



# Intel<sup>®</sup> FPGA SDK for OpenCL<sup>™</sup> Pro Edition

---

## Programming Guide

Updated for Intel<sup>®</sup> Quartus<sup>®</sup> Prime Design Suite: **22.4**



[Online Version](#)



[Send Feedback](#)

UG-OCL002

ID: **683846**

Version: **2022.12.19**

## Contents

---

<b>1. Intel® FPGA SDK for OpenCL™ Overview.....</b>	<b>7</b>
1.1. Intel FPGA SDK for OpenCL Pro Edition Programming Guide Prerequisites.....	7
1.2. Intel FPGA SDK for OpenCL FPGA Programming Flow.....	8
<b>2. Intel FPGA SDK for OpenCL Offline Compiler Kernel Compilation Flows.....</b>	<b>10</b>
2.1. One-Step Compilation for Simple Kernels.....	11
2.2. Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow.....	12
<b>3. Obtaining General Information on Software, Compiler, and Custom Platform.....</b>	<b>16</b>
3.1. Displaying the Software Version (version).....	16
3.2. Displaying the Compiler Version (-version).....	16
3.3. Listing the Intel FPGA SDK for OpenCL Utility Command Options (help).....	17
3.3.1. Displaying Information on an Intel FPGA SDK for OpenCL Utility Command Option (help <command_option>).....	17
3.4. Listing the Intel FPGA SDK for OpenCL Offline Compiler Command Options (no argument, -help, or -h).....	17
3.5. Listing the Available FPGA Boards and Custom Platforms (-list-boards and -list- board-packages).....	17
3.6. Displaying the Compilation Environment of an OpenCL Binary (env).....	18
<b>4. Managing an FPGA Board.....</b>	<b>19</b>
4.1. Installing an FPGA Board (install).....	20
4.2. Uninstalling an FPGA Board (uninstall).....	22
4.3. Querying the Device Name of Your FPGA Board (diagnose).....	22
4.4. Running a Board Diagnostic Test (diagnose <device_name>).....	23
4.5. Programming the FPGA Offline or without a Host (program <device_name>).....	23
4.6. Programming the Flash Memory (flash <device_name>).....	23
<b>5. Structuring Your OpenCL Kernel.....</b>	<b>25</b>
5.1. Guidelines for Naming the Kernel.....	25
5.2. Programming Strategies for Optimizing Data Processing Efficiency.....	26
5.2.1. Unrolling a Loop (unroll Pragma) .....	27
5.2.2. Disabling Pipelining of a Loop (disable_loop_pipelining Pragma) .....	28
5.2.3. Coalescing Nested Loops.....	29
5.2.4. Fusing Adjacent Loops (loop_fuse Pragma).....	30
5.2.5. Marking Loops to Prevent Automatic Fusion (nofusion Pragma).....	32
5.2.6. Specifying a Loop Initiation interval (II).....	33
5.2.7. Loop Concurrency (max_concurrency Pragma) .....	34
5.2.8. Loop Speculation (speculated_iterations Pragma) .....	35
5.2.9. Loop Interleaving Control (max_interleaving Pragma).....	35
5.2.10. Floating Point Optimizations (fp_contract and fp_reassociate Pragma).....	38
5.2.11. Specifying Work-Group Sizes.....	39
5.2.12. Specifying Number of Compute Units.....	40
5.2.13. Specifying Number of SIMD Work-Items.....	40
5.2.14. Specifying the private_copies Memory Attribute.....	41
5.2.15. Specifying the use_stall_enable_clusters Cluster-control Attribute .....	42
5.3. Programming Strategies for Optimizing Pointer-to-Local Memory Size.....	42

5.4. Implementing the Intel FPGA SDK for OpenCL Channels Extension.....	43
5.4.1. Overview of the Intel FPGA SDK for OpenCL Channels Extension.....	43
5.4.2. Channel Data Behavior.....	44
5.4.3. Multiple Work-Item Ordering for Channels.....	45
5.4.4. Restrictions in the Implementation of Intel FPGA SDK for OpenCL Channels Extension .....	46
5.4.5. Enabling the Intel FPGA SDK for OpenCL Channels for OpenCL Kernel.....	48
5.5. Implementing OpenCL Pipes.....	60
5.5.1. Overview of the OpenCL Pipe Functions.....	61
5.5.2. Pipe Data Behavior.....	62
5.5.3. Multiple Work-Item Ordering for Pipes.....	62
5.5.4. Restrictions in OpenCL Pipes Implementation.....	64
5.5.5. Enabling OpenCL Pipes for Kernels.....	65
5.5.6. Direct Communication with Kernels via Host Pipes.....	75
5.6. Implementing Arbitrary Precision Integers.....	82
5.7. Using Predefined Preprocessor Macros in Conditional Compilation.....	83
5.8. Declaring __constant Address Space Qualifiers.....	83
5.9. Including Structure Data Types as Arguments in OpenCL Kernels.....	84
5.9.1. Matching Data Layouts of Host and Kernel Structure Data Types.....	84
5.9.2. Disabling Insertion of Data Structure Padding .....	86
5.9.3. Specifying the Alignment of a Struct.....	86
5.10. Inferring a Register.....	87
5.10.1. Inferring a Shift Register.....	88
5.11. Enabling Double Precision Floating-Point Operations.....	89
5.12. Single-Cycle Floating-Point Accumulator for Single Work-Item Kernels.....	89
5.12.1. Programming Strategies for Inferring the Accumulator.....	90
5.13. Integer Promotion Rules.....	91
<b>6. Designing Your Host Application.....</b>	<b>93</b>
6.1. Host Programming Requirements.....	93
6.1.1. Host Machine Memory Requirements.....	93
6.1.2. Host Binary Requirement.....	93
6.1.3. Multiple Host Threads.....	94
6.1.4. Out-of-order Command Queues.....	94
6.1.5. Requirement for Multiple Command Queues to Execute Kernels Concurrently... ..	94
6.2. Allocating OpenCL Buffers for Manual Partitioning of Global Memory.....	94
6.2.1. Partitioning Buffers Across Multiple Interfaces of the Same Memory Type.....	95
6.2.2. Partitioning Buffers Across Different Memory Types (Heterogeneous Memory)..	96
6.2.3. Creating a Pipe Object in Your Host Application.....	98
6.2.4. Enabling All Global Memory.....	98
6.3. Triggering Collection Profiling Data During Kernel Execution.....	99
6.3.1. Profiling Autorun Kernels .....	101
6.3.2. Profile Data Acquisition.....	103
6.4. Accessing Custom Platform-Specific Functions.....	103
6.5. Modifying Host Program for Structure Parameter Conversion.....	104
6.6. Managing Host Application.....	105
6.6.1. Displaying Example Makefile Fragments (example-makefile or makefile).....	105
6.6.2. Compiling and Linking Your Host Application.....	106
6.6.3. Using OpenCL ICD Extension APIs.....	110
6.6.4. Programming an FPGA via the Host.....	110
6.6.5. Termination of the Runtime Environment and Error Recovery.....	115

6.7. Allocating Shared Memory for OpenCL Kernels Targeting SoCs.....	115
6.8. Sharing Multiple Devices Across Multiple Host Programs.....	117
<b>7. Compiling Your OpenCL Kernel.....</b>	<b>119</b>
7.1. Compiling Your Kernel to Create Hardware Configuration File.....	119
7.2. Compiling Your Kernel without Building Hardware (-c).....	120
7.3. Compiling and Linking Your Kernels or Object Files without Building Hardware (-rtl)....	120
7.4. Specifying the Location of Header Files (-I=<directory>).....	121
7.5. Specifying the Name of an Intel FPGA SDK for OpenCL Offline Compiler Output File (-o <filename>).....	122
7.6. Compiling a Kernel for a Specific FPGA Board and Custom Platform (- board=<board_name>) and (-board-package=<board_package_path>).....	122
7.7. Resolving Hardware Generation Fitting Errors during Kernel Compilation (-high-effort)	124
7.8. Specifying Schedule $F_{max}$ Target for Kernels (-clock=<clock_target>) .....	125
7.9. Defining Preprocessor Macros to Specify Kernel Parameters (-D<macro_name>).....	125
7.10. Generating Compilation Progress Report (-v).....	126
7.11. Displaying the Estimated Resource Usage Summary On-Screen (-report).....	128
7.12. Suppressing Warning Messages from the Intel FPGA SDK for OpenCL Offline Compiler (-W).....	128
7.13. Converting Warning Messages from the Intel FPGA SDK for OpenCL Offline Compiler into Error Messages (-Werror).....	128
7.14. Removing Debug Data from Compiler Reports and Source Code from the .aocx File (-g0).....	129
7.15. Disabling Burst-Interleaving of Global Memory (-no- interleaving=<global_memory_type>).....	129
7.16. Forcing Ring Interconnect for Global Memory (-global-ring).....	130
7.17. Forcing a Single Store Ring to Reduce Area at the Expense of Write Throughput to Global Memory (-force-single-store-ring).....	130
7.18. Forcing Fewer Read Data Reorder Units to Reduce Area at the Expense of Read Throughput to Global Memory (-num-reorder).....	130
7.19. Configuring Constant Memory Cache Size (-const-cache-bytes=<N>).....	130
7.20. Relaxing the Order of Floating-Point Operations (-ffp-reassociate).....	131
7.21. Reducing Floating-Point Rounding Operations (-ffp-contract=fast).....	131
7.22. Speeding Up Your OpenCL Compilation (-fast-compile).....	132
7.23. Compiling Your Kernel Incrementally (-incremental).....	132
7.23.1. The Incremental Compile Report.....	133
7.23.2. Additional Command Options for Incremental Compilation.....	135
7.23.3. Limitations of the Incremental Compilation Feature.....	137
7.24. Compiling Your Kernel with Memory Error Correction Coding (-ecc).....	137
7.25. Disabling Hardware Kernel Invocation Queue (-no-hardware-kernel-invocation- queue).....	138
7.26. Modifying the Handshaking Protocol (-hyper-optimized-handshaking).....	138
7.27. Pipelining Loops in Non-task Kernels (-auto-pipeline).....	139
<b>8. Emulating and Debugging Your OpenCL Kernel.....</b>	<b>140</b>
8.1. Setting up the Emulator.....	140
8.2. Modifying Channels Kernel Code for Emulation.....	141
8.2.1. Emulating a Kernel that Passes Pipes or Channels by Value.....	142
8.2.2. Emulating Channel Depth.....	143
8.2.3. Emulating Applications with a Channel That Reads or Writes to an I/O Channel.....	143
8.3. Compiling a Kernel for Emulation (-march=emulator).....	144

8.4. Emulating Your OpenCL Kernel.....	144
8.5. Debugging Your OpenCL Kernel on Linux.....	146
8.6. Limitations of the Intel FPGA SDK for OpenCL Emulator.....	147
8.7. Discrepancies in Hardware and Emulator Results.....	148
8.8. Emulator Environment Variables.....	150
8.9. Extensions Supported by the Emulator.....	151
8.10. Emulator Known Issues.....	152
<b>9. Developing OpenCL Applications Using Third-party IDEs.....</b>	<b>153</b>
9.1. FPGA Workflows in Microsoft Visual Studio.....	153
9.1.1. Preparing the Visual Studio Environment.....	153
9.1.2. Creating an FPGA OpenCL Template.....	153
9.1.3. Configuring the Build Targets.....	154
9.1.4. Configuring Build Options for a Project.....	155
9.1.5. Generating the High-level Design Report.....	155
9.1.6. Building and Running the FPGA Template.....	156
9.2. FPGA Workflows in Eclipse.....	156
9.2.1. Preparing the Eclipse Environment.....	157
9.2.2. Creating a Simple FPGA application.....	157
9.2.3. Creating a Makefile Project.....	162
9.2.4. Building a Project.....	163
9.3. Limitations.....	164
<b>10. Developing OpenCL Applications Using Intel Code Builder for OpenCL.....</b>	<b>165</b>
10.1. Configuring the Intel Code Builder for OpenCL Offline Compiler Plug-in for Microsoft Visual Studio.....	165
10.2. Configuring the Intel Code Builder for OpenCL Offline Compiler Plug-in for Eclipse.....	165
10.3. Creating a Session in the Intel Code Builder for OpenCL .....	166
10.4. Configuring a Session.....	167
<b>11. Intel FPGA SDK for OpenCL Advanced Features.....</b>	<b>169</b>
11.1. OpenCL Library.....	169
11.1.1. Creating Library Objects From OpenCL Code.....	172
11.1.2. Understanding RTL Modules and the OpenCL Pipeline.....	174
11.1.3. Packaging an OpenCL Helper Function File for an OpenCL Library.....	190
11.1.4. Packaging an RTL Component for an OpenCL Library .....	190
11.1.5. Verifying the RTL Modules.....	192
11.1.6. Specifying an OpenCL Library when Compiling an OpenCL Kernel.....	193
11.1.7. Debugging Your OpenCL Library Through Simulation (Preview).....	194
11.1.8. Using an OpenCL Library that Works with Simple Functions (Example 1).....	197
11.1.9. Using an OpenCL Library that Works with External Memory (Example 2).....	197
11.1.10. OpenCL Library Command-Line Options.....	198
11.2. Memory Attributes for Configuring Kernel Memory Systems.....	199
11.2.1. Restrictions on the Use of Variable-specific Attributes.....	201
11.3. Kernel Attributes for Reducing the Overhead on Hardware Usage.....	202
11.3.1. Hardware for Kernel Interface.....	202
11.4. Kernel Replication Using the num_compute_units(X,Y,Z) Attribute.....	205
11.4.1. Customization of Replicated Kernels Using the get_compute_id() Function...	206
11.4.2. Using Channels with Kernel Copies.....	207
11.5. Intra-Kernel Registered Assignment Built-In Function .....	208

- A. Support Statuses of OpenCL Features ..... 210**
  - A.1. Support Statuses of OpenCL 1.0 Features.....210
    - A.1.1. OpenCL 1.0 C Programming Language Implementation..... 210
    - A.1.2. OpenCL C Programming Language Restrictions.....212
    - A.1.3. Argument Types for Built-in Geometric Functions.....213
    - A.1.4. Numerical Compliance Implementation..... 214
    - A.1.5. Image Addressing and Filtering Implementation..... 214
    - A.1.6. Atomic Functions..... 214
    - A.1.7. Embedded Profile Implementation..... 215
  - A.2. Support Statuses of OpenCL 1.2 Features.....215
    - A.2.1. OpenCL 1.2 Runtime Implementation.....216
    - A.2.2. OpenCL 1.2 C Programming Language Implementation..... 216
  - A.3. Support Statuses of OpenCL 2.0 Features.....217
    - A.3.1. OpenCL 2.0 Headers.....217
    - A.3.2. OpenCL 2.0 Runtime Implementation.....217
    - A.3.3. OpenCL 2.0 C Programming Language Restrictions for Pipes..... 217
  - A.4. Intel FPGA SDK for OpenCL Allocation Limits.....219
- B. Intel FPGA SDK for OpenCL Pro Edition Programming Guide Archives..... 220**
- C. Document Revision History of the Intel FPGA SDK for OpenCL Pro Edition Programming Guide..... 221**

## 1. Intel® FPGA SDK for OpenCL™ Overview

---

The *Intel® FPGA SDK for OpenCL™ Programming Guide* provides descriptions, recommendations and usage information on the Intel Software Development Kit (SDK) for OpenCL compiler and tools. The Intel FPGA SDK for OpenCL<sup>(1)</sup> is an OpenCL<sup>(2)</sup>-based heterogeneous parallel programming environment for Intel FPGA products.

*Tip:* If you are looking for information about oneAPI programming model, including information about Data Parallel C++ (DPC++), then refer to the [Intel oneAPI Programming Guide](#).

### 1.1. Intel FPGA SDK for OpenCL Pro Edition Programming Guide Prerequisites

The *Intel FPGA SDK for OpenCL Pro Edition Programming Guide* assumes that you are knowledgeable in OpenCL concepts and application programming interfaces (APIs). It also assumes that you have experience creating OpenCL applications and are familiar with the OpenCL Specification version 1.0.

Before using the Intel FPGA SDK for OpenCL or the Intel FPGA Runtime Environment (RTE) for OpenCL to program your device, familiarize yourself with the respective getting started guides. This document assumes that you have performed the following tasks:

- For developing and deploying OpenCL kernels, download the tar file and run the installers to install the SDK, the Intel Quartus® Prime Pro Edition software, and device support.
- For deployment of OpenCL kernels, download and install the RTE.
- If you want to use the SDK or the RTE to program an Intel SoC FPGA, you also have to download and install the IntelSoC FPGA Embedded Development Suite (EDS) Pro Edition.
- Install and set up your FPGA board.
- Verify that board installation is successful, and the board functions correctly.

If you have not performed the tasks described above, refer to the SDK's getting starting guides for more information.

---

(1) The Intel FPGA SDK for OpenCL is based on a published Khronos Specification, and has passed the Khronos Conformance Testing Process. Current conformance status can be found at [www.khronos.org/conformance](http://www.khronos.org/conformance).

(2) OpenCL and the OpenCL logo are trademarks of Apple Inc. used by permission of the Khronos Group™.

**Related Information**

- [OpenCL References Pages](#)
- [OpenCL Specification version 1.0](#)
- [Intel FPGA SDK for OpenCL Pro Edition Getting Started Guide](#)
- [Intel FPGA RTE for OpenCL Pro Edition Getting Started Guide](#)

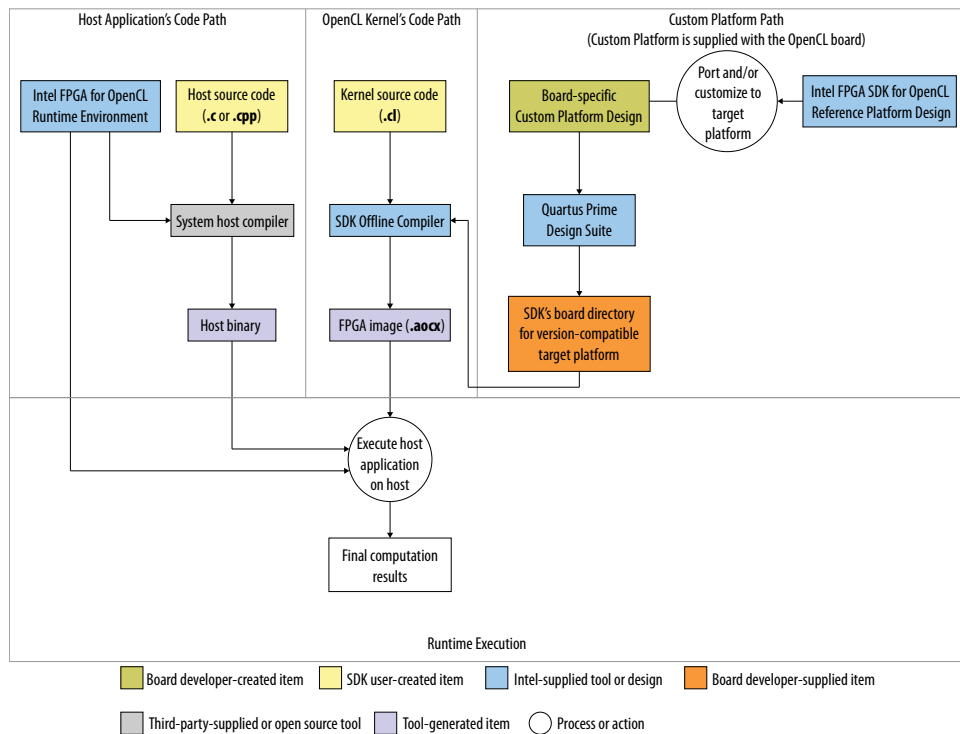
**1.2. Intel FPGA SDK for OpenCL FPGA Programming Flow**

Applications using the Intel FPGA SDK for OpenCL have two main components: the FPGA programming bitstream(s), and the host program that manages the application and FPGA accelerator.

*Tip:* If you are looking for information about oneAPI DPC++ FPGA programming flow, then refer to the [Intel oneAPI Programming Guide](#).

The Intel FPGA SDK for OpenCL Offline Compiler first compiles your OpenCL kernels to an image file that the host program uses to program the FPGA. The host-side C compiler compiles your host program and then links it to the Intel FPGA SDK for OpenCL runtime libraries.

