intel optimized DL frameworks and tools

intel

AI

ANALYTICS TOOLKIT

 Drop-in acceleration for data analytics and machine learning workflows with compute-intensive Python packages



Achieve End-to-End Performance for AI Workloads–Customer Use Cases

Accelerate Training + Inference - Most are optimized by Intel[®] AI Analytics Toolkit

May also use Intel[®] Distribution of OpenVINO[™] toolkit, Intel[®] oneAPI Base Toolkit, or Intel[®] oneAPI Rendering Toolkit



Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

Intel[®] oneAPI Rendering Toolkit

Powerful Libraries for High-Fidelity Visualization Applications

- Deliver high-performance, high-fidelity visualization applications on Intel[®] architecture
- Create amazing visual, hyper-realistic renderings via ray tracing with global illumination
- Access all system memory space to create renderings using the largest data sets
- Flexible, cost-efficient development using open-source libraries

YOUR WORK, MORE VIVID.



Intel oneAPI Rendering & Ray Tracing Libraries

Intel[®] Embree

Rasterizer

Award winning, High-Performance, Feature-Rich Ray Tracing & Photorealistic Rendering

Intel[®] Open Image Denoise Al-Accelerated Denoiser for Superior Visual

Ouality

Intel[®] OpenSWR High-Performance, Scalable, OpenGL-Compatible

Intel[®] Open Volume Kernel Library

Render & Simulate 3D Spatial Data Processing

Intel[®] Open Path Guiding Library Path Tracing Importance Sampling, higher fidelity at same or lower compute cost Beta Released O1'22!

Intel[®] OSPRay Scalable, Portable, Distributed Rendering API HPC, SciVis and Pro Render Path Tracer

Intel[®] OSPRay Studio Real-time rendering through a graphical user interface with this new scene graph application

Intel[®] OSPRay for Hydra Connect the Rendering Toolkit libraries to Universal Scene Description Hydra Rendering subsystem via plugin

Learn More: intel.com/oneAPI-RenderKit











¹ The Addams Family 2 - Digital Domain, Marvel Studios, Chaos Group V-Ray ² Scene © Blender Foundation | cloud.blender.org/spring ³ Model from Leigh Orf at University of Wisconsin. For more tornado visualization, visit Leigh Orf's site

⁴ Smoke volume, data courtesy OpenVDB example repository ⁵ Blender Institute: Nishita Demo (CC-BY-ND) ⁶ Copyright Animal Logic Pty Limited 2021. All rights reserved.

Intel® oneAPI Rendering Toolkit

Render Your Vision in Highest Fidelity: Your Open Path to Advanced Ray Tracing



intel

ONEAPT

RENDERING

TOOLKIT

¹Courtesy Baozou Production in association with Tangent Animation using Blender with Intel[®] Embree. Media courtesy of Netflix, Inc. Now streaming on Netflix. Netflix subscription required. *See slide notes for configuration details.

Refer to software.intel.com/articles/optimization-notice for more information regarding performance & optimization choices in Intel software products.

For workloads and configurations visit www.Intel.com/PerformanceIndex. Results may vary. Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

SCIENTIFIC

& TECHNICAI

AWARDS

Intel[®] Embree

Intel® oneAPI Tools for IoT

Accelerate development of IoT applications for smart, connected devices that run at the network's edge for healthcare, smart homes, industrial, retail, aerospace, and more.

Who needs this product?

- OEMs, ODMs, SIs, ISVs
- C, C++, SYCL, OpenMP, Python* Developers

Top Features/Benefits

- Leverage more cores & built-in technologies in IAbased platforms through optimized compilers & libraries
- Easily connect sensors to devices, and devices to cloud with IoT Connection Tools
- Speed development & maintenance of Yocto Project platform projects
- Develop with confidence with powerful analysis tools to identify threading, memory & offloading optimization opportunities
- DPC++ compatibility tool helps migrate existing code written in CUDA