**Figure 1. Schematic Diagram of the Intel FPGA SDK for OpenCL Programming Model**



The following SDK components work together to program an Intel FPGA:

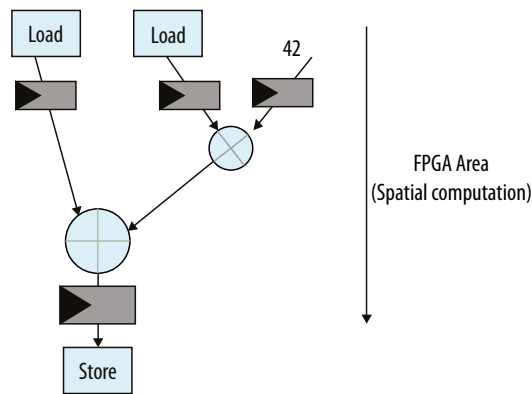
- The host application and the host compiler
- The OpenCL kernel(s) and the offline compiler
- The Custom Platform



The Custom Platform provides the board support package. Typically, the board manufacturer develops the Custom Platform that supports a specific OpenCL board. The offline compiler targets the Custom Platform when compiling an OpenCL kernel to generate a hardware programming image. The host then runs the host application, which usually programs and executes the hardware image onto the FPGA.

In a sequential implementation of a program (for example, on a conventional processor), the program counter controls the sequence of instructions that are executed on the hardware, and the instructions that execute on the hardware across time. In a spatial implementation of a program, such as program implementation within the Intel FPGA SDK for OpenCL, instructions are executed as soon as the prerequisite data is available. Programs are interpreted as a series of connections representing the data dependencies.

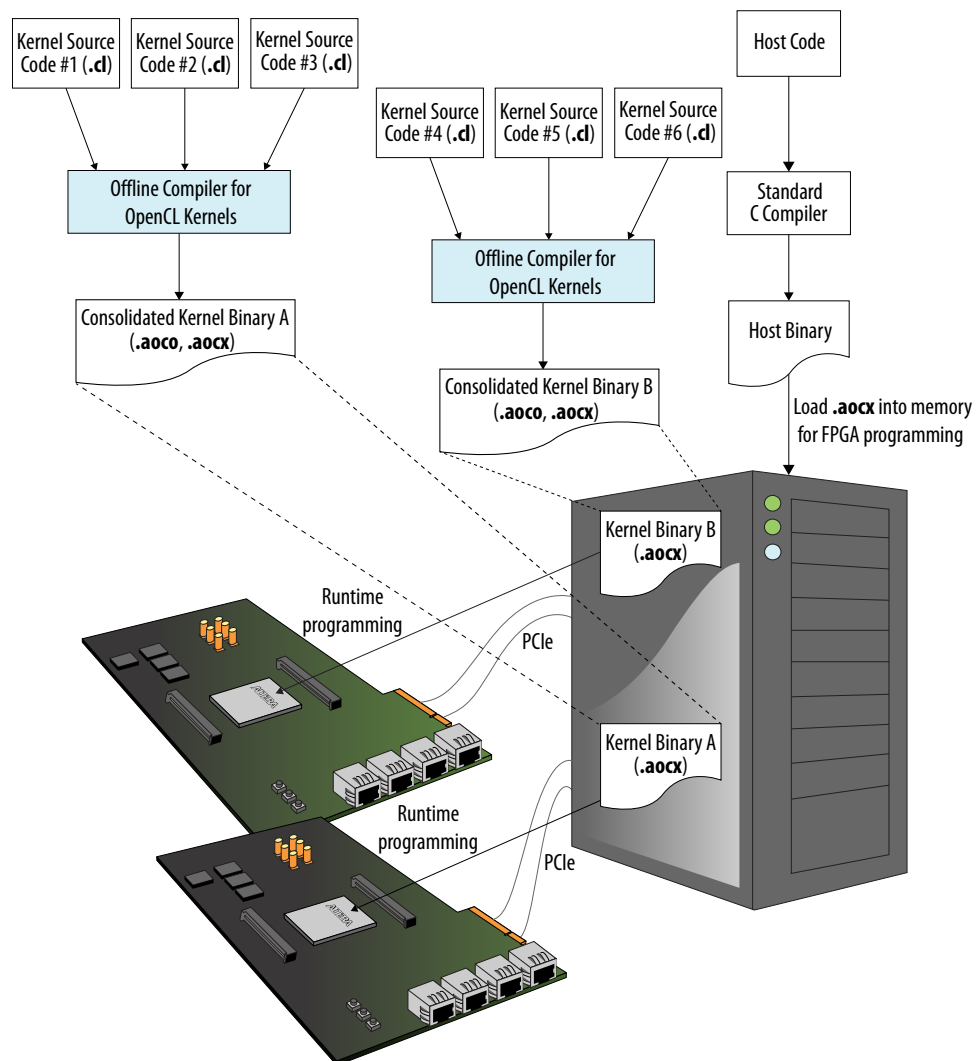
**Figure 2. FPGA Data Flow Architecture**



## 2. Intel FPGA SDK for OpenCL Offline Compiler Kernel Compilation Flows

The Intel FPGA SDK for OpenCL Offline Compiler can create your FPGA programming image file (.aocx file) in a single or multistep process. The complexity of your kernels dictates the preferred compilation.

**Figure 3. The Intel FPGA SDK for OpenCL FPGA Programming Flow**



Intel Corporation. All rights reserved. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Intel warrants performance of its FPGA and semiconductor products to current specifications in accordance with Intel's standard warranty, but reserves the right to make changes to any products and services at any time without notice. Intel assumes no responsibility or liability arising out of the application or use of any information, product, or service described herein except as expressly agreed to in writing by Intel. Intel customers are advised to obtain the latest version of device specifications before relying on any published information and before placing orders for products or services.

\*Other names and brands may be claimed as the property of others.

An OpenCL kernel source file (.cl) contains your OpenCL kernel source code that runs on the FPGA. The offline compiler groups one or more kernels into a temporary file and then compiles this file to generate the following files and folders:

- A .aoco object file is an intermediate object file that contains information for later stages of the compilation. It is not presented unless the aoc option `-save-temps` is specified.
- A .aocx image file is the hardware configuration file and contains information necessary to program the FPGA at runtime.
- The work folder or subdirectory, which contains data necessary to create the .aocx file. By default, the name of the work directory is the name of your .cl file. If you compile multiple kernel source files, the name of the work directory is the name of the last .cl file you list in the aoc command line.

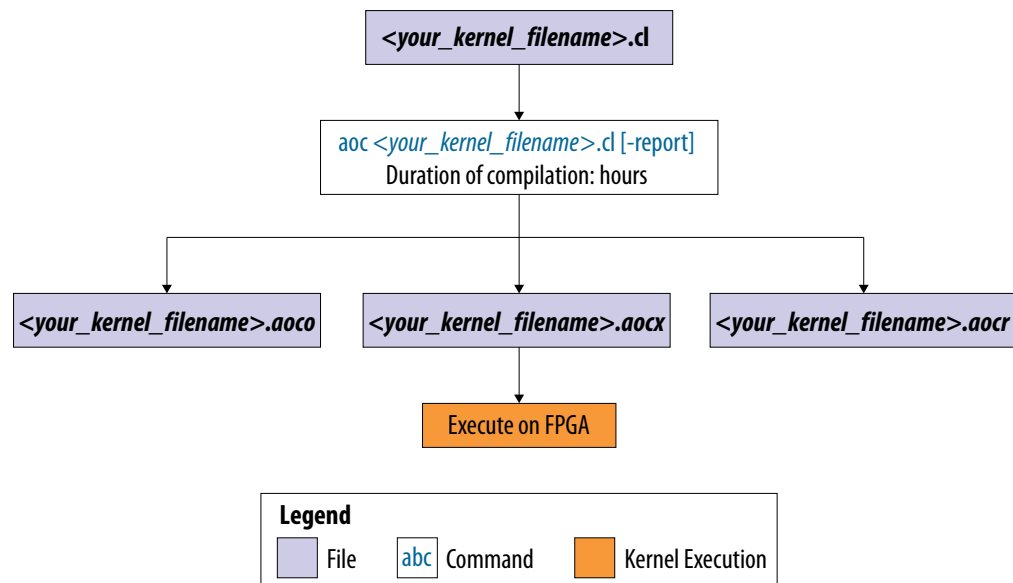
The .aocx file contains data that the host application uses to create program objects, a concept within the OpenCL runtime API, for the target FPGA. The host application first loads these program objects into memory. Then the host runtime uses these program objects to program the target FPGA, as required for kernel launch operations by the host program.

## 2.1. One-Step Compilation for Simple Kernels

By default, the Intel FPGA SDK for OpenCL Offline Compiler compiles one or more OpenCL kernels and creates a hardware configuration file in a single step. Choose this compilation option if you have a single OpenCL kernel in your application, or if you have multiple kernels in your application that require minimal iterations.

The following figure illustrates the OpenCL kernel design flow that has a single compilation step.

Figure 4. One-Step OpenCL Kernel Compilation Flow



A successful compilation results in the following files and reports:

- A .aocr file
- A .aocx file
- In the `<your_kernel_filename>/reports/report.html` file, the estimated resource usage summary provides a preliminary assessment of area usage. If you have a single work-item kernel, the optimization report identifies performance bottlenecks.

**Attention:** It is very time consuming to iterate on your design using the one-step compilation flow. For each iteration, you must perform a full compilation for FPGA hardware, which takes hours. Then you must execute the kernel on the FPGA to measure its performance.

#### Related Information

[Compiling Your Kernel to Create Hardware Configuration File](#) on page 119

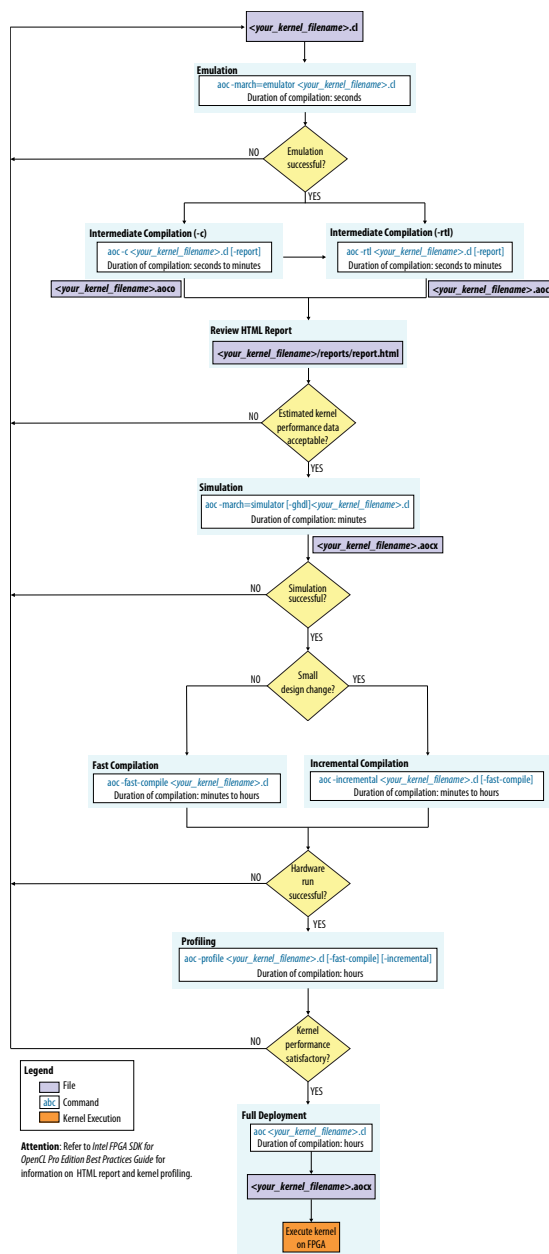
## 2.2. Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow

Choose the multistep Intel FPGA SDK for OpenCL Pro Edition design flow if you want to iterate on your OpenCL kernel design to implement optimizations or other iterative modifications.

**Tip:** If you are looking for or interested in learning about DPC++ FPGA design flows, refer to [Types of DPC++ FPGA Compilation](#) in the *Intel oneAPI Programming Guide*.

The following figure outlines the stages in the SDK's design flow. The steps in the design flow serve as checkpoints for identifying functional errors and performance bottlenecks. They allow you to modify your OpenCL kernel code without performing a full compilation on each iteration. You have the option to perform some or all of the compilations steps.

Figure 5. Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow



The SDK's design flow includes the following steps:

### 1. Emulation

Assess the functionality of your OpenCL kernel by executing it on one or multiple emulation devices on an x86-64 host system. For Linux systems, the Emulator offers symbolic debug support. Symbolic debug allows you to locate the origins of functional errors in your kernel code.

### 2. Intermediate Compilation

There are two available intermediate compilation steps. You have the option to include one or both of these compilation steps in your design flow.

- Compile one or more `.cl` kernel source files using the `-c` flag. Doing so instructs the offline compiler to generate `.aoco` object files that contain the output from the OpenCL parser.
- Compile one or more `.cl` kernel source files or `.aoco` files, but not both, using the `-rtl` flag. Doing so instructs the offline compiler to perform the following tasks:
  - If the input files are `.cl` files, the offline compiler generates an intermediate `.aoco` file for each kernel source file and then links them to generate a `.aocr` file.
  - If the input files are `.aoco` files, the offline compiler links them to generate a `.aocr` file.
  - Creates a `<your_kernel_filename>` directory.

The offline compiler uses the `.aocr` file to generate the final `.aocx` hardware configuration file.

*Note:* If you compile your kernel(s) using the `-c` flag in an environment where the default board is X, and then you compile your `.aoco` files using the `-rtl` flag in an environment where the default board is Y, the offline compiler reads board X from the `.aoco` files and then pass it on to the subsequent compilation stages.

### 3. Review HTML Report

Review the `<your_kernel_filename>/reports/report.html` file of your OpenCL application to determine whether the estimated kernel performance data is acceptable. The HTML report also provides suggestions on how you can modify your kernel to increase performance.

### 4. Simulation (Preview)

Assess the functionality of your OpenCL kernel by running it through simulation. Simulation lets you assess the function correctness and dynamic performance of your kernel without a long compilation time. You can capture and view waveforms for your kernel to help you debug your kernel.

*Note:* The performance of the simulator is very slow when compared to that of hardware. So, Intel recommends using a smaller data set for testing.

### 5. Fast Compilation

Assess the functionality of your OpenCL kernel in hardware. The fast compilation step generates a `.aocx` file in a fraction of the time required to complete a full compilation. The Intel FPGA SDK for OpenCL Offline Compiler reduces compilation time by performing only light optimizations.

*Note:* Fast compilation allows you to compile quickly. However, the performance is inferior than a regular Intel Quartus Prime Pro Edition compilation.

### 6. Incremental Compilation

Assess the functionality of your OpenCL kernel in hardware. The incremental compilation step generates a `.aocx` file by compiling only the kernels you have modified. The Intel FPGA SDK for OpenCL Offline Compiler improves your productivity by scaling compilation times with the size of your design changes rather than the size of your overall design.

#### 7. Profiling

Instruct the Intel FPGA SDK for OpenCL Offline Compiler to insert performance counters in the FPGA programming image. During execution, the counters collect performance information that you can then review in the Intel FPGA dynamic profiler for OpenCL GUI.

#### 8. Full deployment

When you are satisfied with the performance of your OpenCL kernel throughout the design flow, perform a full compilation. The resulting `.aocx` file is suitable for deployment.

For more information about the HTML report and kernel profiling, refer to the *Intel FPGA SDK for OpenCL Pro Edition Best Practices Guide*.

#### Related Information

- [Compiling Your Kernel Incrementally \(-incremental\)](#) on page 132
- [Speeding Up Your OpenCL Compilation \(-fast-compile\)](#) on page 132
- [Reviewing Your Kernel's report.html File](#)
- [Compiling Your OpenCL Kernel](#) on page 119
- [Emulating and Debugging Your OpenCL Kernel](#) on page 140
- [Debugging Your OpenCL Library Through Simulation \(Preview\)](#) on page 194
- [Profiling Your Kernel to Identify Performance Bottlenecks](#)

## 3. Obtaining General Information on Software, Compiler, and Custom Platform

---

The Intel FPGA SDK for OpenCL includes two sets of command options: the SDK utility commands (`aocl <command_option>`) and the Intel FPGA SDK for OpenCL Offline Compiler commands (`aoc <command_option>`). Each set of commands includes options you can invoke to obtain general information on the software, the compiler, and the Custom Platform.

[Displaying the Software Version \(`version`\) on page 16](#)

[Displaying the Compiler Version \(`-version`\) on page 16](#)

[Listing the Intel FPGA SDK for OpenCL Utility Command Options \(`help`\) on page 17](#)

[Listing the Intel FPGA SDK for OpenCL Offline Compiler Command Options \(`no argument`, `-help`, or `-h`\) on page 17](#)

[Listing the Available FPGA Boards and Custom Platforms \(`-list-boards` and `-list-board-packages`\) on page 17](#)

[Displaying the Compilation Environment of an OpenCL Binary \(`env`\) on page 18](#)

### 3.1. Displaying the Software Version (`version`)

To display the version of the Intel FPGA SDK for OpenCL, invoke the `version` utility command.

- At the command prompt, invoke the `aocl version` command.  
Example output:

```
aocl <version>.<build> (Intel(R) FPGA SDK for OpenCL(TM),  
Version <version> Build <build>, Copyright (C) <year> Intel  
Corporation)
```

### 3.2. Displaying the Compiler Version (`-version`)

To display the version of the Intel FPGA SDK for OpenCL Offline Compiler, invoke the `-version` compiler command.

- At a command prompt, invoke the `aoc -version` command.  
Example output:

```
Intel(R) FPGA SDK for OpenCL(TM), 64-Bit Offline Compiler  
Version <version> Build <build>  
Copyright (C) <year> Intel Corporation
```



### 3.3. Listing the Intel FPGA SDK for OpenCL Utility Command Options (help)

To display information on the Intel FPGA SDK for OpenCL utility command options, invoke the `help` utility command.

- At a command prompt, invoke the `aocl help` command.  
The SDK categorizes the utility command options based on their functions. It also provides a description for each option.

#### 3.3.1. Displaying Information on an Intel FPGA SDK for OpenCL Utility Command Option (help <command\_option>)

To display information on a specific Intel FPGA SDK for OpenCL utility command option, include the command option as an argument of the `help` utility command.

- At a command prompt, invoke the `aocl help <command_option>` command.  
For example, to obtain more information on the `install` utility command option, invoke the `aocl help install` command.

Example output:

```
aocl install - Installs a board onto your host system.  
  
Usage: aocl install  
  
Description:  
This command installs a board's drivers and other necessary software for the  
host operating system to communicate with the board.  
For example this might install PCIe drivers.
```

### 3.4. Listing the Intel FPGA SDK for OpenCL Offline Compiler Command Options (no argument, -help, or -h)

To display information on the Intel FPGA SDK for OpenCL Offline Compiler command options, invoke the compiler command without an argument, or invoke the compiler command with the `-help` or `-h` command option.

- At a command prompt, invoke one of the following commands:
  - `aoc`
  - `aoc -help`
  - `aoc -h`

The SDK categorizes the offline compiler command options based on their functions. It also provides a description for each option.

### 3.5. Listing the Available FPGA Boards and Custom Platforms (-list-boards and -list-board-packages)

To list the FPGA boards available in your Custom Platform, include the `-list-boards` option in the `aoc` command.

- At a command prompt, invoke the `aoc -list-boards` command.

The Intel FPGA SDK for OpenCL Offline Compiler generates an output that resembles the following:

```
Board list:
<board_name_1>
<board_name_2>
...
```

Where `<board_name_N>` is the board name you use in your `aoc` command to target a specific FPGA board.

- When multiple Custom Platforms are installed, to list FPGA boards available in a specific Custom Platform, include the `-board-package=<custom_platform_path>` option in the `aoc` command. At the command prompt, invoke the command as:

```
aoc -board-package=<custom_platform_path> -list-boards
```

The Intel FPGA SDK for OpenCL Offline Compiler lists available boards within the Custom Platform.

- To list Custom Platforms available in the system, include the `-list-board-packages` option in the `aoc` command. At a command prompt, invoke the `aoc -list-board-packages` command.

*Note:* Starting from the 20.3 release, Intel has removed the support for Windows and Linux OpenCL BSPs. Use version 20.2 or older OpenCL BSPs available at [FPGA Software Download Center](#) as a reference. If you want to migrate your OpenCL BSP to a newer version, follow the recommended steps provided in the Reference Platform Porting Guides available under [Intel FPGA SDK for OpenCL documentation](#).

The Intel FPGA SDK for OpenCL Offline Compiler generates an output that resembles the following:

```
Installed board packages:
<board_package_1>
<board_package_2>
...
```

Where `<board_package_N>` is the board package of the Custom Platform installed in your system or shipped within the Intel FPGA SDK for OpenCL.

### 3.6. Displaying the Compilation Environment of an OpenCL Binary (env)

To display the Intel FPGA SDK for OpenCL Offline Compiler's input arguments and the environment for a compiled OpenCL design, invoke the `env` utility command.

- At the command prompt, invoke the `aocl env <object_file_name>` or the `aocl env <executable_file_name>` command, where `<object_file_name>` is the name of the `.aoco` file of your OpenCL kernel, and the `<executable_file_name>` is the name of the `.aocx` file of your kernel.

## 4. Managing an FPGA Board

---

The Intel FPGA SDK for OpenCL includes utility commands you can invoke to install, uninstall, diagnose, and program your FPGA board.

*Note:* Starting from the 20.3 release, Intel has removed the support for Windows and Linux OpenCL BSPs. Use version 20.2 or older OpenCL BSPs available at [FPGA Software Download Center](#) as a reference. If you want to migrate your OpenCL BSP to a newer version, follow the recommended steps provided in the Reference Platform Porting Guides available under [Intel FPGA SDK for OpenCL documentation](#).

You can install multiple Custom Platforms simultaneously on the same system with the `aocl install` utility. The Custom Platform subdirectory contains the `board_env.xml` file.

In a system with multiple Custom Platforms, ensure that the host program uses the FPGA Client Driver (FCD), formerly Altera Client Driver (ACD), to discover the boards rather than linking to the Custom Platform memory-mapped device (MMD) libraries directly.

FCD is set up for you when you run the `aocl install` utility. The installed BSP is registered on the system so the runtime and SDK utilities can find the necessary BSP files.

*Important:* Do not move a BSP to a different directory after you install it. To move a BSP:

1. Uninstall the BSP from its current location with the `aocl uninstall` utility.
2. Change the BSP directory.
3. Reinstall the BSP in the new location with the `aocl install` utility.

[Installing an FPGA Board \(install\)](#) on page 20

[Uninstalling an FPGA Board \(uninstall\)](#) on page 22

[Querying the Device Name of Your FPGA Board \(diagnose\)](#) on page 22

[Running a Board Diagnostic Test \(diagnose <device\\_name>\)](#) on page 23

[Programming the FPGA Offline or without a Host \(program <device\\_name>\)](#) on page 23

[Programming the Flash Memory \(flash <device\\_name>\)](#) on page 23

### Related Information

- [Installing an FPGA Board \(install\)](#) on page 20
- [Linking Your Host Application to the Khronos ICD Loader Library](#) on page 107

## 4.1. Installing an FPGA Board (install)

Before creating an OpenCL application for an FPGA accelerator board or SoC device, you must first download and install the Intel Reference Platform or the Custom Platform from your board vendor. Most Custom Platform installers require administrator privileges.

To install your board into the host system, invoke the `aocl install <path_to_customplatform>` utility command.

The steps below outline the board installation procedure. Some Custom Platforms require additional installation tasks. Consult your board vendor's documentation for further information on board installation.

**Attention:**

- If you are installing the Intel Arria® 10 SoC Development Kit for use with the Intel Arria 10 SoC Development Kit Reference Platform (a10soc), refer to the *Installing the Intel Arria 10 Development Kit* in [AN 807: Configuring the Intel Arria 10 GX FPGA Development Kit for the Intel FPGA SDK for OpenCL](#) for more information.
- If you want to use Intel FPGA SDK for OpenCL with the Intel Arria 10 GX FPGA Development Kit, refer to the Application Note [AN 807: Configuring the Intel Arria 10 GX FPGA Development Kit for the Intel FPGA SDK for OpenCL](#) for more information.

1. Follow your board vendor's instructions to connect the FPGA board to your system.
2. Download the Custom Platform for your FPGA board from your board vendor's website.

*Note:* Starting from the 20.3 release, Intel has removed the support for Windows OpenCL BSPs. Use 20.2 or older OpenCL BSPs available at [FPGA Software Download Center](#) as a reference. If you want to migrate your OpenCL BSP to a newer version, follow the recommended steps provided in the Reference Platform Porting Guides available under [Intel FPGA SDK for OpenCL documentation](#).

3. Install the Custom Platform in a folder that you own (that is, not a system folder). You can install multiple Custom Platforms simultaneously on the same system using the SDK utilities, such as `aocl diagnose` with multiple Custom Platforms. The Custom Platform subdirectory contains the `board_env.xml` file.

In a system with multiple Custom Platforms, ensure that the host program uses the FPGA Client Driver (FCD) to discover the boards rather than linking to the Custom Platforms' memory-mapped device (MMD) libraries directly. If you have not set up FCD correctly for Custom Platform, FCD finds all the installed boards at runtime.

4. Set the `QUARTUS_ROOTDIR_OVERRIDE` user environment variable to point to the Intel Quartus Prime Pro Edition software installation directory.
5. Add the paths to the Custom Platform libraries (for example, path to the MMD library of the board support package resembles `<path_to_customplatform>/windows64/bin`) to the `PATH` (Windows) or `LD_LIBRARY_PATH` (Linux) environment variable setting.

The *Intel FPGA SDK for OpenCL Pro Edition Getting Started Guide* contains more information on the `init_openc1` script. For information on setting user environment variables and running the `init_openc1` script, refer to the *Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables* section.

6. Invoke the command `aocl install <path_to_customplatform>` at a command prompt.

Invoking `aocl install <path_to_customplatform>` installs both the FCD and a board driver that allows communication between host applications and hardware kernel programs.

**Remember:**

- You need administrative rights to install a board. To run a Windows command prompt as an administrator, click **Start > All Programs > Accessories**. Under **Accessories**, right click **Command Prompt**. In the right-click menu, click **Run as Administrator**.

On Windows 8.1 or Windows 10 systems, you might also need to disable signed driver verification. For details, see the following articles:

- Windows 8: <https://www.intel.com/content/www/us/en/support/programmable/articles/000085519.html>
- Windows 10: <https://www.intel.com/content/www/us/en/support/programmable/articles/000086373.html>
- If the system already has the driver installed and you need to install FCD without the administrative rights, you can invoke the `aocl install` command with the flag `-fcd-only` as shown below and follow the prompt for FCD installation:

```
aocl install <path_to_customplatform> -fcd-only
```

7. Query a list of FPGA devices installed in your machine by invoking the `aocl diagnose` command.  
The software generates an output that includes the `<device_name>`, which is an acl number that ranges from `acl0` to `acl127`.  
**Attention:** For possible errors after implementing the `aocl diagnose` utility, refer to [Possible Errors After Running the diagnose Utility](#) section in the *Intel Arria 10 GX FPGA Development Kit Reference Platform Porting Guide*. For more information on querying the `<device_name>` of your accelerator board, refer to the *Querying the Device Name of Your FPGA Board* section.
8. Verify the successful installation of the FPGA board by invoking the command `aocl diagnose <device_name>` to run any board vendor-recommended diagnostic test.

#### Related Information

- [Querying the Device Name of Your FPGA Board \(diagnose\)](#) on page 22
- [Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables \(Windows\)](#)
- [Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables \(Linux\)](#)
- [Intel FPGA SDK for OpenCL FPGA Platforms page](#)

## 4.2. Uninstalling an FPGA Board (uninstall)

To uninstall an FPGA board, invoke the `uninstall` utility command, uninstall the Custom Platform, and unset the relevant environment variables.

To uninstall your FPGA board, perform the following tasks:

1. Disconnect the board from your machine by following the instructions provided by your board vendor.
2. Invoke the `aocl uninstall <path_to_customplatform>` utility command to remove the current host computer drivers (for example, PCIe® drivers). The Intel FPGA SDK for OpenCL uses these drivers to communicate with the FPGA board.

*Remember:*

- You need `root` privileges to uninstall the Custom Platform. If you want to keep the driver while removing the installed FCD, you can invoke the `aocl uninstall` command with the flag `-fcd-only` as shown below and follow the prompt for FCD uninstall:

```
aocl uninstall <path_to_customplatform> -fcd-only
```

- For Linux systems, if you had installed the FCD to a specific directory, then prior to uninstalling, you need to ensure that you have set an environment variable `ACL_BOARD_VENDOR_PATH` that points to that specific FCD installation directory.
3. Uninstall the Custom Platform.
  4. Unset the `LD_LIBRARY_PATH` (for Linux) or `PATH` (for Windows) environment variable.

## 4.3. Querying the Device Name of Your FPGA Board (diagnose)

Some OpenCL software utility commands require you to specify the device name (`<device_name>`). The `<device_name>` refers to the `acl` number (e.g., `acl0` to `acl127`) that corresponds to the FPGA device. When you query a list of accelerator boards, the OpenCL software produces a list of installed devices on your machine in the order of their device names.

- To query a list of installed devices on your machine, type `aocl diagnose` at a command prompt.

The software generates an output that resembles the example shown below:

```
aocl diagnose: Running diagnostic from <board_package_path>/<board_name>/<platform>/libexec
```

```
Verified that the kernel mode driver is installed on the host machine.
```

```
Using board package from vendor: <board_vendor_name>
```

```
Querying information for all supported devices that are installed on the host machine ...
```

device_name	Status	Information
acl0	Passed	<descriptive_board_name> PCIe dev_id = <device_ID>, bus:slot.func = 02:00.00, at Gen 2 with 8 lanes. FPGA temperature = 43.0 degrees C.
acl1	Passed	<descriptive_board_name> PCIe dev_id = <device_ID>, bus:slot.func = 03:00.00,

```
at Gen 2 with 8 lanes.  
FPGA temperature = 35.0 degrees C.  
  
Found 2 active device(s) installed on the host machine, to perform a full  
diagnostic on a specific device, please run aocl diagnose <device_name>  
  
DIAGNOSTIC_PASSED
```

#### Related Information

Probing the OpenCL FPGA Devices on page 112

### 4.4. Running a Board Diagnostic Test (diagnose <device\_name>)

To perform a detailed diagnosis on a specific FPGA board, include <device\_name> as an argument of the `diagnose` utility command.

- At a command prompt, invoke the `aocl diagnose <device_name>` command, where <device\_name> is the acl number (for example, `acl0` to `acl127`) that corresponds to your FPGA device.

You can identify the <device\_name> when you query the list of installed boards in your system.

Consult your board vendor's documentation for more board-specific information on using the `diagnose` utility command to run diagnostic tests on multiple FPGA boards.

### 4.5. Programming the FPGA Offline or without a Host (program <device\_name>)

To program an FPGA device offline or without a host, invoke the `program` utility command.

- At a command prompt, invoke the `aocl program <device_name> <your_kernel_filename>.aocx` command

where:

<device\_name> refers to the acl number (for example, `acl0` to `acl127`) that corresponds to your FPGA device, and

<your\_kernel\_filename>.aocx is the executable file you use to program the hardware.

**Note:** To program an SoC, specify the full path of the device when invoking the `program` utility command. For example, `aocl program /dev/<device_name> <your_kernel_filename>.aocx`.

### 4.6. Programming the Flash Memory (flash <device\_name>)

If supported by a Custom Platform, invoke the `flash` utility command to initialize the FPGA with a specified startup configuration.

**Note:** For example instructions on programming the micro SD flash card on an SoC board such as the Intel Arria 10 SoC Development Kit, refer to the *Building the SD Card Image* section of the *Intel FPGA SDK for OpenCL Intel Arria 10 SoC Development Kit Reference Platform Porting Guide*.

- At a command prompt, invoke the `aocl flash <device_name> <your_kernel_filename>.aocx` command  
where:  
`<device_name>` refers to the acl number (for example, acl0 to acl127) that corresponds to your FPGA device, and  
`<your_kernel_filename>.aocx` is the executable file you use to program the hardware.

#### Related Information

[Building the SD Card Image for the Intel Arria 10 SoC Development Kit](#)



## 5. Structuring Your OpenCL Kernel

---

Intel offers recommendations on how to structure your OpenCL kernel code. Consider implementing these programming recommendations when you create a kernel or modify a kernel written originally to target another architecture.

[Guidelines for Naming the Kernel](#) on page 25

[Programming Strategies for Optimizing Data Processing Efficiency](#) on page 26

[Programming Strategies for Optimizing Pointer-to-Local Memory Size](#) on page 42

[Implementing the Intel FPGA SDK for OpenCL Channels Extension](#) on page 43

[Implementing OpenCL Pipes](#) on page 60

[Implementing Arbitrary Precision Integers](#) on page 82

[Using Predefined Preprocessor Macros in Conditional Compilation](#) on page 83

[Declaring \\_\\_constant Address Space Qualifiers](#) on page 83

[Including Structure Data Types as Arguments in OpenCL Kernels](#) on page 84

[Inferring a Register](#) on page 87

[Enabling Double Precision Floating-Point Operations](#) on page 89

[Single-Cycle Floating-Point Accumulator for Single Work-Item Kernels](#) on page 89

[Integer Promotion Rules](#) on page 91

### 5.1. Guidelines for Naming the Kernel

Intel recommends that you include only alphanumeric characters in your file names.

- Begin a file name with an alphanumeric character.  
If the file name of your OpenCL application begins with a nonalphanumeric character, compilation fails with the following error message:

```
Error: Quartus compilation FAILED
See quartus_sh_compile.log for the output log.
```

- Ensure that the kernel file name only contains alphanumeric character, dash, underscore, or dot.  
The Intel FPGA SDK for OpenCL application only accepts file name containing alphanumeric character, dash, underscore, or dot. File name having other characters is treated as an invalid file name and triggers the following compilation error message:

```
aoc foo\*1.cl
Error: File: foo*1.cl contains invalid characters.
Ensure the file name only contains alphanumeric characters, dash, underscore
or dot.
```

- For Windows systems, ensure that the combined length of the kernel file name and its file path does not exceed 260 characters.  
64-bit Windows 7 and Windows 8.1 have a 260-character limit on the length of a file path. If the combined length of the kernel file name and its file path exceeds 260 characters, the offline compiler generates the following error message:

```
The filename or extension is too long.
The system cannot find the path specified.
```

In addition to the compiler error message, the following error message appears in the `<your_kernel_filename>/quartus_sh_compile.log` file:

```
Error: Can't copy <file_type> files: Can't open
<your_kernel_filename> for write: No such file or directory
```

For Windows 10, you can remove the 260-character limit. For more information, see your Windows 10 documentation.

- Do not name your `.cl` OpenCL kernel source file "kernel", "Verilog", or "VHDL" as they are reserved keywords.  
Naming the source file `kernel.cl`, `Verilog.cl`, or `VHDL.cl` causes the offline compiler to generate intermediate design files that have the same names as certain internal files, which leads to a compilation error.

## 5.2. Programming Strategies for Optimizing Data Processing Efficiency

Optimize the data processing efficiency of your kernel by implementing strategies such as unrolling loops, setting work-group sizes, and specifying compute units and work-items.

[Unrolling a Loop \(unroll Pragma\)](#) on page 27

[Disabling Pipelining of a Loop \(disable\\_loop\\_pipelining Pragma\)](#) on page 28

[Coalescing Nested Loops](#) on page 29

- [Fusing Adjacent Loops \(`loop\_fuse` Pragma\) on page 30](#)
- [Marking Loops to Prevent Automatic Fusion \(`nofusion` Pragma\) on page 32](#)
- [Specifying a Loop Initiation interval \(II\) on page 33](#)
- [Loop Concurrency \(`max\_concurrency` Pragma\) on page 34](#)
- [Loop Speculation \(`speculated\_iterations` Pragma\) on page 35](#)
- [Loop Interleaving Control \(`max\_interleaving` Pragma\) on page 35](#)
- [Floating Point Optimizations \(`fp\_contract` and `fp\_reassociate` Pragma\) on page 38](#)
- [Specifying Work-Group Sizes on page 39](#)
- [Specifying Number of Compute Units on page 40](#)
- [Specifying Number of SIMD Work-Items on page 40](#)
- [Specifying the `private\_copies` Memory Attribute on page 41](#)
- [Specifying the `use\_stall\_enable\_clusters` Cluster-control Attribute on page 42](#)

### 5.2.1. Unrolling a Loop (`unroll` Pragma)

Loop unrolling involves replicating a loop body multiple times, and reducing the trip count of a loop. Unroll loops to reduce or eliminate loop control overhead on the FPGA. In cases where there are no loop-carried dependencies and the offline compiler can perform loop iterations in parallel, unrolling loops can also reduce latency and overhead on the FPGA.

To direct the offline compiler to unroll a loop, or explicitly not to unroll a loop, insert an `unroll` kernel pragma in the kernel code preceding a loop you want to unroll.

*Note:* Unrolling of nested loops with large bounds might generate huge number of instructions, which can lead to very long compile times.

- Provide an unroll factor whenever possible. To specify an unroll factor  $N$ , insert the `#pragma unroll <N>` directive before a loop in your kernel code.

The offline compiler attempts to unroll the loop at most  $<N>$  times.

Consider the code fragment below. By assigning a value of 2 as the unroll factor, you direct the offline compiler to unroll the loop twice.

```
#pragma unroll 2
for(size_t k = 0; k < 4; k++)
{
    mac += data_in[(gid * 4) + k] * coeff[k];
}
```

- To prevent a loop from unrolling, specify an unroll factor of 1 (that is, `#pragma unroll 1`).
- To unroll a loop fully, you may omit the unroll factor by simply inserting the `#pragma unroll` directive before a loop in your kernel code.

The offline compiler attempts to unroll the loop fully if it understands the trip count. The offline compiler issues a warning if it cannot execute the unroll request.

Consider the following code fragment where the unroll factor is 2:

```
float data[N];
#pragma unroll 2
for (int i = 0; i < N; i++){
    data[i] = function(i, a);
}
```

The offline compiler partially unrolls the loop as shown in the following code fragment:

```
float data[N];
for (int i = 0; i < N; i += 2){
    data[i + 0] = function(i + 0, a);
    if (i + 1 < N){
        data[i + 1] = function(i + 1, a);
    }
}
```

### 5.2.2. Disabling Pipelining of a Loop (`disable_loop_pipelining` Pragma)

If loop-carried dependencies result in an initiation interval (II) that is equal or close to the latency of a given iteration (effectively inducing serial execution of the pipelined loop), disable pipelining of the loop to generate a simpler datapath and reduce area utilization.

Use the `disable_loop_pipelining` pragma to direct the Intel FPGA SDK for OpenCL Offline Compiler to disable pipelining of a loop. This pragma applies to single work-item kernels (that is, single-threaded kernels) in which loops are pipelined. Refer to the *Single Work-Item Kernel versus NDRange Kernel* section of the *Intel FPGA SDK for OpenCL Pro Edition Best Practices Guide* for information about loop pipelining and kernel properties that drive the offline compiler's decision about whether to treat a kernel as single-threaded.

Unless otherwise specified, the compiler always attempts to generate a pipelined loop datapath where possible. When generating a pipelined circuit, resources of the loop must be duplicated to execute multiple iterations simultaneously, leading to an increased silicon area utilization. In cases where loop pipelining does not result in an improvement in throughput, avoid the area overhead by applying the

`disable_loop_pipelining` pragma to the loop, as shown in the following code snippet. When you apply this pragma, the offline compiler generates a simple sequential loop datapath.

```
#pragma disable_loop_pipelining
for (int i = 1; i < N; i++) {
    int j = a[i-1];
    // Memory dependency induces a high-latency loop feedback path
    a[i] = foo(j)
}
```

In the above example, the offline compiler fails to schedule this loop with a small II due to memory dependency (as reported in the Details pane of the Loops Analysis section of the HTML report). In such cases, loop pipelining is unlikely to be beneficial.

#### Related Information

- [Single Work-Item Kernel versus NDRange Kernel](#)
- [Reviewing Your Kernel's report.html File](#)

### 5.2.3. Coalescing Nested Loops

Use the `loop_coalesce` pragma to direct the Intel FPGA SDK for OpenCL Offline Compiler to coalesce nested loops into a single loop without affecting the loop functionality. Coalescing loops can help reduce your kernel area usage by directing the compiler to reduce the overhead needed for loop control.

Coalescing nested loops also reduces the latency of the kernels, which can further reduce your kernel area usage. However, in some cases, coalescing loops might lengthen the critical loop initiation interval path, so coalescing loops might not be suitable for all kernels.

For NDRange kernels, the compiler automatically attempts to coalesce loops even if they are not annotated by the `loop_coalesce` pragma. Coalescing loops in NDRange kernels improves throughput as well as reducing kernel area usage. You can use the `loop_coalesce` pragma to prevent the automatic coalescing of loops in NDRange kernels.

To coalesce nested loops, specify the pragma as follows:

```
#pragma loop_coalesce <loop_nesting_level>
```

The `<loop_nesting_level>` parameter is optional and is an integer that specifies how many nested loop levels that you want the compiler to attempt to coalesce. If you do not specify the `<loop_nesting_level>` parameter, the compiler attempts to coalesce all of the nested loops.

For example, consider the following set of nested loops:

```
for (A)
  for (B)
    for (C)
      for (D)
    for (E)
```

If you place the pragma before loop (A), then the loop nesting level for these loops is defined as:

- Loop (A) has a loop nesting level of 1.
- Loop (B) has a loop nesting level of 2.
- Loop (C) has a loop nesting level of 3.
- Loop (D) has a loop nesting level of 4.
- Loop (E) has a loop nesting level of 3.