Learn More & Download



Direct Programming	API-Based Programming	Analysis & debug Tools
Intel® C++ Compiler Classic	Intel® oneAPI DPC++ Library oneDPL	Intel [®] Inspector
Eclipse IDE	Intel® oneAPI Math Kernel Library oneMKL	Intel® VTune™ Profiler
Linux Kernel Build Tools	Intel® oneAPI Data Analytics Library - oneDAL	Intel® Advisor
Intel® oneAPI DPC++ /C++ Compiler	Intel® oneAPI Threading Building Blocks - oneTBB	Intel [®] Distribution for GDB
Intel [®] DPC++ Compatibility Tool	Intel® oneAPI Video Processing Library - oneVPL	
Intel [®] Distribution for Python	Intel [®] oneAPI Collective Communications Library oneCCL	
Intel® FPGA Add-on for oneAPI Base Toolkit	Intel® oneAPI Deep Neural Network Library - oneDNN	
	Intel [®] Integrated Performance	
Intel® oneAPI IoT Toolkit +	Primitives – Intel® IPP	intel.
Intel® oneAPI Base Toolkit +		oneAPI

Intel[®] Distribution of OpenVINO[™] toolkit Powered by oneAPI

Deliver High-Performance Deep Learning Inference

A toolkit to accelerate development of high-performance deep learning inference & computer vision in vision/AI applications used from edge to cloud. It enables deep learning on hardware accelerators & easy deployment across Intel[®] CPUs, GPUs, FPGAs, VPUs.

Who needs this product?

- Computer vision, deep learning software developers
- Data scientists
- OEMs, ISVs, System Integrators

Usages

Security surveillance, robotics, retail, healthcare, Al, office automation, transportation, non-vision use cases (speech, NLP, Audio, text) & more

Learn More & Download





Accelerate Development of Smart, Connected Devices Customer Use Cases

May be optimized one or a combination of the Intel® oneAPI Base & IoT Toolkits, Intel® AI Analytics Toolkit, & Intel® Distribution of OpenVINO[™] toolkit





Samsung Medison Uses oneAPI to Power Obstetric Ultrasound Systems Intel® oneAPI Base Toolkit & Intel® Distribution of OpenVINO™ toolkit help accelerate image processing at the edge for consistent measurement accuracy & improved workflows.¹ Intel PR News Byte 09/10/20 | Video

SAMPLE USE CASES & PROOF POINTS



United Imaging Successfully Ported CUDA* Code to oneAPI

<u>United Imaging</u> develops advanced medical imaging diagnosis and treatment products and innovative medical IT solutions. It used <u>Intel</u> <u>oneAPI Base Toolkit</u> for code migration and optimizations.



Optimized by Intel® oneAPI Analytics Toolkit & Intel® Distribution of OpenVINO™ toolkit

Accrad Al-based Solution Helps Accelerate Lung Disease Diagnosis – Acceleration for training + inference.

Learn more in the <u>solution brief</u>

AbbVie Machine Translation Solution

accelerates natural language processing inference models using processor optimized capabilities.

Notices & Disclaimers

Performance varies by use, configuration and other factors. Learn more at www.Intel.com/PerformanceIndex. Results may vary.

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure.

Texas Advanced Computing Center (TACC) Frontera references

Article: <u>HPCWire: Visualization & Filesystem Use Cases Show Value of Large Memory Fat Notes on Frontera</u>. <u>www.intel.com/content/dam/support/us/en/documents/memory-and-storage/data-center-persistent-mem/Intel-Optane-DC-Persistent-Memory-Quick-Start-Guide.pdf</u> <u>software.intel.com/content/www/us/en/develop/articles/introduction-to-programming-with-persistent-memory-from-intel.html</u> <u>wreda.github.io/papers/assise-osdi20.pdf</u>