Depending on the loop nesting level that you specify, the compiler attempts to coalesce loops differently:

- If you specify `#pragma loop_coalesce 1` on loop (A), the compiler does not attempt to coalesce any of the nested loops.
- If you specify `#pragma loop_coalesce 2` on loop (A), the compiler attempts to coalesce loops (A) and (B).
- If you specify `#pragma loop_coalesce 3` on loop (A), the compiler attempts to coalesce loops (A), (B), (C), and (E).
- If you specify `#pragma loop_coalesce 4` on loop (A), the compiler attempts to coalesce all of the loops [loop (A) - loop (E)].

**Important:** If you specify `#pragma loop_coalesce 1` for a loop in an NDRange kernel, you prevent automatic loop coalescing for that loop.

### Example

The following simple example shows how the compiler coalesces two loops into a single loop.

Consider a simple nested loop written as follows:

```
#pragma loop_coalesce
for (int i = 0; i < N; i++)
    for (int j = 0; j < M; j++)
        sum[i][j] += i+j;
```

The compiler coalesces the two loops together so that they run as if they were a single loop written as follows:

```
int i = 0;
int j = 0;
while(i < N){

    sum[i][j] += i+j;
    j++;

    if (j == M){
        j = 0;
        i++;
    }
}
```

## 5.2.4. Fusing Adjacent Loops (`loop_fuse` Pragma)

Use the `loop_fuse` pragma to direct the Intel FPGA SDK for OpenCL Offline Compiler to fuse adjacent loops into a single loop without affecting either loop's functionality. The `loop_fuse` construct defines a region of code where the compiler always attempts to fuse adjacent loops when it is safe to do so.

Fusing adjacent loops can help reduce your kernel area use by reducing the overhead required for loop control and increasing the performance of your kernel by executing both original loops concurrently as one (fused) loop.

To specify a block of program code within which the compiler attempts to fuse loops, specify the pragma as follows:

```
#pragma loop_fuse [clause[[,]clause]...] new-line  
    structured_block
```

where `clause` is one of the following:

*depth(constant-integer-expression)*      If a *depth* clause is present, the *constant-integer-expression* clause parameter defines the number of nesting depths at which the fusion of adjacent loops is attempted. The *depth* clause extends the applicability of the `loop_fuse` construct to all loops nested in top-level loops contained in the construct at nesting depth less than or equal to the clause parameter, including loops that become adjacent as a result of fusion of their corresponding containing loops. In the absence of a *depth* clause, only loops at the top-level of the `loop_fuse` construct are attempted to be fused (that is, loops not contained in other loops defined within the construct). The *depth* clause with a parameter of 1 is equivalent to the absence of a *depth* clause.

*independent*      If an *independent* clause is present, adjacent loops that are fusion candidates within a `loop_fuse` construct are assumed to have no negative-distance data access dependencies. That is, for two adjacent loops considered for fusion, iterations of the logically-second loop does not access data elements produced in a later iteration of the logically-first loop. The *independent* clause overrides the offline compiler's static analysis during loop fusion safety analysis.

If a function call is present in a `loop_fuse` construct at any of the applicable nesting depths and inlining the function call materializes a loop, then the resulting loop is considered to be a candidate for fusion.

**Caution:** Default clauses are none, making the `loop_fuse` construct unable to introduce a functional error on its own. Introduction of an *independent* clause is a guarantee from you that bypasses a respective aspect of the compiler's safety analysis and might lead to functional errors.

### Nested Depth Clauses

In programs where `loop_fuse` constructs are nested and their implied sets of fusion candidates overlap, the overall set of fusion candidates comprises a union of all loops covered by the distinct `loop_fuse` regions. The `loop_fuse` attribute clauses apply only to the fusion candidates implied by the directive to which the clauses apply.

```
#pragma loop_fuse depth(2) independent
{
  L1: for(...) {}
  L2: for(...) {
    #pragma loop_fuse depth(2)
    {
      L3: for(...) {}
      L4: for(...) {
        L5: for(...) {}
        L6: for(...) {}
      }
    }
  }
}
```

In this example, loops L1, L2, L3, L4, L5, and L6 are considered for fusion and loops L1, L2, L3, L4 are considered for fusion overriding the compiler's dependence analysis.

### Related Information

- [Reviewing Your Kernel's report.html File](#)
- [Loop Fusion](#)

## 5.2.5. Marking Loops to Prevent Automatic Fusion (nofusion Pragma)

Use the `nofusion` pragma to direct the Intel FPGA SDK for OpenCL Offline Compiler to avoid fusing the annotated loop with any of the adjacent loops.

To specify a loop to not be eligible for fusion, specify the `nofusion` pragma as follows:

```
#pragma nofusion
for (int i = 0; i < N; ++i) {
  ...
}
```

In the following example, the compiler does not apply the loop fusion transformation to loops L1 and L2.

```
#pragma nofusion
L1: for (int j=0; j < N; ++j) {
  data[j] += Q;
}
L2: for (int i = 0; i < N; ++i) {
  output[i] = Q * data[i];
}
```

### Related Information

- [Reviewing Your Kernel's report.html File](#)
- [Loop Fusion](#)



## 5.2.6. Specifying a Loop Initiation Interval (II)

The initiation interval, or II, is the number of clock cycles between the launch of successive loop iterations. Use the `ii` pragma to direct the Intel FPGA SDK for OpenCL Offline Compiler to attempt to set the II for the loop that follows the pragma declaration. If the offline compiler cannot achieve the specified II for the loop, then the compilation errors out.

The `ii` pragma applies to single work-item kernels (that is, single-threaded kernels) in which loops are pipelined. Refer to the *Single Work-Item Kernel versus NDRange Kernel* section of the *Intel FPGA SDK for OpenCL Best Practices Guide* for information on loop pipelining, and on kernel properties that drive the offline compiler's decision on whether to treat a kernel as single-threaded.

The higher the II value, the longer the wait before the subsequent loop iteration starts executing. Refer to the *Reviewing Your Kernel's report.html File* section of the *Intel FPGA SDK for OpenCL Best Practices Guide* for information on II, and on the compiler reports that provide you with details on the performance implications of II on a specific loop.

For some loops in your kernel, specifying a higher II value with the `ii` pragma than the value the compiler chooses by default can increase the maximum operating frequency ( $f_{MAX}$ ) of your kernel without a decrease in throughput.

A loop is a good candidate to have the `ii` pragma applied to it if the loop meets the following conditions:

- The loop is pipelined because the kernel is single-threaded.
- The loop is not critical to the throughput of your kernel.
- The running time of the loop is small compared to other loops it might contain.

To specify a loop initiation interval for a loop, specify the pragma before the loop as follows:

```
#pragma ii <desired_initiation_interval>
```

The `<desired_initiation_interval>` parameter is required and is an integer that specifies the number of clock cycles to wait between the beginning of execution of successive loop iterations.

### Example

Consider a case where your kernel has two distinct, pipelineable loops: a short-running initialization loop that has a loop-carried dependence and a long-running loop that does the bulk of your processing. In this case, the compiler does not know that the initialization loop has a much smaller impact on the overall throughput of your design. If possible, the compiler attempts to pipeline both loops with an II of 1.

Because the initialization loop has a loop-carried dependence, it has a feedback path in the generated hardware. To achieve an II with such a feedback path, some clock frequency might be sacrificed. Depending on the feedback path in the main loop, the rest of your design can run at a higher operating frequency.

If you specify `#pragma ii 2` on the initialization loop, you tell the compiler that it can be less aggressive in optimizing II for this loop. Less aggressive optimization permits the compiler to pipeline the path limiting the  $f_{max}$  and allow your overall kernel design to achieve a higher  $f_{max}$ .

The initialization loop takes longer to run with its new II. However, the decrease in the running time of the long-running loop due to higher  $f_{\max}$  compensates for the increased length in running time of the initialization loop.

#### Related Information

- [Single Work-Item Kernel versus NDRange Kernel](#)
- [Reviewing Your Kernel's report.html File](#)

### 5.2.7. Loop Concurrency (`max_concurrency` Pragma)

You can use the `max_concurrency` pragma to limit the concurrency of a loop in your component.

The concurrency of a loop is how many iterations of that loop can be in progress at one time. By default, the Intel FPGA SDK for OpenCL tries to maximize the concurrency of loops so that your component runs at peak throughput.

The `max_concurrency` pragma applies to single work-item kernels (that is, single-threaded kernels) in which loops are pipelined. Refer to the *Single Work-Item Kernel versus NDRange Kernel* section of the *Intel FPGA SDK for OpenCL Pro Edition Best Practices Guide* for information on loop pipelining, and on kernel properties that drive the offline compiler's decision on whether to treat a kernel as single-threaded.

The `max_concurrency` pragma enables you to control the on-chip memory resources required to pipeline your loop. To achieve simultaneous execution of loop iterations, the offline compiler must create copies of any memory that is private to a single iteration. These copies are called *private copies*. The greater the permitted concurrency, the more private copies the compiler must create.

The kernel's HTML report (`report.html`) provides the following information pertaining to loop concurrency:

- Maximum concurrency that the offline compiler has chosen

This information is available in the Loop Analysis report and Kernel Memory viewer:

- In the Loop Analysis report, a message in the Details pane reports as the maximum number of simultaneous executions has been limited to  $N$ .

*Note:* The value of unsigned  $N$  can be greater than or equal to zero. A value of  $N = 0$  indicates unlimited concurrency.

- In the Kernel Memory Viewer, the bank view of your local memory graphically shows the number of private copies.

- Impact to memory usage

This information is available in the Area Analysis report. A message in the Details pane reports that the offline compiler has created  $N$  independent copies of the memory to enable simultaneous execution of  $N$  loop iterations.

If you want to exchange some performance for physical memory savings, apply `#pragma max_concurrency <N>` to the loop, as shown below. When you apply this pragma, the offline compiler limits the number of simultaneously-executed loop iterations to  $N$ . The number of private copies of loop memories is also reduced to  $N$ .

```
#pragma max_concurrency 1
for (int i = 0; i < N; i++) {
    int arr[M];
    // Doing work on arr
}
```

You can also control the number of private copies (created for a local memory and accessed within a loop) by using `__attribute__((private_copies(N)))`. Refer to [Memory Attributes for Configuring Kernel Memory Systems](#) on page 199 for more details about the attribute. If a local memory with `__attribute__((private_copies(N)))` is accessed with a loop that has `#pragma max_concurrency M`, the offline compiler limits the number of simultaneously-executed loop iterations to  $\min(M, N)$ .

#### Related Information

- [Single Work-Item Kernel versus NDRange Kernel](#)
- [Reviewing Your Kernel's report.html File](#)

### 5.2.8. Loop Speculation (`speculated_iterations` Pragma)

Use the `speculated_iterations` pragma to direct the Intel FPGA SDK for OpenCL Offline Compiler to improve the performance of pipelined loops.

The `speculated_iterations` pragma is applied to loops and hence, it must appear directly before the loop (the same place as other loop pragmas) as shown in the following:

```
#pragma speculated_iterations k // where k >= 0
```

The Intel FPGA SDK for OpenCL Offline Compiler generates hardware to run  $k$  extra iterations of the loop while ensuring that the extra iterations do not affect anything. This allows either reducing the II of the loop or increasing the  $f_{\max}$ . The deciding factor is how quickly the exit condition of the loop is calculated. If the calculation takes many cycles, it is better to have `speculated_iterations` larger.

*Remember:* Extra iterations increase the time before the next invocation of the loop can begin. This may be a factor if the actual number of iterations of the loop is very small (less than 5 to 10 or similar). In this case, specify the `#pragma speculated_iterations` value as 0 to allow subsequent loop iterations to start immediately but at the cost of a larger II to allow more time to evaluate the exit condition. Refer to the Loop Analysis report in the HLD report to identify whether the exit condition is a bottleneck for II.

### 5.2.9. Loop Interleaving Control (`max_interleaving` Pragma)

The Intel FPGA SDK for OpenCL Offline Compiler attempts to maximize the throughput and hardware resource occupancy of pipelined inner loops in a loop nest by issuing new inner loop iterations as frequently as possible (minimizing the loop initiation

interval). When the compiler cannot achieve a loop II of 1 for an inner loop, the compiler configures the loop nest to interleave iterations of one invocation of the inner loop with iterations of other invocations of the inner loop.

As an example, consider the loop nest in the following code snippet:

```
// Loop j is pipelined with ii=1
for (int j = 0; j < M; j++) {
    int a[N];

    // Loop i is pipelined with ii=2
    for (int i = 1; i < N; i++) {
        a[i] = foo(i)
    }
}
```

In this example, the inner *i* loop is pipelined with a loop II of 2. Under normal pipelining, this means that the inner loop hardware only achieves 50% utilization since one *i* iteration is initiated every other cycle. To take advantage of these idle cycles, the compiler interleaves a second invocation of the *i* loop from the next iteration of the outer *j* loop. Here, a loop invocation means to start pipelined execution of a loop body. In this example, since the *i* loop resides inside the *j* loop, and the *j* loop has a trip count of *M*, the *i* loop is invoked *M* times. Since the *j* loop is an outermost loop, it is invoked once. The following table illustrates the difference between normal pipelined execution of the *i* loop and interleaved execution for this example where *N*=5:

**Table 1. Difference Between Normal Pipelined Execution and Interleaved Execution**

Cycle	Pipelined	Interleaved
0	(0,0)	(0,0)
1	---	(1,0)
2	(0,1)	(0,1)
3	---	(1,1)
4	(0,2)	(0,2)
5	---	(1,2)
6	(0,3)	(0,3)
7	---	(1,3)
8	(0,4)	(0,4)
9	---	(1,4)
10	(1,0)	(2,0)
11	---	(3,0)
12	(1,1)	(2,1)
13	---	(3,1)
14	(1,2)	(2,2)
15	---	(3,2)
16	(1,3)	(2,3)
<i>continued...</i>		

Cycle	Pipelined	Interleaved
17	---	(3,3)
18	(1,4)	(2,4)
19	---	(3,4)

The table shows the values  $(j, i)$  for each inner loop iteration that is initiated at each cycle. At cycle 0, both modes of execution initiate the  $(0,0)$ <sup>th</sup> iteration of the  $i$  loop. Under normal pipelined execution, no  $i$  loop iteration is initiated at cycle 1. Under interleaved execution, the  $(1,0)$ <sup>th</sup> iteration of the innermost loop, that is, the first iteration of the next  $(j=1)$  invocation of the  $i$  loop is initiated. By cycle 10, interleaved execution has initiated all of the iterations of both the  $j=0$  invocation of the  $i$  loop and the  $j=1$  invocation of the  $i$  loop. This represents twice the efficiency of the normal pipelined execution.

In some cases, you may decide that the performance benefit from interleaving is not equal to the area cost associated with enabling interleaving. In these cases, you may want to limit or restrict the amount of interleaving to reduce FPGA area utilization. To limit the number of interleaved invocations of an inner loop that can be executed simultaneously, annotate the inner loop with the `max_interleaving` pragma. The annotated loop must be contained inside another pipelined loop. The required parameter  $(n)$  specifies an upper bound on the degree of interleaving allowed, that is, how many invocations of the containing loop can execute the annotated loop at a given time.

Specify the `max_interleaving` pragma in one of the following ways:

- `#pragma max_interleaving 1`  
The compiler restricts the annotated (inner) loop to be invoked only once per outer loop iteration. That is, all iterations of the inner loop travels the pipeline before the next invocation of the inner loop can occur.
- `#pragma max_interleaving 0`  
The compiler allows the pipeline to contain a number of simultaneous invocations of the inner loop equal to the loop initiation interval (II) of the inner loop. For example, an inner loop with an II of 2 can have iterations from two invocations in the pipeline at a time. This behavior is the default behavior for the compiler if you do not specify the `max_interleaving` pragma.

In the following code snippet, the compiler restricts the pipelined execution of the  $i$  loop. A new invocation of the  $i$  loop corresponds only to the subsequent iteration of the  $j$  loop.

```
// Loop j is pipelined with ii=1
for (int j = 0; j < M; j++) {
    int a[N];
    // Loop i is pipelined with ii=2
    #pragma max_interleaving 1
    for (int i = 1; i < N; i++) {
        a[i] = foo(i)
    }
    ...
}
```

## 5.2.10. Floating Point Optimizations (`fp contract` and `fp reassociate Pragma`)

Use `fp contract` and `fp reassociate` pragmas to influence the intermediate rounding and conversions of floating-point operations and the ordering of arithmetic operations in your kernel at finer granularity than the Intel FPGA SDK for OpenCL Offline Compiler command options.

**Note:** You can place `fp contract` and `fp reassociate` pragma statements inside any compound statement (brace-enclosed sequences of statements), outside of all functions (file scope), or at the start of a function within the curly braces.

For example:

```
{
  #pragma clang fp reassociate(on)
  T=(1.0f-a)*(1.0f-b)*Ti0j0+a*(1.0f-b)*Ti1j0+(1.0f-a)*b*Ti0j1+a*b*Ti1j1;
}
```

### `fp contract` Pragma

The `fp contract` pragma controls whether the compiler can skip intermediate rounding and conversions mainly between double-precision arithmetic operations. If multiple occurrences of this pragma affect the same scope of your code, the pragma with the narrowest scope takes precedence.

This pragma has the following syntax:

```
#pragma clang fp contract(state)
```

where, the `state` parameter can be one of the following values:

**Table 2.** `state` Parameter Values

Value	Description
<code>off</code>	Turns off any permissions to fuse instructions into FMAs.
<code>fast</code>	Allows the fusing of multiply and add instructions into an FMA, but it might violate the language standard.

### `fp reassociate` Pragma

The `fp reassociate` pragma controls the relaxing of the order of floating point arithmetic operations within the code block that this pragma is applied to. With reordering, the compiler can optimize the hardware structure, which improves the performance of your kernel. If multiple occurrences of this pragma affect the same scope of your code, the pragma with the narrowest scope takes precedence.

This pragma has the following syntax:

```
#pragma clang fp reassociate(state)
```

where, the *state* parameter can be one of the following values:

**Table 3. state Parameter Values**

Value	Description
on	Enables the compiler to reorder floating-point operations to improve performance and on-chip area.
off	Prohibits the compiler from reordering floating-point operations to improve performance and on-chip area.

### 5.2.11. Specifying Work-Group Sizes

Specify a maximum or required work-group size whenever possible. The Intel FPGA SDK for OpenCL Offline Compiler relies on this specification to optimize hardware usage of the OpenCL kernel without involving excess logic.

If you do not specify a `max_work_group_size` or a `reqd_work_group_size` attribute in your kernel, the work-group size assumes a default value depending on compilation time and runtime constraints.

- If your kernel contains a barrier, the offline compiler sets a default maximum scalarized work-group size of 128 work-items.
- If your kernel does not query any OpenCL intrinsics that allow different threads to behave differently (that is, local or global thread IDs, or work-group ID), the offline compiler infers a single-threaded execution mode and sets the maximum work-group size to (1,1,1). In this case, the OpenCL runtime also enforces a global enqueue size of (1,1,1), and loop pipelining optimizations are enabled within the offline compiler.

To specify the work-group size, modify your kernel code in the following manner:

- To specify the maximum number of work-items that the offline compiler provisions for a work-group in a kernel, insert the `max_work_group_size(X, Y, Z)` attribute in your kernel source code.

For example:

```
__attribute__((max_work_group_size(512,1,1)))
__kernel void sum (__global const float * restrict a,
                  __global const float * restrict b,
                  __global float * restrict answer)
{
    size_t gid = get_global_id(0);
    answer[gid] = a[gid] + b[gid];
}
```

- To specify the required number of work-items that the offline compiler provisions for a work-group in a kernel, insert the `reqd_work_group_size(X, Y, Z)` attribute in your kernel source code.

For example:

```
__attribute__((reqd_work_group_size(64,1,1)))
__kernel void sum (__global const float * restrict a,
                  __global const float * restrict b,
                  __global float * restrict answer)
{
    size_t gid = get_global_id(0);
    answer[gid] = a[gid] + b[gid];
}
```

### 5.2.12. Specifying Number of Compute Units

To increase the data-processing efficiency of an OpenCL kernel, you can instruct the Intel FPGA SDK for OpenCL Offline Compiler to generate multiple kernel compute units. Each compute unit is capable of executing multiple work-groups simultaneously.

**Caution:** Multiplying the number of kernel compute units increases data throughput at the expense of FPGA resource consumption and global memory bandwidth contention between compute units.

- To specify the number of compute units for a kernel, insert the `num_compute_units(N)` attribute in the kernel source code.

For example, the code fragment below directs the offline compiler to instantiate two compute units in a kernel:

```
__attribute__((num_compute_units(2)))
__kernel void test(__global const float * restrict a,
                  __global const float * restrict b,
                  __global float * restrict answer)
{
    size_t gid = get_global_id(0);
    answer[gid] = a[gid] + b[gid];
}
```

The offline compiler dynamically distributes work-groups across the specified number of compute units.

*Note:* To identify the specific compute unit on which a work-item is executing, call the `get_compute_id()` intrinsic function. Refer to *Customization of Replicated Kernels Using the `get_compute_id()` Function* for more information.

#### Related Information

[Customization of Replicated Kernels Using the `get\_compute\_id\(\)` Function](#) on page 206

### 5.2.13. Specifying Number of SIMD Work-Items

You have the option to increase the data-processing efficiency of an OpenCL kernel by executing multiple work-items in a single instruction multiple data (SIMD) manner without manually vectorizing your kernel code. Specify the number of work-items within a work-group that the Intel FPGA SDK for OpenCL Offline Compiler should execute in an SIMD or vectorized manner.

**Important:** Introduce the `num_simd_work_items(N)` attribute in conjunction with the `reqd_work_group_size(X, Y, Z)` attribute. The `num_simd_work_items` attribute you specify must evenly divide the first argument that you specify to the `cl::reqd_work_group_size` attribute.

- To specify the number of SIMD work-items in a work-group, insert the `num_simd_work_item(N)` attribute in the kernel source code.

For example, the code fragment below assigns a fixed work-group size of 64 work-items to a kernel. It then consolidates the work-items within each work-group into four SIMD vector lanes:

```
__attribute__((num_simd_work_items(4)))
__attribute__((reqd_work_group_size(64,1,1)))
__kernel void test(__global const float * restrict a,
```



```

        __global const float * restrict b,
        __global float * restrict answer)
    {
        size_t gid = get_global_id(0);
        answer[gid] = a[gid] + b[gid];
    }

```

The offline compiler vectorizes the kernel datapath according to the value you specify for `num_simd_work_items` whenever possible.

### 5.2.14. Specifying the `private_copies` Memory Attribute

You have the option to apply the `private_copies` memory attribute to a variable declaration inside an OpenCL kernel as follows:

```
int __attribute__((private_copies(k))) local_A[M];
```

where,  $k$  is an unsigned integer. When this attribute is applied to a variable declared or accessed inside a pipelined loop, the Intel FPGA SDK for OpenCL Offline Compiler creates  $k$  independent copies of the memory implementing this variable. This allows up to  $k$  iterations of the pipelined loop to run in parallel, where each iteration accesses its own copy of the memory. If this attribute is not applied or if  $k$  is set to 0, then the compiler chooses an appropriate number of copies, up to a maximum of 16 to maximize throughput.

Consider the following example where the outer loop declares four local arrays:

```

for (int i = 0; i < N; i++) {
    int local_A[M];
    int local_B[M];
    int local_C[M];
    int local_D[M];

    // Step 1
    for (int j = 0; j < M; j++) {
        local_A[j] = initA();
    }

    // Step 2
    for (int j = 0; j < M; j++) {
        local_B[j] = initB(local_A[j]);
    }

    // Step 3
    for (int j = 0; j < M; j++) {
        local_C[j] = initC(local_B[j]);
    }

    // Step 4
    for (int j = 0; j < M; j++) {
        local_D[j] = initD(local_C[j]);
    }
}

```

In this example, the outer loop contains four steps, where each step corresponds to an inner loop. In Step 1, the first local array `local_A` is initialized. In Step 2, `local_A` is read from, but not written to. This is the last use of `local_A` in the outer loop. Similarly, `local_B` is first used in Step 2, where it is initialized. In Step 3, `local_B` is read from, but not written to, and this is the last use of `local_B`. Similarly, `local_C` is used only in Steps 3 and 4. The Intel FPGA SDK for OpenCL Offline Compiler privatizes each array by making 16 copies. These copies are enough to support concurrency of 16 on the outer loop. However, because the live ranges of these local arrays do not

span the entire outer loop, all 16 copies are not required to maximize throughput of the outer loop. This means that the amount of area consumed by making these copies is higher than necessary. In this case, applying the `private_copies` attribute to control the number of copies of these local arrays can reduce the area used while maintaining the throughput of the outer loop.

### 5.2.15. Specifying the `use_stall_enable_clusters` Cluster-control Attribute

You can apply the `use_stall_enable_clusters` cluster-control attribute to your OpenCL kernel to reduce the area of your kernel while possibly decreasing kernel  $f_{MAX}$  and throughput, as follows:

```
__kernel void __attribute__((use_stall_enable_clusters))
example(__global int * restrict input,
        __global int * restrict output, int size){
    for(int i = 0; i < size; ++i){
        output[i] = input[i];
    }
    ...
}
```

The Intel FPGA SDK for OpenCL Offline Compiler typically groups related operations into clusters. In several scenarios, the clusters are stall-free clusters. A stall-free cluster executes the operations without any stalls and contains a FIFO at the end of the cluster that holds the results if the cluster is stalled. This FIFO adds area and latency to the kernel, but might allow a higher  $f_{MAX}$  and increased throughput.

If you prefer lower FPGA area use and lower latency over higher throughput, use the `__attribute__((use_stall_enable_clusters))` attribute to bias the compiler to produce stall-enabled clusters. Stall-enabled clusters lack an exit FIFO to buffer all data in the event that the whole SFC is stalled, which reduces area and latency, but passes stall signals to the contained operations. Passing stall signals might reduce  $f_{MAX}$ .

Not all operations support stall, and these operations cannot be contained in a stall-enabled cluster. The compiler generates a warning if some operations cannot be placed into a stall-enabled cluster.

The compiler automatically uses stall-free clusters for kernels as they are generally more beneficial. This attribute requests the compiler to form stall-enabled clusters if possible.

*Intel Stratix 10 Restriction:* This attribute does not apply to designs that target Intel Stratix® 10 devices unless you specify the `-hyper-optimized-handshaking=off` option of the `aoc` command.

## 5.3. Programming Strategies for Optimizing Pointer-to-Local Memory Size

When using a pointer-to-local kernel argument to create a local memory allocation, the Intel FPGA SDK for OpenCL Offline Compiler must decide at compilation time (not runtime) on the size of the local memory system to build on the FPGA. This

specification allows the offline compiler to build the correctly sized local memory system for the pointer argument. If you do not specify a size, the offline compiler uses the default size.

- To specify a size other than the default of 16 kilobytes (kB), include the `local_mem_size(N)` attribute in the pointer declaration within your kernel source code, where the value  $N$  specifies the desired memory size in bytes. The  $N$  value should be a power of 2.

For example:

```
__kernel void myLocalMemoryPointer(  
    __local float * A,  
    __attribute__((local_mem_size(1024))) __local float * B,  
    __attribute__((local_mem_size(32768))) __local float * C)  
{  
    //statements  
}
```

In the `myLocalMemoryPointer` kernel, 16 kB of local memory (default) is allocated to pointer A, 1 kB is allocated to pointer B, and 32 kB is allocated to pointer C.

**Attention:** Instead of using pointer-to-local kernel arguments, Intel recommends that you define local memory systems within the kernel scope.

## 5.4. Implementing the Intel FPGA SDK for OpenCL Channels Extension

The Intel FPGA SDK for OpenCL channels extension provides a mechanism for passing data between kernels and synchronizing kernels with high efficiency and low latency.

**Attention:** If you want to take advantage of the capabilities of channels but have the ability to run your kernel program using other SDKs, implement your design using OpenCL pipes instead.

### Related Information

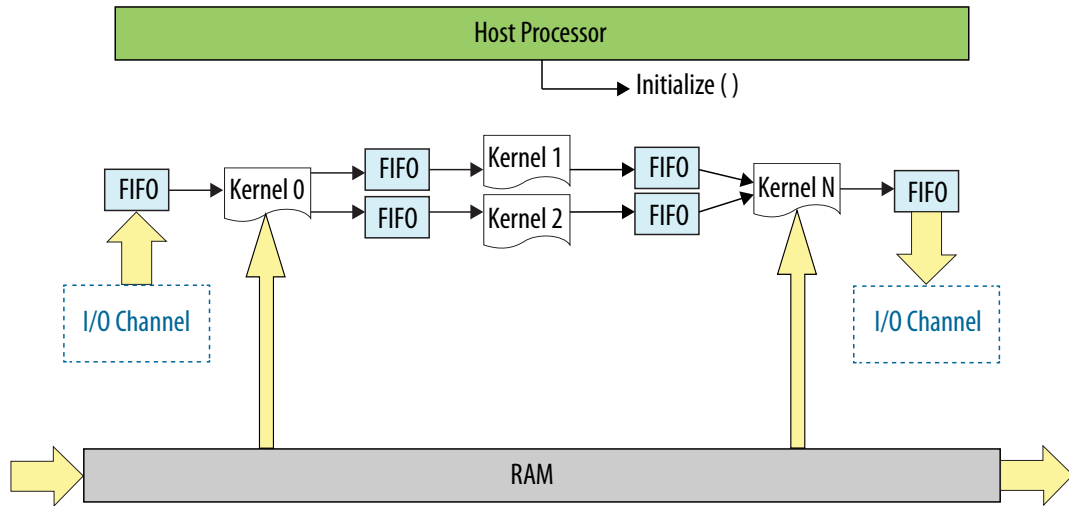
[Implementing OpenCL Pipes](#) on page 60

### 5.4.1. Overview of the Intel FPGA SDK for OpenCL Channels Extension

The Intel FPGA SDK for OpenCL channels extension allows kernels to communicate directly with each other through FIFO buffers.

Implementation of channels decouples data movement between concurrently executing kernels from the host processor.

Figure 6. Overview of Channels Implementation



### 5.4.2. Channel Data Behavior

Data written to a channel remains in a channel provided that the kernel program remains loaded on the FPGA device. In other words, data written to a channel persists across multiple work-groups and NDRange invocations. However, data is not persistent across multiple or different invocations of kernel programs that lead to FPGA device reprogramming.

Data in channels does not persist between context, program, device, kernel, or platform releases, even if the OpenCL implementation performs optimizations that avoid reprogramming operations on a device. For example, if you run a host program twice using the same .aocx file, or if a host program releases and reacquires a context, the data in the channel might or might not persist across the operation. FPGA device reset operations might happen behind the scenes on object releases that purge data in any channels

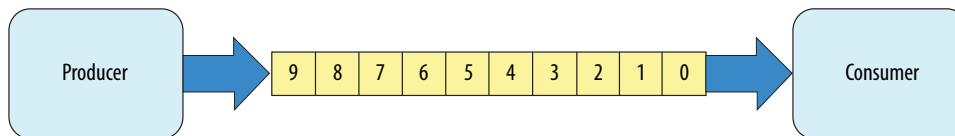
Consider the following code example:

```
channel int c0;

__kernel void producer() {
    for (int i = 0; i < 10; i++) {
        write_channel_intel (c0, i);
    }
}

__kernel void consumer (__global uint * restrict dst) {
    for (int i = 0; i < 5; i++) {
        dst[i] = read_channel_intel(c0);
    }
}
```

Figure 7. Channel Data FIFO Ordering



The kernel `producer` writes ten elements (`[0, 9]`) to the channel. The kernel consumer does not contain any work-item identifier queries; therefore, it receives an implicit `reqd_work_group_size` attribute of `(1,1,1)`. The implied `reqd_work_group_size(1,1,1)` attribute means that consumer must be launched as a single work-item kernel. In the example above, consumer reads five elements from the channel per invocation. During the first invocation, the kernel `consumer` reads values 0 to 4 from the channel. Because the data persists across NDRange invocations, the second time you execute the kernel `consumer`, it reads values 5 to 9.

For this example, to avoid a deadlock from occurring, you need to invoke the kernel `consumer` twice for every invocation of the kernel `producer`. If you call `consumer` less than twice, `producer` stalls because the channel becomes full. If you call `consumer` more than twice, `consumer` stalls because there is insufficient data in the channel.

### 5.4.3. Multiple Work-Item Ordering for Channels

The OpenCL specification does not define a work-item ordering. The Intel FPGA SDK for OpenCL enforces a work-item order to make it easier to rationalize channel read and write operations.

Multiple work-item accesses to a channel can be useful in some scenarios. For example, they are useful when data words in the channel are independent, or when the channel is implemented for control logic. The main concern regarding multiple work-item accesses to a channel is the order in which the kernel writes data to and reads data from the channel. If possible, the SDK's channels extension processes work-item read and write operations to the channel in a deterministic order. As such, the read and write operations remain consistent across kernel invocations.

#### Requirements for Deterministic Multiple Work-Item Ordering

To guarantee deterministic ordering, the SDK checks that a channel access is work-item invariant based on the following characteristics:

- All paths through the kernel must execute the channel access.
- If the first requirement is not satisfied, none of the branch conditions that reach the channel call should execute in a work-item-dependent manner.
- The kernel is not inferred as a single work-item kernel.

If the SDK cannot guarantee deterministic ordering of multiple work-item accesses to a channel, it warns you that the channels might not have well-defined ordering and therefore might exhibit nondeterministic execution. Primarily, the SDK fails to provide deterministic ordering if you have work-item-variant code on loop executions with channel calls, as illustrated below:

```
__kernel void ordering (__global int * restrict check,  
                      __global int * restrict data) {  
    int condition = check[get_global_id(0)];
```

```
if (condition) {
    for (int i = 0; i < N, i++) {
        process(data);
        write_channel_intel (req, data[i]);
    }
} else {
    process(data);
}
```

### 5.4.3.1. Work-Item Serial Execution of Channels

Work-item serial execution refers to an ordered execution behavior where work-item sequential IDs determine their execution order in the compute unit.

When you implement channels in a kernel, the Intel FPGA SDK for OpenCL Offline Compiler enforces that kernel behavior is equivalent to having at most one work-group in flight within the compute unit at a time. The compiler also ensures that the kernel executes channels in work-item serial execution, where the kernel executes work-items with smaller IDs first. A work-item has the identifier  $(x, y, z, \text{group})$ , where  $x, y, z$  are the local 3D identifiers, and  $\text{group}$  is the work-group identifier.

The work-item ID  $(x_0, y_0, z_0, \text{group}_0)$  is considered to be smaller than the ID  $(x_1, y_1, z_1, \text{group}_1)$  if one of the following conditions is true:

- $\text{group}_0 < \text{group}_1$
- $\text{group}_0 = \text{group}_1$  and  $z_0 < z_1$
- $\text{group}_0 = \text{group}_1$  and  $z_0 = z_1$  and  $y_0 < y_1$
- $\text{group}_0 = \text{group}_1$  and  $z_0 = z_1$  and  $y_0 = y_1$  and  $x_0 < x_1$

Work-items with incremental IDs execute in a sequential order. For example, the work-item with an ID  $(x_0, y_0, z_0, \text{group}_0)$  executes the write channel call first. Then, the work-item with an ID  $(x_1, y_0, z_0, \text{group}_0)$  executes the call, and so on. Defining this order ensures that the system is verifiable with external models.

### Channel Execution in Loop with Multiple Work-Items

When channels exist in the body of a loop with multiple work-items, as shown below, each loop iteration executes prior to subsequent iterations. This implies that loop iteration 0 of each work-item in a work-group executes before iteration 1 of each work-item in a work-group, and so on.

```
__kernel void ordering (__global int * data, int X) {
    int n = 0;
    while (n < X)
    {
        write_channel_intel (req, data[get_global_id(0)]);
        n++;
    }
}
```

### 5.4.4. Restrictions in the Implementation of Intel FPGA SDK for OpenCL Channels Extension

There are certain design restrictions to the implementation of channels in your OpenCL application.

## Multiple Channel Call Site

A kernel can read from the same channel multiple times. However, multiple kernels cannot read from the same channel. Similarly, a kernel can write to the same channel multiple times but multiple kernels cannot write to the same channel.

```
__kernel void k1() {  
    read_channel_intel (channell);  
    read_channel_intel (channell);  
    read_channel_intel (channell);  
}
```

The Intel FPGA SDK for OpenCL Offline Compiler cannot compile the following code and issues an error:

```
__kernel void k1(){  
    write_channel_intel (channell, 1);  
}  
  
__kernel void k2() {  
    write_channel_intel (channell, 2);  
}
```

## Feedback and Feed-forward Channels

Performance of a kernel that has multiple accesses (reads or writes) to the same channel might be poor.

## Static Indexing

The Intel FPGA SDK for OpenCL channels extension does support indexing into arrays of channel IDs, but it leads to inefficient hardware.

Consider the following example:

```
channel int ch[WORKGROUP_SIZE];  
  
__kernel void consumer()  
{  
  
    int gid = get_global_id(0);  
    int value = read_channel_intel(ch[gid]);  
  
    //statements  
}
```

Compilation of this example generates the following warning message:

```
Compiler Warning: Dynamic access into channel array ch was expanded into  
predicated  
static accesses on every channel of the array.
```

If the access is dynamic and you know that only a subset of the channels in the array can be accessed, you can generate slightly more efficient hardware with a switch statement:

```
channel int ch[WORKGROUP_SIZE];  
  
__kernel void consumer() {  
  
    int gid = get_global_id(0);  
    int value;
```

```
switch(gid)
{
    case 0: value = read_channel_intel(ch[0]); break;

    case 2: value = read_channel_intel(ch[2]); break;
    case 3: value = read_channel_intel(ch[3]); break;
    //statements

    case WORKGROUP_SIZE-1:read_channel_intel(ch[WORKGROUP_SIZE-1]); break;
}
//statements
}
```

### Kernel Vectorization Support

You cannot vectorize kernels that use channels; that is, do not include the `num_simd_work_items` kernel attribute in your kernel code. Vectorizing a kernel that uses channels creates multiple channel accesses inside the same kernel and requires arbitration, which negates the advantages of vectorization. As a result, the SDK's channel extension does not support kernel vectorization.

### Instruction-Level Parallelism on `read_channel_intel` and `write_channel_intel` Calls

If no data dependencies exist between `read_channel_intel` and `write_channel_intel` calls, the offline compiler attempts to execute these instructions in parallel. As a result, the offline compiler might execute these `read_channel_intel` and `write_channel_intel` calls in an order that does not follow the sequence expressed in the OpenCL kernel code.

Consider the following code sequence:

```
in_data1 = read_channel_intel(channel1);
in_data2 = read_channel_intel(channel2);
in_data3 = read_channel_intel(channel3);
```

Because there are no data dependencies between the `read_channel_intel` calls, the offline compiler can execute them in any order.

## 5.4.5. Enabling the Intel FPGA SDK for OpenCL Channels for OpenCL Kernel

To implement the Intel FPGA SDK for OpenCL channels extension, modify your OpenCL kernels to include channels-specific pragma and API calls.

To enable the channel extension, use the following pragma:

```
#pragma OPENCL EXTENSION cl_intel_channels : enable
```

Channel declarations are unique within a given OpenCL kernel program. Also, channel instances are unique for every OpenCL kernel program device pair. If the runtime loads a single OpenCL kernel program onto multiple devices, each device will have a single copy of the channel. However, these channel copies are independent and do not share data across the devices.

### 5.4.5.1. Declaring the Channel Handle

Use a channel variable to define the connectivity between kernels or between kernels and I/O.



To read from and write to a channel, the kernel must pass the channel variable to each of the corresponding API calls.

- Declare the channel handle as a file scope variable in the kernel source code using the following convention: `channel <type> <variable_name>`

For example: `channel int c;`

- The Intel FPGA SDK for OpenCL channel extension supports simultaneous channel accesses by multiple variables declared in a data structure. Declare a `struct` data structure for a channel in the following manner:

```
typedef struct type_ {
    int a;
    int b;
} type_t;

channel type_t foo;
```

### 5.4.5.2. Implementing Blocking Channel Writes

The `write_channel_intel` API call allows you to send data across a channel.

- To implement a blocking channel write, use the following `write_channel_intel` function signature:

```
void write_channel_intel (channel <type> channel_id, const
<type> data);
```

Where:

`channel_id` identifies the buffer to which the channel connects, and it must match the `channel_id` of the corresponding read channel (`read_channel_intel`).

`data` is the data that the channel write operation writes to the channel.

`<type>` defines a channel data width. Follow the OpenCL conversion rules to ensure that data the kernel writes to a channel is convertible to `<type>`.

The following code snippet demonstrates the implementation of the `write_channel_intel` API call:

```
//Defines chan, a kernel file-scope channel variable.
channel long chan;

/*Defines the kernel which reads eight bytes (size of long) from global memory,
and passes this data to the channel.*/
__kernel void kernel_write_channel( __global const long * src ) {
    for (int i = 0; i < N; i++) {
        //Writes the eight bytes to the channel.
        write_channel_intel(chan, src[i]);
    }
}
```

#### **Caution:**

When you send data across a channel using the `write_channel_intel` API call, keep in mind that if the channel is full (that is, if the FIFO buffer is full of data), your kernel stalls and waits until at least one data slot becomes available in the FIFO buffer. Use the Intel FPGA dynamic profiler for OpenCL to check for channel stalls.