KFBIO

KFBIO m. tuberculosis screening detectron2 model throughput performance on 2nd Intel® Xeon® Gold 6252 processor: NEW: Test 1 (single instance with PyTorch 1.6: Tested by Intel as of 5/22/2020. 2-socket 2nd Gen Intel® Xeon® Gold 6252 processor; 24 cores, HT On, Turbo ON, Total Memory 192 GB (12 slots/16 GB/2666 MHz), BIOS: SSE5C620.86B.02.01.0008.031920191559 (ucode: 0x500002c), Ubuntu 18.04.4 LTS, kernel 5.3.0-51-generic, mitigated Test 2 (24 instances with PyTorch 1.6: Tested by Intel as of 5/22/2020. 2-socket 2nd Gen Intel Xeon Gold 6252 Processor, 24 cores, HT On, Turbo ON, Total Memory 192 GB (12 slots/16 GB/2666 MHz), BIOS: SSE5C620.86B.02.01.0008.031920191559 (ucode: 0x500002c), Ubuntu 18.04.4 LTS, kernel 5.3.0-51-generic, mitigated BASELINE: (single instance with PyTorch 1.4): Tested by Intel as of 5/22/2020. 2-socket 2nd Gen Intel Xeon Gold 6252 Processor, 24 cores, HT On, Turbo ON, Total Memory 192 GB (12 slots/16 GB/2666 MHz), BIOS: SSE5C620.86B.02.01.0008.031920191559 (ucode: 0x500002c), Ubuntu 18.04.4 LTS, kernel 5.3.0-51-generic, mitigated BASELINE: (single instance with PyTorch 1.4): Tested by Intel as of 5/22/2020. 2-socket 2nd Gen Intel Xeon Gold 6252 Processor, 24 cores, HT On, Turbo ON, Total Memory 192 GB (12 slots/16 GB/2666 MHz), BIOS: SSE5C620.86B.02.01.0008.031920191559 (ucode: 0x500002c), Ubuntu 18.04.4 LTS, kernel 5.3.0-51-generic, mitigated BASELINE: (single instance with PyTorch 1.4): Tested by Intel as of 5/22/2020. 2-socket 2nd Gen Intel Xeon Gold 6252 Processor, 24 cores, HT On, Turbo ON, Total Memory 192 GB (12 slots/16 GB/2666 MHz), BIOS: SSE5C620.86B.02.01.0008.031920191559 (ucode: 0x500002c), Ubuntu 18.04.4 LTS, kernel 5.3.0-51-generic, mitigated.

Tangent Studios

Configurations for Render Times with Intel® Embree, testing conducted by Tangent Animation Labs. Render farm: 8x Intel® Core™ processors +hyperthread*2 + 128gig. In-office workstations: Intel® Xeon® processors HP blade c7000 chassis, with HP460 gen8 blades - 2x Intel Xeon E5-2650 V2, Eight Core 2.6GHz-128GB. Software: Blender 2.78 with custom build using Intel® Embree. For more information on Tangent's work with Embree, watch this video: www.youtube.com/watch?time_continue=251&v=_2la4h8q3xs&feature=emb_logo

Recreation of the performance numbers can be recreated using Agent327, Blender and Embree.

Chaos Group - Up to 90% Memory Reduction for Displacement

Testing conducted by Chaos Group with Intel[®] Embree 2020. Software Corona Renderer 5 with Intel Embree. Up to 90% memory reduction calculated using Corona Renderer 5 with regular displacement grids per triangle of 154 bytes versus Corona Renderer 5 with Intel Embree, which has a displacement capability grid of 12 bytes per grid triangle. (12/154 = 7.8% usage or >90% memory reduction.) Recreation of the performance numbers can be accomplished using Corona Renderer 5 and Embree. For more information, visit the Corona Renderer Blog: <u>blog.corona-renderer.com/corona-renderer-5-for-3ds-max-released/</u>

The Addams Family 2 - Gained a 10% to 20%-and sometimes 25%-efficiency in rendering, saving thousands of hours in rendering production time.

Testing Date: Results are based on data conducted by Cinesite 2020-21. 10% to up to 25% rendering efficiency/thousands of hours saved in rendering production time/15 hrs per frame per shot to 12-13 hrs. Cinesite Configuration: 18-core Intel® Xeon® Scalable processors (W-2295) used in render farm, 2nd gen Intel Xeon processor-based workstations (W-2135 and -2195) used. Rendering tools: Gaffer, Arnold, along with optimizations by Intel® Open Image Denoise.