## Related Information

[Profiling Your Kernel to Identify Performance Bottlenecks](#)

### 5.4.5.2.1. Implementing Nonblocking Channel Writes

Perform nonblocking channel writes to facilitate applications where writes to a full FIFO buffer should not cause the kernel to stall until a slot in the FIFO buffer becomes free. A nonblocking channel write returns a boolean value that indicates whether data was written successfully to the channel (that is, the channel was not full).

Consider a scenario where your application has one data producer with two identical workers. Assume the time each worker takes to process a message varies depending on the contents of the data. In this case, there might be situations where one worker is busy while the other is free. A nonblocking write can facilitate work distribution such that both workers are busy.

- To implement a nonblocking channel write, include the following `write_channel_nb_intel` function signature:

```
bool write_channel_nb_intel(channel <type> channel_id, const
<type> data);
```

The following code snippet of the kernel `producer` facilitates work distribution using the nonblocking channel write extension:

```
channel long worker0, worker1;
__kernel void producer( __global const long * src ) {
    for(int i = 0; i < N; i++) {

        bool success = false;
        do {
            success = write_channel_nb_intel(worker0, src[i]);
            if(!success) {
                success = write_channel_nb_intel(worker1, src[i]);
            }
        }
        while(!success);
    }
}
```

### 5.4.5.3. Implementing Blocking Channel Reads

The `read_channel_intel` API call allows you to receive data across a channel.

- To implement a blocking channel read, include the following `read_channel_intel` function signature:

```
<type> read_channel_intel(channel <type> channel_id);
```

Where:

`channel_id` identifies the buffer to which the channel connects, and it must match the `channel_id` of the corresponding write channel (`write_channel_intel`).

`<type>` defines a channel data width. Ensure that the variable the kernel assigns to read the channel data is convertible from `<type>`.

The following code snippet demonstrates the implementation of the `read_channel_intel` API call:

```
//Defines chan, a kernel file-scope channel variable.
channel long chan;

/*Defines the kernel, which reads eight bytes (size of long) from the channel
and writes it back to global memory.*/
__kernel void kernel_read_channel (__global long * dst); {
    for (int i = 0; i < N; i++) {
        //Reads the eight bytes from the channel.
        dst[i] = read_channel_intel(chan);
    }
}
```

**Caution:** If the channel is empty (that is, if the FIFO buffer is empty), you cannot receive data across a read channel using the `read_channel_intel` API call. Doing so causes your kernel to stall until at least one element of data becomes available from the FIFO buffer.

#### 5.4.5.3.1. Implementing Nonblocking Channel Reads

Perform nonblocking reads to facilitate applications where data is not always available and the operation should not wait for data to become available. The nonblocking read signature is similar to a blocking read. However, it populates the address pointed to by the bool pointer `valid` indicating whether a read operation successfully read data from the channel.

On a successful read (`valid` set to true), the value read from the channel is returned by the `read_channel_nb_intel` function. On a failed read (`valid` set to false), the return value of the `read_channel_nb_intel` function is not defined.

- To implement a blocking channel read, use the following `read_channel_nb_intel` function signature:

```
<type> read_channel_nb_intel(channel <type> channel_id, bool *
valid);
```

The following code snippet demonstrates the use of the nonblocking channel read extension:

```
channel long chan;

__kernel void kernel_read_channel (__global long * dst) {
    int i = 0;
    while (i < N) {
        bool valid0, valid1;
        long data0 = read_channel_nb_intel(chan, &valid0);
        long data1 = read_channel_nb_intel(chan, &valid1);
        if (valid0) {
            process(data0);
        }
        if (valid1) {
            process(data1);
        }
    }
}
```

#### 5.4.5.4. Implementing I/O Channels Using the io Channels Attribute

Include an `io` attribute in your channel declaration to declare a special I/O channel to interface with input or output features of an FPGA board. These features might include network interfaces, PCIe, cameras, or other data capture or processing devices or protocols.

The `io("chan_id")` attribute specifies the I/O feature of an accelerator board with which a channel is connected, where `chan_id` is the name of the I/O interface listed in the `board_spec.xml` file of your Custom Platform.

Because peripheral interface usage might differ for each device type, consult your board vendor's documentation when you implement I/O channels in your kernel program. Your OpenCL kernel code must be compatible with the type of data generated by the peripheral interfaces.

**Caution:**

- Implicit data dependencies might exist for channels that connect to the board directly and communicate with peripheral devices via I/O channels. These implicit data dependencies might lead to unexpected behavior because the Intel FPGA SDK for OpenCL Offline Compiler does not have visibility into these dependencies.
  - External I/O channels communicating with the same peripherals do not obey any sequential ordering. Ensure that the external device does not require sequential ordering because unexpected behavior might occur.
1. Consult the `board_spec.xml` file in your Custom Platform to identify the input and output features available on your FPGA board.

For example, a `board_spec.xml` file might include the following information on I/O features:

```
<channels>
  <interface name="udp_0" port="udp0_out" type="streamsource" width="256"
    chan_id="eth0_in"/>
  <interface name="udp_0" port="udp0_in" type="streamsink" width="256"
    chan_id="eth0_out"/>
  <interface name="udp_0" port="udp1_out" type="streamsource" width="256"
    chan_id="eth1_in"/>
  <interface name="udp_0" port="udp1_in" type="streamsink" width="256"
    chan_id="eth1_out"/>
</channels>
```

The `width` attribute of an `interface` element specifies the width, in bits, of the data type used by that channel. For the example above, both the `uint` and `float` data types are 32 bits wide. Other bigger or vectorized data types must match the appropriate bit width specified in the `board_spec.xml` file.

2. Implement the `io` channel attribute as demonstrated in the following code example. The `io` channel attribute names must match those of the I/O channels (`chan_id`) specified in the `board_spec.xml` file.

```
channel QUDPWord udp_in_IO __attribute__((depth(0)))
    __attribute__((io("eth0_in")));
channel QUDPWord udp_out_IO __attribute__((depth(0)))
    __attribute__((io("eth0_out")));

__kernel void io_in_kernel (__global ulong4 *mem_read,
    uchar read_from,
    int size)
{
    int index = 0;
    ulong4 data;
```

```

int half_size = size >> 1;
while (index < half_size)
{
    if (read_from & 0x01)
    {
        data = read_channel_intel(udp_in_IO);
    }
    else
    {
        data = mem_read[index];
    }
    write_channel_intel(udp_in, data);
    index++;
}
}

__kernel void io_out_kernel (__global ulong2 *mem_write,
                             uchar write_to,
                             int size)
{
    int index = 0;
    ulong4 data;
    int half_size = size >> 1;
    while (index < half_size)
    {
        ulong4 data = read_channel_intel(udp_out);
        if (write_to & 0x01)
        {
            write_channel_intel(udp_out_IO, data);
        }
        else
        {
            //only write data portion
            ulong2 udp_data;
            udp_data.s0 = data.s0;
            udp_data.s1 = data.s1;
            mem_write[index] = udp_data;
        }
        index++;
    }
}

```

**Attention:** Declare a unique `io("chan_id")` handle for each I/O channel specified in the channels eXtensible Markup Language (XML) element within the `board_spec.xml` file.

#### 5.4.5.5. Emulating I/O Channels

When you emulate a kernel that has a channel declared with the `io` attribute, I/O channel input is emulated by reading from a file, and channel output is emulated by writing to a file.

When you run your emulation, the file name used for reading or writing matches the name in the `io` attribute. For example, if you have a channel declaration as follows, your emulation reads or writes (but not both) to a file called `myIOChannel`.

```
channel uint chanA __attribute__((io("myIOChannel")));
```

I/O channels are unidirectional. You can either read from a channel or write to a channel, but not both. However, you can have separate read channels and write channels with the same `io` attribute value.

```
channel uint readChannel __attribute__((io("myIOChannel")));
channel uint writeChannel __attribute__((io("myIOChannel")));
```

### Emulating Reading from an I/O Channel

If a read is issued from a channel with an `io` attribute called `myfile`, a read attempt is made on a file on the disk called `myfile`. If the `myfile` file does not exist or there is insufficient data to read from the file, the behavior depends on the call type:

*Non-blocking read*      If the file does not exist or there is insufficient data, the read attempt returns with a failure message.

*Blocking read*      If the file does not exist or there is insufficient data, the read attempt blocks your program until the file is created on the disk, or the file contains sufficient data.

### Emulating Writing to an I/O Channel

If a write is issued to a channel with an `io` attribute called `myfile`, a write attempt is made to a file on the disk called `myfile`. If the `myfile` file does not exist, a regular file is created and written to. If the `myfile` file exists, it is overwritten. If the write fails, the behavior depends on the call type:

*Non-blocking write*    If the write attempt fails, an error is returned.

*Blocking write*      If the write attempt fails, further write attempts are made.

#### 5.4.5.6. Use Models of Intel FPGA SDK for OpenCL Channels Implementation

Concurrent kernel execution can improve the efficiency of channels on an FPGA. To achieve concurrent execution, the host launches the kernels in parallel. The kernels can communicate with each other through channels where applicable.

The following use models provide an overview on how to exploit concurrent execution safely and efficiently.

##### Feed-Forward Design Model

Implement the feed-forward design model to send data from one kernel to the next without creating any cycles between them. Consider the following code example:

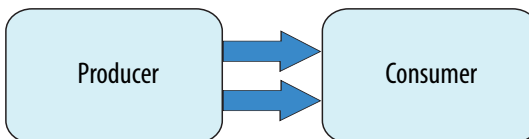
```
__kernel void producer (__global const uint * src,
                       const uint iterations)
{
    for (int i = 0; i < iterations; i++)
    {
        write_channel_intel(c0, src[2*i]);
        write_channel_intel(c1, src[2*i+1]);
    }
}

__kernel void consumer (__global uint * dst,
                       const uint iterations)
{
    for (int i = 0; i < iterations; i++)
    {
        dst[2*i] = read_channel_intel(c0);
    }
}
```

```
dst[2*i+1] = read_channel_intel(c1);  
}  
}
```

The `producer` kernel writes data to channels `c0` and `c1`. The `consumer` kernel reads data from `c0` and `c1`. The figure below illustrates the feed-forward data flow between the two kernels:

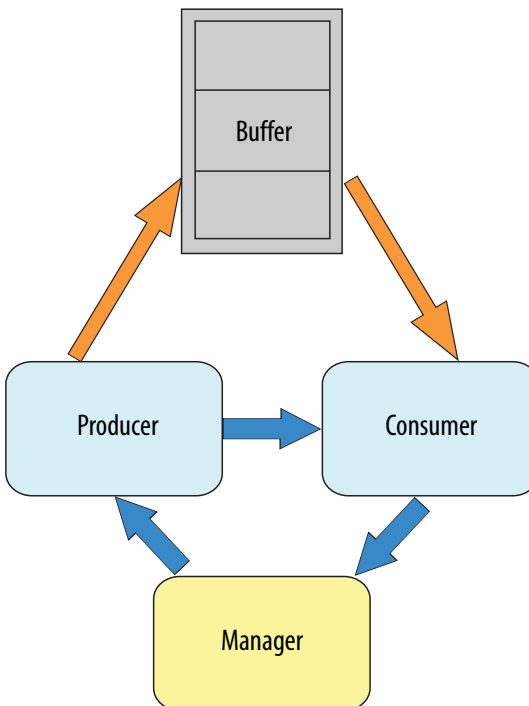
**Figure 8. Feed-Forward Data Flow**



### Buffer Management

In the feed-forward design model, data traverses between the `producer` and `consumer` kernels one word at a time. To facilitate the transfer of large data messages consisting of several words, you can implement a ping-pong buffer, which is a common design pattern found in applications for communication. The figure below illustrates the interactions between kernels and a ping-pong buffer:

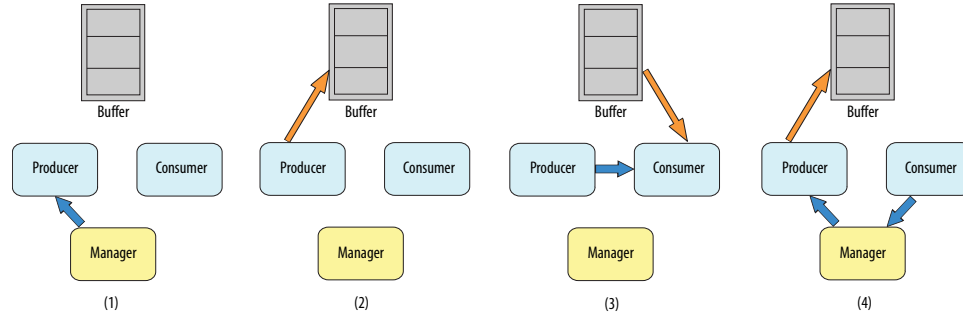
**Figure 9. Feed-Forward Design Model with Buffer Management**



The `manager` kernel manages circular buffer allocation and deallocation between the `producer` and `consumer` kernels. After the `consumer` kernel processes data, the `manager` receives memory regions that the `consumer` frees up and sends them to the `producer` for reuse. The `manager` also sends to the `producer` kernel the initial set of free locations, or tokens, to which the `producer` can write data.

The following figure illustrates the sequence of events that take place during buffer management:

**Figure 10. Kernels Interaction during Buffer Management**



1. The manager kernel sends a set of tokens to the producer kernel to indicate initially which regions in memory are free for producer to use.
2. After manager allocates the memory region, producer writes data to that region of the ping-pong buffer.
3. After producer completes the write operation, it sends a synchronization token to the consumer kernel to indicate what memory region contains data for processing. The consumer kernel then reads data from that region of the ping-pong buffer.

*Note:* When consumer is performing the read operation, producer can write to other free memory locations for processing because of the concurrent execution of the producer, consumer, and manager kernels.

4. After consumer completes the read operation, it releases the memory region and sends a token back to the manager kernel. The manager kernel then recycles that region for producer to use.

### Implementation of Buffer Management for OpenCL Kernels

To ensure that the SDK implements buffer management properly, the ordering of channel read and write operations is important. Consider the following kernel example:

```
__kernel void producer (__global const uint * restrict src,
                       __global volatile uint * restrict shared_mem,
                       const uint iterations)
{
    int base_offset;

    for (uint gID = 0; gID < iterations; gID++)
    {
        // Assume each block of memory is 256 words
        uint lID = 0x0ff & gID;

        if (lID == 0)
        {
            base_offset = read_channel_intel(req);
        }

        shared_mem[base_offset + lID] = src[gID];

        // Make sure all memory operations are committed before
        // sending token to the consumer
        mem_fence(CLK_GLOBAL_MEM_FENCE | CLK_CHANNEL_MEM_FENCE);
    }
}
```



```
        if (lID == 255)
        {
            write_channel_intel(c, base_offset);
        }
    }
}
```

In this kernel, because the following lines of code are independent, the Intel FPGA SDK for OpenCL Offline Compiler can schedule them to execute concurrently:

```
shared_mem[base_offset + lID] = src[gID];
```

and

```
write_channel_intel(c, base_offset);
```

Writing data to `base_offset` and then writing `base_offset` to a channel might be much faster than writing data to global memory. The consumer kernel might then read `base_offset` from the channel and use it as an index to read from global memory. Without synchronization, consumer might read data from producer before `shared_mem[base_offset + lID] = src[gID];` finishes executing. As a result, consumer reads in invalid data. To avoid this scenario, the synchronization token must occur after the producer kernel commits data to memory. In other words, a consumer kernel cannot consume data from the producer kernel until producer stores its data in global memory successfully.

To preserve this ordering, include an OpenCL `mem_fence` token in your kernels. The `mem_fence` construct takes two flags: `CLK_GLOBAL_MEM_FENCE` and `CLK_CHANNEL_MEM_FENCE`. The `mem_fence` effectively creates a control flow dependence between operations that occur before and after the `mem_fence` call. The `CLK_GLOBAL_MEM_FENCE` flag indicates that global memory operations must obey the control flow. The `CLK_CHANNEL_MEM_FENCE` indicates that channel operations must obey the control flow. As a result, the `write_channel_intel` call in the example cannot start until the global memory operation is committed to the shared memory buffer.

#### 5.4.5.7. Implementing Buffered Channels Using the depth Channels Attribute

You may have buffered or unbuffered channels in your kernel program. If there are imbalances in channel read and write operations, create buffered channels to prevent kernel stalls by including the `depth` attribute in your channel declaration. Buffered channels decouple the operation of concurrent work-items executing in different kernels.

You may use a buffered channel to control data traffic, such as limiting throughput or synchronizing accesses to shared memory. In an unbuffered channel, a write operation cannot proceed until the read operation reads a data value. In a buffered channel, a write operation cannot proceed until the data value is copied to the buffer. If the buffer is full, the operation cannot proceed until the read operation reads a piece of data and removes it from the channel.

- If you expect any temporary mismatch between the consumption rate and the production rate to the channel, set the buffer size using the `depth` channel attribute.

The following example demonstrates the use of the `depth` channel attribute in kernel code that implements the Intel FPGA SDK for OpenCL channels extension. The `depth(N)` attribute specifies the minimum depth of a buffered channel, where  $N$  is the number of data values.

```
channel int c __attribute__((depth(10)));

__kernel void producer (__global int * in_data)
{
    for (int i = 0; i < N; i++)
    {
        if (in_data[i])
        {
            write_channel_intel(c, in_data[i]);
        }
    }
}

__kernel void consumer (__global int * restrict check_data,
                       __global int * restrict out_data)
{
    int last_val = 0;

    for (int i = 0; i < N, i++)
    {
        if (check_data[i])
        {
            last_val = read_channel_intel(c);
        }
        out_data[i] = last_val;
    }
}
```

In this example, the write operation can write ten data values to the channel without blocking. Once the channel is full, the write operation cannot proceed until an associated read operation to the channel occurs.

Because the channel read and write calls are conditional statements, the channel might experience an imbalance between read and write calls. You may add a buffer capacity to the channel to ensure that the `producer` and `consumer` kernels are decoupled. This step is particularly important if the `producer` kernel is writing data to the channel when the `consumer` kernel is not reading from it.

#### 5.4.5.8. Enforcing the Order of Channel Calls

To enforce the order of channel calls, introduce memory fence or barrier functions in your kernel program to control memory accesses. A memory fence function is necessary to create a control flow dependence between the channel accesses before and after the fence.

When the Intel FPGA SDK for OpenCL Offline Compiler generates a compute unit, it does not always create instruction-level parallelism on all instructions that are independent of each other. As a result, channel read and write operations might not execute independently of each other even if there is no control or data dependence between them. When channel calls interact with each other, or when channels write data to external devices, deadlocks might occur.

For example, the code snippet below consists of a `producer` kernel and a `consumer` kernel. Channels `c0` and `c1` are unbuffered channels. The schedule of the channel read operations from `c0` and `c1` might occur in the reversed order as the channel write operations to `c0` and `c1`. That is, the `producer` kernel writes to `c0` but the `consumer` kernel might read from `c1` first. This rescheduling of channel calls might cause a deadlock because the `consumer` kernel is reading from an empty channel.

```
__kernel void producer (__global const uint * src,
                       const uint iterations)
{
    for (int i = 0; i < iterations; i++)
    {
        write_channel_intel(c0, src[2*i]);
        write_channel_intel(c1, src[2*i+1]);
    }
}

__kernel void consumer (__global uint * dst,
                      const uint iterations)
{
    for (int i = 0; i < iterations; i++)
    {
        /*During compilation, the AOC might reorder the way the consumer kernel
        writes to memory to optimize memory access. Therefore, c1 might be read
        before c0, which is the reverse of what appears in code.*/

        dst[2*i+1] = read_channel_intel(c0);
        dst[2*i] = read_channel_intel(c1);
    }
}
```

- To prevent deadlocks from occurring by enforcing the order of channel calls, include memory fence functions (`mem_fence`) in your kernel.

Inserting the `mem_fence` call with each kernel's channel flag forces the sequential ordering of the write and read channel calls. The code snippet below shows the modified `producer` and `consumer` kernels:

```
channel uint c0 __attribute__((depth(0)));
channel uint c1 __attribute__((depth(0)));

__kernel void producer (__global const uint * src,
                       const uint iterations)
{
    for (int i = 0; i < iterations; i++)
    {
        write_channel_intel(c0, src[2*i]);
        mem_fence(CLK_CHANNEL_MEM_FENCE);
        write_channel_intel(c1, src[2*i+1]);
    }
}

__kernel void consumer (__global uint * dst;
                      const uint iterations)
{
    for (int i = 0; i < iterations; i++)
    {
        dst[2*i+1] = read_channel_intel(c0);
        mem_fence(CLK_CHANNEL_MEM_FENCE);
    }
}
```

```

        dst[2*i] = read_channel_intel(c1);
    }
}

```

In this example, `mem_fence` in the producer kernel ensures that the channel write operation to `c0` occurs before that to `c1`. Similarly, `mem_fence` in the consumer kernel ensures that the channel read operation from `c0` occurs before that from `c1`.