Your costs and results may vary.

Intel technologies may require enabled hardware, software or service activation.

Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

© Intel Corporation. Intel, the Intel logo, Xeon, Core, VTune, OpenVINO, Agilex, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

oneAPI Ecosystem Support



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.

Intel[®] oneAPI Toolkits – Proven Performance **Top Customer Use Cases**

HPC Cross-architecture

- Argonne National Labs is running exascale-class applications efficiently on current and future Intel CPUs and GPUs.
- Zuse Institute Berlin (ZIB) ported the tsunami simulation easy Wave application from CUDA to DPC++ delivering performance across multiple architectures and vendors.
- GROMACS, accelerated by oneAPI open programming and multiarchitecture tools, runs on Intel Xe architecture-based GPUs with strong performance.
- HPC & AI CERN uses Intel[®] DL Boost and oneAPI to speed simulations with inference acceleration by nearly 2x without accuracy loss*
- High-Fidelity Visualization Using Advanced Ray Tracing Bentley Motors Limited's Al-based car configurator processes 1.7M+ images with up to 10B possible configurations per model. The Addams Family 2, Cinesite used Intel's denoiser achieving 25% rendering efficiency.
- IoT Samsung Medison accelerates ultrasound image processing at the edge on multiple Intel® architectures for improved accuracy and fast diagnosis
- Major CSPs & Framework endorse oneAPI Microsoft Azure, Google Cloud, TensorFlow
- FPGA Using oneAPI, Bittware had its application running in days vs. what typically would take several weeks using Verilog or VHDL*
- And more... 250+ applications developed using oneAPI tools > view catalog







Intel[®] oneAPI Cross-architecture Tools





Video [3:45]

Intel[®] oneAPI Toolkits Free Availability

Get Started Quickly

Code Samples, Quick-start Guides, Webinars, Training

software.intel.com/oneapi



Intel® oneAPI Priority Support

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



With All Paid Licenses

Intel® oneAPI Community

Intel[®] oneAPI Base Toolkit Free Intel[®] Community Support Q This board ~ Support via the Intel public Community Forum Access to only the latest version of oneAPI Toolkits Intel Communities / Developer Software Forums / Toolkits & SDKs / Intel® oneAPI Base Toolkit 261 Discussion Access to online tutorials and self-help forums Support for core tools and libraries to build and deploy high-performance datacentric applications Discussions Post a question

Why SYCL?

```
sycl::queue q(cpu selector{});
                                                  sycl::queue q(gpu selector{});
auto A = sycl::malloc shared<float>(n, q);
                                                  auto A = sycl::malloc shared<float>(n, q);
auto B = sycl::malloc shared<float>(n, q);
                                                  auto B = sycl::malloc shared<float>(n, q);
q.parallel for( sycl::range<1>{n},
                                                  q.parallel for( sycl::range<1>{n},
  [=] (sycl::id<1> i) {
                                                    [=] (sycl::id<1> i) {
      B[i] += A[i] * A[i];
                                                        B[i] += A[i] * A[i];
).wait();
                                                  ).wait();

    Multi-Architecture (CPU, GPU,

Standard C++17 aids time
                                    Unified shared
                                                                FPGA & other targets)
 to developer productivity
                                      memory
```

 Syntax for accelerators (device selection, offload, memory transfer)

- Single source (host & device code)
- Stack based on standards & open specifications (CLANG, LLVM, SPIR-V, Level Zero)

Data Parallel C++: oneAPI's implementation of SYCL

DPC++ = ISO C++ and Khronos SYCL and community extensions

Freedom of Choice: Future-Ready Programming Model

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

DPC++ = ISO C++ and Khronos SYCL and community extensions

- Designed for data parallel programming productivity
- Provides full native high-level language performance on par with standard C++ and broad compatibility
- Adds SYCL from the Khronos Group for data parallelism and heterogeneous programming

Community Project Drives Language Enhancements

- Provides extensions to simplify data parallel programming
- Continues evolution through open and cooperative development