#### 5.4.5.8.1. Defining Memory Consistency Across Kernels When Using Channels

According to the OpenCL Specification version 1.0, memory behavior is undefined unless a kernel completes execution. A kernel must finish executing before other kernels can visualize any changes in memory behavior. However, kernels that use channels can share data through common global memory buffers and synchronized memory accesses. To ensure that data written to a channel is visible to the read channel after execution passes a memory fence, define memory consistency across kernels with respect to memory fences.

- To create a control flow dependency between the channel synchronization calls and the memory operations, add the `CLK_GLOBAL_MEM_FENCE` flag to the `mem_fence` call.

For example:

```

__kernel void producer( __global const uint * src,
                       const uint iterations )
{
    for(int i=0; i < iterations; i++)
    {
        write_channel_intel(c0, src[2*i]);
        mem_fence(CLK_CHANNEL_MEM_FENCE | CLK_GLOBAL_MEM_FENCE);
        write_channel_intel(c1, src[2*i+1]);
    }
}

```

In this kernel, the `mem_fence` function ensures that the write operation to `c0` and memory access to `src[2*i]` occur before the write operation to `c1` and memory access to `src[2*i+1]`. This allows data written to `c0` to be visible to the read channel before data is written to `c1`.

## 5.5. Implementing OpenCL Pipes

The Intel FPGA SDK for OpenCL provides preliminary support for OpenCL pipe functions. OpenCL pipes are part of the OpenCL Specification version 2.0. They provide a mechanism for passing data to kernels and synchronizing kernels with high efficiency and low latency.

Implement pipes if it is important that your OpenCL kernel is compatible with other SDKs.

Refer to the *OpenCL Specification version 2.0* for OpenCL C programming language specification and general information about pipes.

The Intel FPGA SDK for OpenCL implementation of pipes does not encompass the entire pipes specification. As such, it is not fully conformant to the OpenCL Specification version 2.0. The goal of the SDK's pipes implementation is to provide a

solution that works seamlessly on a different OpenCL 2.0-conformant device. To enable pipes for Intel FPGA products, your design must satisfy certain additional requirements.

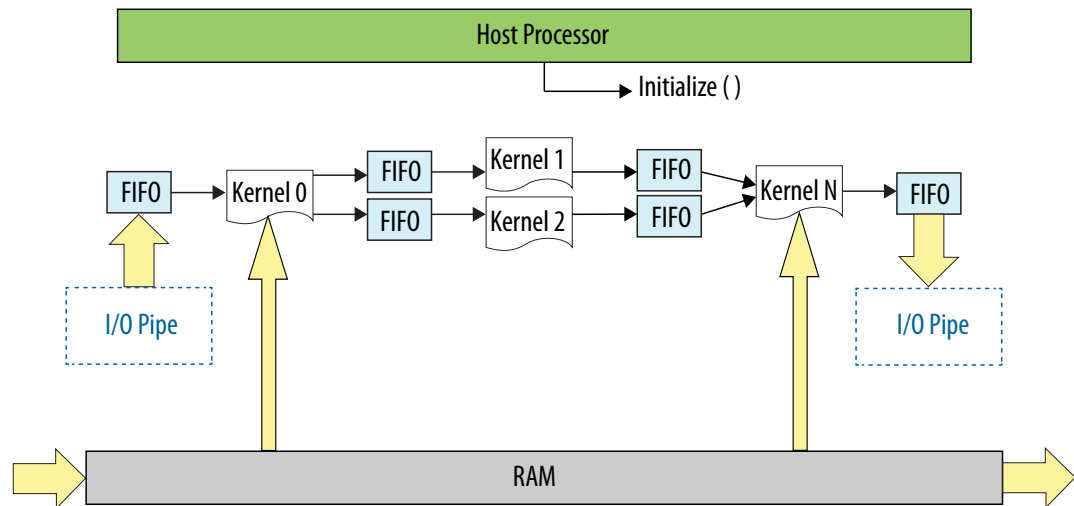
### Related Information

[OpenCL Specification version 2.0 \(API\)](#)

## 5.5.1. Overview of the OpenCL Pipe Functions

OpenCL pipes allow kernels to communicate directly with each other via FIFO buffers.

**Figure 11. Overview of a Pipe Network Implementation**



Implementation of pipes decouples kernel execution from the host processor. The foundation of the Intel FPGA SDK for OpenCL pipes support is the SDK's channels extension. However, the syntax for pipe functions differs from the channels syntax.

**Important:** Unlike channels, pipes have a default nonblocking behavior.

For more information about blocking and nonblocking functions, refer to the corresponding documentation on channels.

### Related Information

- [Implementing Blocking Channel Writes](#) on page 49
- [Implementing Nonblocking Channel Writes](#) on page 50
- [Implementing Nonblocking Channel Reads](#) on page 51
- [Implementing Blocking Channel Reads](#) on page 50

### 5.5.2. Pipe Data Behavior

Data written to a pipe remains in a pipe provided that the kernel program remains loaded on the FPGA device. In other words, data written to a pipe persists across multiple work-groups and NDRange invocations. However, data is not persistent across multiple or different invocations of kernel programs that result in FPGA reprogramming operations.

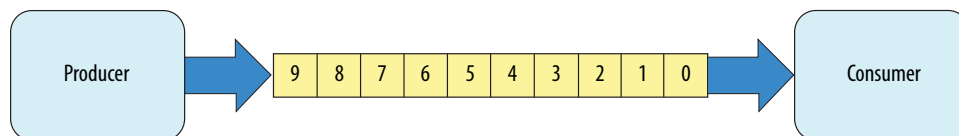
Consider the following code example:

```
__kernel void producer (write_only pipe uint __attribute__((blocking)) c0)
{
    for (uint i = 0; i < 10; i++)
    {
        write_pipe (c0, &i);
    }
}

__kernel void consumer (__global uint * restrict dst,
                        read_only pipe uint __attribute__((blocking))
                        __attribute__((depth(10))) c0)
{
    for (int i = 0; i < 5; i++)
    {
        read_pipe (c0, &dst[i]);
    }
}
```

A read operation to a pipe reads the *least* recent piece of data written to the pipe first. Pipe data maintains a FIFO ordering within the pipe.

**Figure 12. Pipe Data FIFO Ordering**



The kernel `producer` writes ten elements (`[0, 9]`) to the pipe. The kernel `consumer` reads five elements from the pipe per NDRange invocation. During the first invocation, the kernel `consumer` reads values 0 to 4 from the pipe. Because the data persists across NDRange invocations, the second time you execute the kernel `consumer`, it reads values 5 to 9.

For this example, to avoid a deadlock from occurring, you need to invoke the kernel `consumer` twice for every invocation of the kernel `producer`. If you call `consumer` less than twice, `producer` stalls because the pipe becomes full. If you call `consumer` more than twice, `consumer` stalls because there is insufficient data in the pipe.

### 5.5.3. Multiple Work-Item Ordering for Pipes

The OpenCL specification does not define a work-item ordering. The Intel FPGA SDK for OpenCL enforces a work-item order to maintain the consistency in pipe read and write operations.

Multiple work-item accesses to a pipe can be useful in some scenarios. For example, they are useful when data words in the pipe are independent, or when the pipe is implemented for control logic. The main concern regarding multiple work-item accesses to a pipe is the order in which the kernel writes data to and reads data from

the pipe. If possible, the OpenCL pipes process work-items read and write operations to a pipe in a deterministic order. As such, the read and write operations remain consistent across kernel invocations.

### Requirements for Deterministic Multiple Work-Item Ordering

To guarantee deterministic ordering, the SDK checks that the pipe call is work-item invariant based on the following characteristics:

- All paths through the kernel must execute the pipe call.
- If the first requirement is not satisfied, none of the branch conditions that reach the pipe call should execute in a work-item-dependent manner.

If the SDK cannot guarantee deterministic ordering of multiple work-item accesses to a pipe, it warns you that the pipes might not have well-defined ordering with nondeterministic execution. Primarily, the SDK fails to provide deterministic ordering if you have work-item-variant code on loop executions with pipe calls, as illustrated below:

```
__kernel void
ordering (__global int * check, global int * data,
         write_only pipe int __attribute__((blocking)) req)
{
    int condition = check[get_global_id(0)];

    if (condition)
    {
        for (int i = 0; i < N; i++)
        {
            process(data);
            write_pipe (req, &data[i]);
        }
    }
    else
    {
        process(data);
    }
}
```

#### 5.5.3.1. Work-item Serial Execution of Pipes

Work-item serial execution refers to an ordered execution behavior where work-item sequential IDs determine their execution order in the compute unit.

When you implement pipes in a kernel, the Intel FPGA SDK for OpenCL Offline Compiler enforces that kernel behavior is equivalent to having at most one work-group in flight. The offline compiler also ensures that the kernel executes pipes in work-item serial execution, where the kernel executes work items with smaller IDs first. A work-item has the identifier ( $x$ ,  $y$ ,  $z$ ,  $group$ ), where  $x$ ,  $y$ ,  $z$  are the local 3D identifiers, and  $group$  is the work-group identifier.

The work-item ID ( $x_0$ ,  $y_0$ ,  $z_0$ ,  $group_0$ ) is considered to be smaller than the ID ( $x_1$ ,  $y_1$ ,  $z_1$ ,  $group_1$ ) if one of the following conditions is true:

- $group_0 < group_1$
- $group_0 = group_1$  and  $z_0 < z_1$
- $group_0 = group_1$  and  $z_0 = z_1$  and  $y_0 < y_1$
- $group_0 = group_1$  and  $z_0 = z_1$  and  $y_0 = y_1$  and  $x_0 < x_1$

Work items with incremental IDs execute in a sequential order. For example, the work-item with an ID ( $x_0, y_0, z_0, group_0$ ) executes the write channel call first. Then, the work-item with an ID ( $x_1, y_0, z_0, group_0$ ) executes the call, and so on. Defining this order ensures that the system is verifiable with external models.

### Pipe Execution in a Loop with Multiple Work Items

Only in kernels compiled as an NDRange kernel, when pipes exist in the body of a loop with multiple work items, as shown below, each loop iteration executes prior to subsequent iterations. This implies that loop iteration 0 of each work item in a work group executes before iteration 1 of each work item in a work group, and so on.

```
__kernel void ordering (__global int * data,
                      write_only pipe int __attribute__((blocking)) req)
{
    write_pipe (req, &data[get_global_id(0)]);
}
```

## 5.5.4. Restrictions in OpenCL Pipes Implementation

There are certain design restrictions to the implementation of pipes in your OpenCL application.

### Default Behavior

By default, pipes exhibit nonblocking behavior. If you want the pipes in your kernel to exhibit blocking behavior, specify the blocking attribute (`__attribute__((blocking))`) when you declare the read and write pipes.

### Emulation Support

The Intel FPGA SDK for OpenCL Emulator supports emulation of kernels that contain pipes. The level of Emulator support aligns with the subset of OpenCL pipes support that is implemented for the FPGA hardware.

### Pipes API Support

Currently, the SDK's implementation of pipes does not support all the built-in pipe functions in the OpenCL Specification version 2.0. For a list of supported and unsupported pipe APIs, refer to *OpenCL 2.0 C Programming Language Restrictions for Pipes*.

### Single Call Site

Because the pipe read and write operations do not function deterministically, for a given kernel, you can only assign one call site per pipe ID. For example, the Intel FPGA SDK for OpenCL Offline Compiler cannot compile the following code example:

```
read_pipe(pipe1, &in_data1);
read_pipe(pipe2, &in_data2);
read_pipe(pipe1, &in_data3);
```

The second `read_pipe` call to `pipe1` causes compilation failure because it creates a second call site to `pipe1`.



To gather multiple data from a given pipe, divide the pipe into multiple pipes, as shown below:

```
read_pipe(pipe1, &in_data1);  
read_pipe(pipe2, &in_data2);  
read_pipe(pipe3, &in_data3);
```

Because you can only assign a single call site per pipe ID, you cannot unroll loops containing pipes. Consider the following code:

```
#pragma unroll 4  
for (int i = 0; i < 4; i++)  
{  
    read_pipe (pipe1, &in_data1);  
}
```

The offline compiler issues the following warning message during compilation:

Compiler Warning: Unroll is required but the loop cannot be unrolled.

### Feedback and Feed-Forward Pipes

Pipes within a kernel can be either `read_only` or `write_only`. Performance of a kernel that reads and writes to the same pipe is poor.

### Kernel Vectorization Support

You cannot vectorize kernels that use pipes; that is, do not include the `num_simd_work_items` kernel attribute in your kernel code. Vectorizing a kernel that uses pipes creates multiple pipe masters and requires arbitration, which OpenCL pipes specification does not support.

### Instruction-Level Parallelism on `read_pipe` and `write_pipe` Calls

If no data dependencies exist between `read_pipe` and `write_pipe` calls, the offline compiler attempts to execute these instructions in parallel. As a result, the offline compiler might execute these `read_pipe` and `write_pipe` calls in an order that does not follow the sequence expressed in the OpenCL kernel code.

Consider the following code sequence:

```
in_data1 = read_pipe(pipe1);  
in_data2 = read_pipe(pipe2);  
in_data3 = read_pipe(pipe3);
```

Because there are no data dependencies between the `read_pipe` calls, the offline compiler can execute them in any order.

### Related Information

[OpenCL 2.0 C Programming Language Restrictions for Pipes](#) on page 217

## 5.5.5. Enabling OpenCL Pipes for Kernels

To implement pipes, modify your OpenCL kernels to include pipes-specific API calls.

Pipes declarations are unique within a given OpenCL kernel program. Also, pipe instances are unique for every OpenCL kernel program-device pair. If the runtime loads a single OpenCL kernel program onto multiple devices, each device will have a single copy of each pipe. However, these pipe copies are independent and do not share data across the devices.

### 5.5.5.1. Ensuring Compatibility with Other OpenCL SDKs

Currently, Intel's implementation of OpenCL pipes is partially conformant to the OpenCL Specification version 2.0. If you port a kernel that implements pipes from another OpenCL SDK to the Intel FPGA SDK for OpenCL, you must modify the host code and the kernel code. The modifications do not affect subsequent portability of your application to other OpenCL SDKs.

#### Original Program Code

Below is an example of an OpenCL application:

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include "CL/opencl.h"
#define SIZE 1000

const char *kernel_source = "__kernel void pipe_writer(__global int *in,"
                             "                               write_only pipe int p_in)
\n"
                             "{\n"
                             "    int gid = get_global_id(0);\n"
                             "    write_pipe(p_in, &in[gid]);\n"
                             "}\n"
                             "__kernel void pipe_reader(__global int *out,"
                             "                               read_only pipe int p_out)
\n"
                             "{\n"
                             "    int gid = get_global_id(0);\n"
                             "    read_pipe(p_out, &out[gid]);\n"
                             "}\n";

int main()
{
    int *input = (int *)malloc(sizeof(int) * SIZE);
    int *output = (int *)malloc(sizeof(int) * SIZE);
    memset(output, 0, sizeof(int) * SIZE);
    for (int i = 0; i != SIZE; ++i)
    {
        input[i] = rand();
    }

    cl_int status;
    cl_platform_id platform;
    cl_uint num_platforms;
    status = clGetPlatformIDs(1, &platform, &num_platforms);

    cl_device_id device;
    cl_uint num_devices;
    status = clGetDeviceIDs(platform,
                            CL_DEVICE_TYPE_ALL,
                            1,
                            &device,
                            &num_devices);

    cl_context context = clCreateContext(0, 1, &device, NULL, NULL, &status);
    cl_command_queue queue = clCreateCommandQueue(context, device, 0, &status);
```

```

size_t len = strlen(kernel_source);
cl_program program = clCreateProgramWithSource(context,
                                              1,
                                              (const char **)&kernel_source,
                                              &len,
                                              &status);

status = clBuildProgram(program, num_devices, &device, "", NULL, NULL);

cl_kernel pipe_writer = clCreateKernel(program, "pipe_writer", &status);
cl_kernel pipe_reader = clCreateKernel(program, "pipe_reader", &status);

cl_mem in_buffer = clCreateBuffer(context,
                                  CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
                                  sizeof(int) * SIZE,
                                  input,
                                  &status);
cl_mem out_buffer = clCreateBuffer(context,
                                   CL_MEM_WRITE_ONLY,
                                   sizeof(int) * SIZE,
                                   NULL,
                                   &status);

cl_mem pipe = clCreatePipe(context, 0, sizeof(cl_int), SIZE, NULL, &status);

status = clSetKernelArg(pipe_writer, 0, sizeof(cl_mem), &in_buffer);
status = clSetKernelArg(pipe_writer, 1, sizeof(cl_mem), &pipe);
status = clSetKernelArg(pipe_reader, 0, sizeof(cl_mem), &out_buffer);
status = clSetKernelArg(pipe_reader, 1, sizeof(cl_mem), &pipe);

size_t size = SIZE;
cl_event sync;
status = clEnqueueNDRangeKernel(queue,
                                pipe_writer,
                                1,
                                NULL,
                                &size,
                                &size,
                                0,
                                NULL,
                                &sync);
status = clEnqueueNDRangeKernel(queue,
                                pipe_reader,
                                1,
                                NULL,
                                &size,
                                &size,
                                1,
                                &sync,
                                NULL);

status = clFinish(queue);

status = clEnqueueReadBuffer(queue,
                             out_buffer,
                             CL_TRUE,
                             0,
                             sizeof(int) * SIZE,
                             output,
                             0,
                             NULL,
                             NULL);

int golden = 0, result = 0;
for (int i = 0; i != SIZE; ++i)
{
    golden += input[i];
    result += output[i];
}

int ret = 0;
if (golden != result)

```

```
    {
        printf("FAILED!");
        ret = 1;
    } else
    {
        printf("PASSED!");
    }
    printf("\n");

    return ret;
}
```

### Host Code Modification

If the original host code runs on OpenCL SDKs that conforms to the OpenCL Specification version 2.0, you must modify it before running it on the Intel FPGA SDK for OpenCL. To modify the host code, perform the following changes:

1. Use the `clCreateProgramWithBinary` function instead of the `clCreateProgramWithSource` function to create the program.
2. Move the contents of the `kernel_source` string into a separate source file. Refer to [Kernel Code Modification](#) on page 68 for more information.

### Kernel Code Modification

If your kernel code runs on OpenCL SDKs that conforms to the OpenCL Specification version 2.0, you must modify it before running it on the Intel FPGA SDK for OpenCL. To modify the kernel code, perform the following changes:

1. Create a separate source (`.cl`) file for the kernel code.
2. Rename the pipe arguments so that they are the same in both kernels. For example, rename `p_in` and `p_out` to `p`.
3. Specify the `depth` attribute for the pipe arguments. Assign a `depth` attribute value that equals to the maximum number of packets that the pipe creates to hold in the host.
4. Build the kernel program in the offline compilation mode because the Intel FPGA SDK for OpenCL has an offline compiler.

The modified kernel code appears as follows:

```
#define SIZE 1000

__kernel void pipe_writer(__global int *in,
                          write_only pipe int __attribute__((depth(SIZE))) p)
{
    int gid = get_global_id(0);
    write_pipe(p, &in[gid]);
}

__kernel void pipe_reader(__global int *out,
                          read_only pipe int __attribute__((depth(SIZE))) p)
{
    int gid = get_global_id(0);
    read_pipe(p, &out[gid]);
}
```

#### 5.5.5.2. Declaring the Pipe Handle

Use the `pipe` variable to define the static pipe connectivity between kernels or between kernels and I/O.

To read from and write to a pipe, the kernel must pass the pipe variable to each of the corresponding API call.

- Declare the pipe handle as a file scope variable in the kernel source code in the following convention: `<access_qualifier> pipe <type> <variable_name>`

The `<type>` of the pipe may be any OpenCL built-in scalar or vector data type with a scalar size of 1024 bits or less. It may also be any user-defined type that is comprised of scalar or vector data type with a scalar size of 1024 bits or less.

Consider the following pipe handle declarations:

```
__kernel void first (pipe int c)
__kernel void second (write_only pipe int c)
```

The first example declares a read-only pipe handle of type `int` in the kernel `first`. The second example declares a write-only pipe in the kernel `second`. The kernel `first` may only read from pipe `c`, and the kernel `second` may only write to pipe `c`.

**Important:** The Intel FPGA SDK for OpenCL Offline Compiler statically infers the connectivity of pipes in your system by matching the names of the pipe arguments. In the example above, the kernel `first` is connected to the kernel `second` by the pipe `c`.

In an Intel OpenCL system, only one kernel may read to a pipe. Similarly, only one kernel may write to a pipe. If a non-I/O pipe does not have at least one corresponding reading operation and one writing operation, the offline compiler issues an error.

For more information in the Intel FPGA SDK for OpenCL I/O pipe implementation, refer to *Implementing I/O Pipes Using the io Attribute*.

### Related Information

[Implementing I/O Pipes Using the io Attribute](#) on page 72

#### 5.5.5.3. Implementing Pipe Writes

The `write_pipe` API call allows you to send data across a pipe.

Intel only supports the convenience version of the `write_pipe` function. By default, `write_pipe` calls are nonblocking. Pipe write operations are successful only if there is capacity in the pipe to hold the incoming packet.

- To implement a pipe write, include the following `write_pipe` function signature:

```
int write_pipe (write_only pipe <type> pipe_id, const <type>
*data);
```

Where:

`pipe_id` identifies the buffer to which the pipe connects, and it must match the `pipe_id` of the corresponding read pipe (`read_pipe`).

`data` is the data that the pipe write operation writes to the pipe. It is a pointer to the packet type of the pipe. Note that writing to the pipe might lead to a global or local memory load, depending on the source address space of the data pointer.

`<type>` defines a pipe data width. The return value indicates whether the pipe write operation is successful. If successful, the return value is 0. If pipe write is unsuccessful, the return value is -1.

The following code snippet demonstrates the implementation of the `write_pipe` API call:

```
/*Declares the writable nonblocking pipe, p, which contains packets of type int*/
__kernel void kernel_write_pipe (__global const long *src,
                                write_only pipe int p)
{
    for (int i = 0; i < N; i++)
    {
        //Performs the actual writing
        //Emulates blocking behavior via the use of a while loop
        while (write_pipe(p, &src[i]) < 0) { }
    }
}
```

The `while` loop is unnecessary if you specify a *blocking attribute*. To facilitate better hardware implementations, Intel provides facility for blocking `write_pipe` calls by specifying the blocking attribute (that is, `__attribute__((blocking))`) on the pipe argument declaration for the kernel. Blocking `write_pipe` calls always return success.

**Caution:** When you send data across a blocking write pipe using the `write_pipe` API call, keep in mind that if the pipe is full (that is, if the FIFO buffer is full of data), your kernel stalls. Use the Intel FPGA dynamic profiler for OpenCL to check for pipe stalls.

### Related Information

[Profiling Your Kernel to Identify Performance Bottlenecks](#)

#### 5.5.5.4. Implementing Pipe Reads

The `read_pipe` API call allows you to receive data across a pipe.

Intel only supports the convenience version of the `read_pipe` function. By default, `read_pipe` calls are non-blocking.

- To implement a pipe read, include the following `read_pipe` function signature:

```
int read_pipe (read_only_pipe <type> pipe_id, <type> *data);
```

Where:

- `pipe_id` identifies the buffer to which the pipe connects, and it must match the `pipe_id` of the corresponding pipe write operation (`write_pipe`).
- `data` is the data that the pipe read operation reads from the pipe. It is a pointer to the location of the data.

*Note:* `read_pipe` call might lead to a global or local memory load, depending on the source address space of the data pointer.

- `<type>` defines the packet size of the data.

The following code snippet demonstrates the implementation of the `read_pipe` API call:

```
/*Declares the read_only_pipe that contains packets
of type long.*/
/*Declares that read_pipe calls within the kernel will exhibit
blocking behavior*/
__kernel void kernel_read_pipe(__global long *dst,
                              read_only pipe long __attribute__((blocking)) p)
{
```

```

for (int i = 0; i < N; i++)
{
    /*Reads from a long from the pipe and stores it
    into global memory at the specified location*/
    read_pipe(p, &dst[i]);
}

```

To facilitate better hardware implementations, Intel provides facility for blocking `read_pipe` calls by specifying the blocking attribute (that is, `__attribute__((blocking))`) on the pipe argument declaration for the kernel. Blocking `read_pipe` calls always return success.

**Caution:** If the pipe is empty (that is, if the FIFO buffer is empty), you cannot receive data across a blocking read pipe using the `read_pipe` API call. Doing so causes your kernel to stall.

### 5.5.5.5. Implementing Buffered Pipes Using the `depth` Attribute

You may have buffered or unbuffered pipes in your kernel program. If there are imbalances in pipe read and write operations, create buffered pipes to prevent kernel stalls by including the `depth` attribute in your pipe declaration. Buffered pipes decouple the operation of concurrent work-items executing in different kernels.

You may use a buffered pipe to control data traffic, such as limiting throughput or synchronizing accesses to shared memory. In an unbuffered pipe, a write operation can only proceed when the read operation is expecting to read data. Use unbuffered pipes in conjunction with blocking read and write behaviors in kernels that execute concurrently. The unbuffered pipes provide self-synchronizing data transfers efficiently.

In a buffered pipe, a write operation can only proceed if there is capacity in the pipe to hold the incoming packet. A read operation can only proceed if there is at least one packet in the pipe.

Use buffered pipes if pipe calls are predicated differently in the writer and reader kernels, and the kernels do not execute concurrently.

- If you expect any temporary mismatch between the consumption rate and the production rate to the pipe, set the buffer size using the `depth` attribute.

The following example demonstrates the use of the `depth` attribute in kernel code that implements the OpenCL pipes. The `depth(N)` attribute specifies the minimum depth of a buffered pipe, where  $N$  is the number of data values. If the read and write kernels specify different depths for a given buffered pipe, the Intel FPGA SDK for OpenCL Offline Compiler uses the larger depth of the two.

```

__kernel void
producer (__global int *in_data,
          write_only pipe int __attribute__((blocking))
          __attribute__((depth(10))) c)
{
    for (i = 0; i < N; i++)
    {
        if (in_data[i])
        {
            write_pipe( c, &in_data[i] );
        }
    }
}

```

```

__kernel void
consumer ( __global int *check_data,
           __global int *out_data,
           read_only pipe int __attribute__((blocking)) c )
{
    int last_val = 0;
    for ( i = 0; i < N; i++)
    {
        if (check_data[i])
        {
            read_pipe( c, &last_val );
        }
        out_data[i] = last_val;
    }
}

```

In this example, the write operation can write ten data values to the pipe successfully. After the pipe is full, the write kernel returns failure until a read kernel consumes some of the data in the pipe.

Because the pipe read and write calls are conditional statements, the pipe might experience an imbalance between read and write calls. You may add a buffer capacity to the pipe to ensure that the `producer` and `consumer` kernels are decoupled. This step is particularly important if the `producer` kernel is writing data to the pipe when the `consumer` kernel is not reading from it.

### 5.5.5.6. Implementing I/O Pipes Using the `io` Attribute

Include an `io` attribute in your OpenCL pipe declaration to declare a special I/O pipe to interface with input or output features of an FPGA board. These features might include network interfaces, PCIe, cameras, or other data capture or processing devices or protocols.

In the Intel FPGA SDK for OpenCL channels extension, the `io("chan_id")` attribute specifies the I/O feature of an accelerator board with which a channel interfaces. The `chan_id` argument is the name of the I/O interface listed in the `board_spec.xml` file of your Custom Platform. The same I/O features can be used to identify I/O pipes.

Because peripheral interface usage might differ for each device type, consult your board vendor's documentation when you implement I/O pipes in your kernel program. Your OpenCL kernel code must be compatible with the type of data generated by the peripheral interfaces. If there is a difference in the byte ordering between the external I/O pipes and the kernel, the Intel FPGA SDK for OpenCL Offline Compiler converts the byte ordering seamlessly upon entry and exit.

#### **Caution:**

- Implicit data dependencies might exist for pipes that connect to the board directly and communicate with peripheral devices via I/O pipes. These implicit data dependencies might lead to compilation issues because the offline compiler cannot identify these dependencies.
  - External I/O pipes communicating with the same peripherals do not obey any sequential ordering. Ensure that the external device does not require sequential ordering because unexpected behavior might occur.
1. Consult the `board_spec.xml` file in your Custom Platform to identify the input and output features available on your FPGA board.



For example, a `board_spec.xml` file might include the following information on I/O features:

```
<channels>
  <interface name="udp_0" port="udp0_out" type="streamsource" width="256"
    chan_id="eth0_in"/>
  <interface name="udp_0" port="udp0_in" type="streamsink" width="256"
    chan_id="eth0_out"/>
  <interface name="udp_0" port="udp1_out" type="streamsource" width="256"
    chan_id="eth1_in"/>
  <interface name="udp_0" port="udp1_in" type="streamsink" width="256"
    chan_id="eth1_out"/>
</channels>
```

The `width` attribute of an `interface` element specifies the width, in bits, of the data type used by that pipe. For the example above, both the `uint` and `float` data types are 32 bits wide. Other bigger or vectorized data types must match the appropriate bit width specified in the `board_spec.xml` file.

2. Implement the `io` attribute as demonstrated in the following code example. The `io` attribute names must match those of the I/O channels (`chan_id`) specified in the `board_spec.xml` file.

```
__kernel void test (pipe uint pkt  __attribute__((io("enet"))),;
                  pipe float data __attribute__((io("pcie"))));
```

**Attention:** Declare a unique `io("chan_id")` handle for each I/O pipe specified in the channels XML element within the `board_spec.xml` file.

### 5.5.5.7. Enforcing the Order of Pipe Calls

To enforce the order of pipe calls, introduce memory fence or barrier functions in your kernel program to control memory accesses. A memory fence function is necessary to create a control flow dependence between the pipe synchronization calls before and after the fence.

When the Intel FPGA SDK for OpenCL Offline Compiler generates a compute unit, it does not create instruction-level parallelism on all instructions that are independent of each other. As a result, pipe read and write operations might not execute independently of each other even if there is no control or data dependence between them. When pipe calls interact with each other, or when pipes write data to external devices, deadlocks might occur.

For example, the code snippet below consists of a producer kernel and a consumer kernel. Pipes `c0` and `c1` are unbuffered pipes. The schedule of the pipe read operations from `c0` and `c1` might occur in the reversed order as the pipe write operations to `c0` and `c1`. That is, the producer kernel writes to `c0` but the consumer kernel might read from `c1` first. This rescheduling of pipe calls might cause a deadlock because the consumer kernel is reading from an empty pipe.

```
__kernel void producer (__global const uint * restrict src,
                      const uint iterations,
                      write_only pipe uint __attribute__((blocking)) c0,
                      write_only pipe uint __attribute__((blocking)) c1)
{
  for (int i = 0; i < iterations; i++) {
    write_pipe (c0, &src[2*i  ]);
    write_pipe (c1, &src[2*i+1]);
  }
}
```

```

__kernel void consumer (__global uint * restrict dst,
                        const uint iterations,
                        read_only_pipe uint __attribute__((blocking)) c0,
                        read_only_pipe uint __attribute__((blocking)) c1)
{
    for (int i = 0; i < iterations; i++) {
        read_pipe (c0, &dst[2*i+1]);
        read_pipe( c1, &dst[2*i]);
    }
}

```

- To prevent deadlocks from occurring by enforcing the order of pipe calls, include memory fence functions (`mem_fence`) in your kernel.

Inserting the `mem_fence` call with each kernel's pipe flag forces the sequential ordering of the write and read pipe calls. The code snippet below shows the modified producer and consumer kernels:

```

__kernel void producer (__global const uint * src,
                        const uint iterations,
                        write_only_pipe uint __attribute__((blocking)) c0,
                        write_only_pipe uint __attribute__((blocking)) c1)
{
    for (int i = 0; i < iterations; i++)
    {
        write_pipe(c0, &src[2*i  ]);
        mem_fence(CLK_CHANNEL_MEM_FENCE);
        write_pipe(c1, &src[2*i+1]);
    }
}

__kernel void consumer (__global uint * dst;
                        const uint iterations,
                        read_only_pipe uint __attribute__((blocking)) c0,
                        read_only_pipe uint __attribute__((blocking)) c1)
{
    for(int i = 0; i < iterations; i++)
    {
        read_pipe(c0, &dst[2*i  ]);
        mem_fence(CLK_CHANNEL_MEM_FENCE);
        read_pipe(c1, &dst[2*i+1]);
    }
}

```

In this example, `mem_fence` in the producer kernel ensures that the pipe write operation to `c0` occurs before that to `c1`. Similarly, `mem_fence` in the consumer kernel ensures that the pipe read operation from `c0` occurs before that from `c1`.

#### 5.5.5.7.1. Defining Memory Consistency Across Kernels When Using Pipes

According to the OpenCL Specification version 2.0, memory behavior is undefined unless a kernel completes execution. A kernel must finish executing before other kernels can visualize any changes in memory behavior. However, kernels that use pipes can share data through common global memory buffers and synchronized memory accesses. To ensure that data written to a pipe is visible to the read pipe after execution passes a memory fence, define memory consistency across kernels with respect to memory fences.

- To create a control flow dependency between the pipe synchronization calls and the memory operations, add the `CLK_GLOBAL_MEM_FENCE` flag to the `mem_fence` call.

For example:

```
__kernel void producer (__global const uint * restrict src,
                       const uint iterations,
                       write_only pipe uint __attribute__((blocking)) c0,
                       write_only pipe uint __attribute__((blocking)) c1)
{
    for (int i = 0; i < iterations; i++)
    {
        write_pipe(c0, &src[2*i]);
        mem_fence(CLK_CHANNEL_MEM_FENCE | CLK_GLOBAL_MEM_FENCE);
        write_pipe(c1, &src[2*i+1]);
    }
}
```

In this kernel, the `mem_fence` function ensures that the write operation to `c0` and memory access to `src[2*i]` occur before the write operation to `c1` and memory access to `src[2*i+1]`. This allows data written to `c0` to be visible to the read pipe before data is written to `c1`.

### 5.5.6. Direct Communication with Kernels via Host Pipes

The `cl_intel_fpga_host_pipe` extension enables point-to-point pipe communication between a kernel and the host program. Without the extension, pipes within OpenCL can only be used to communicate between kernels, and not with the host program directly.

*Note:* A BSP with support for host pipes is required to use the `cl_intel_fpga_host_pipe` extension. Without the BSP support, the `clCreatePipe` call fails. You can use the `a10gx_hostpipe` board variant in the `a10_ref` BSP for Intel Arria 10 GX Development Kit with host pipes. If you are using a custom BSP, consult your board vendor for support for host pipes.

The extension provides two new values in the `flags` argument of `clCreatePipe` to make a pipe host accessible, and adds four new API functions (`clReadPipeIntelFPGA`, `clWritePipeIntelFPGA`, `clMapHostPipeIntelFPGA`, and `clUnmapHostPipeIntelFPGA`) to allow the host to read from and write to a pipe that was created with host access enabled. A new optional kernel argument attribute is added to specify in the kernel language that the opposing end of a pipe kernel argument is the host program, and consequently that the pipe is not connected to another kernel. A pipe kernel argument is specialized in the kernel definition to connect to either a host pipe or another kernel, and cannot dynamically switch between the two at runtime.

When a pipe kernel argument is marked for host accessibility, the kernel language pipe accessors are restricted to a subset of the 2.x functions (reservations are not supported), and memory consistency or visibility guarantees are made beyond OpenCL synchronization points.

Support for host accessible pipes is a device property, advertised as `cl_intel_fpga_host_pipe`.

**Note:**

- Data written to a host pipe is eventually made visible on the read side without any OpenCL synchronization point.
- A restriction of our implementation of host pipes is that the platform only supports two host pipes. One for read and one for write. Furthermore, the compiler accepts a pipe of only 32-bytes width, and hence `ulong4` is used in the *Example Use of `cl_intel_fpga_host_pipe` Extension* section.
- Host programs using the `cl_intel_fpga_host_pipe` extension must include the `CL/cl_ext_intelfpga.h` header file and enable OpenCL 2.0 support, as described in [Support Statuses of OpenCL 2.0 Features](#) on page 217.

**Related Information**

[Example Use of the `cl\_intel\_fpga\_host\_pipe` Extension](#) on page 81

**5.5.6.1. Optional `intel_host_accessible` Kernel Argument Attribute**

The `cl_intel_fpga_host_pipe` extension adds an optional `intel_host_accessible` kernel argument attribute. Applying this attribute to a kernel pipe argument specifies that the host connects the kernel pipe argument to a host-accessible pipe, and not to another kernel pipe argument.

```
__attribute__((intel_host_accessible))
```

**5.5.6.2. API Functions for Interacting with `cl_mem` Pipe Objects Bound to Host-Accessible Pipe Kernel Arguments**

Including the `clReadPipeIntelFPGA`, `clWritePipeIntelFPGA`, `clMapHostPipeIntelFPGA`, and `clUnmapHostPipeIntelFPGA` host API functions allows the host program to read from and write to the `cl_mem` pipe objects that have been bound (using `clSetKernelArg` argument) to host-accessible pipe kernel arguments.

- `clReadPipeIntelFPGA` and `clWritePipeIntelFPGA` functions operate on single words of the pipe's width.
- `clMapHostPipeIntelFPGA` function is an advanced mechanism to reduce latency and overhead when performing many word reads or writes on a host pipe.
- `clUnmapHostPipeIntelFPGA` function allows the host program to signal to the OpenCL runtime that it has written to or read from either a portion of or the entire mapped region that was created through a previous `clMapHostPipeIntelFPGA` function call.

The following sections describe the API functions for bound `cl_mem` objects:

**`clReadPipeIntelFPGA` Function**

This function reads a data packet from a pipe with the following characteristics:

- Created with the `CL_MEM_HOST_READ_ONLY` flag.
- Bound to a kernel argument that has the `write_only` definition and the `intel_host_accessible` kernel argument attribute.

Each `clReadPipeIntelFPGA` function call reads one packet from the pipe. The operation is non-blocking. It does not wait until data is available in the pipe to successfully read before returning.

### Syntax

```
cl_int clReadPipeIntelFPGA (cl_mem pipe,  
                           gentype *ptr);
```

### Return Values

The `clReadPipeIntelFPGA` function returns `CL_SUCCESS` if the read is successful. Otherwise, it returns one of the following errors:

**Table 4. Return Values for `clReadPipeIntelFPGA` Function**

Return Value	Description
<code>CL_INVALID_MEM_OBJECT</code>	Pipe was not created with the <code>clCreatePipe</code> function or if the <code>CL_MEM_HOST_READ_ONLY</code> flag is not used when creating the pipe.
<code>CL_INVALID_KERNEL</code>	Pipe is not bound to a kernel argument using the <code>clSetKernelArg</code> function.
<code>CL_INVALID_VALUE</code>	The <code>ptr</code> attribute is <code>NULL</code> .
<code>CL_PIPE_EMPTY</code>	Pipe is empty with no data, and therefore there were no valid packets to read.

### `clWritePipeIntelFPGA` Function

This function writes a data packet to a pipe with the following characteristics:

- Created with the `CL_MEM_HOST_WRITE_ONLY` flag.
- Bound to a kernel argument that has the `read_only` definition and the `intel_host_accessible` kernel argument attribute.

Each `clWritePipeIntelFPGA` function call writes one packet to the pipe. The operation is non-blocking. It does not wait until there is a capacity in the pipe to successfully write before returning. A return status of `CL_SUCCESS` does not imply that data is available to the kernel for reading. Data eventually is available for reading by the kernel, assuming that any previously mapped buffers on the host pipe are unmapped.

### Syntax

```
cl_int clWritePipeIntelFPGA (cl_mem pipe,  
                             gentype *ptr);
```

### Return Values

The `clWritePipeIntelFPGA` function returns `CL_SUCCESS` if the write is successful. Otherwise, it returns one of the following errors:

**Table 5. Return Values for `clWritePipeIntelFPGA` Function**

Return Value	Description
<code>CL_INVALID_MEM_OBJECT</code>	Pipe is not created with <code>clCreatePipe</code> function or if the <code>CL_MEM_HOST_WRITE_ONLY</code> flag is not used when creating the pipe.
<code>CL_INVALID_KERNEL</code>	Pipe is not bound to a kernel argument using the <code>clSetKernelArg</code> function.
<code>CL_INVALID_VALUE</code>	The <code>ptr</code> attribute is <code>NULL</code> .
<code>CL_PIPE_FULL</code>	Pipe is already full of data and cannot accept further packets without one being read by a kernel. The packet capacity of the pipe is specified as an argument to the <code>clCreatePipe</code> function.

### `clMapHostPipeIntelFPGA` Function

This function returns a `void *` in the host address space. The pipe can write data to this address space if it was created with the `CL_MEM_HOST_WRITE_ONLY` flag. The pipe can read data from this address space if it was created with the `CL_MEM_HOST_READ_ONLY` flag.

The `mapped_size` argument specifies the maximum number of bytes that the host can access, as determined by the runtime in the memory. The value specified by `mapped_size` argument is less than or equal to the value of the `requested_size` argument that the caller specifies.

After writing to or reading from the returned `void *`, the host must execute one or more `clUnmapHostPipeIntelFPGA` function calls to signal to the runtime that data is ready for transfer to the device (on a write) and that the runtime can reclaim the memory for reuse (on a read or write). If the `clMapHostPipeIntelFPGA` function is called before the `clUnmapHostPipeIntelFPGA` function unmaps all memory mapped by a previous `clMapHostPipeIntelFPGA` function call, the buffer returned by the second `clMapHostPipeIntelFPGA` call does not overlap with that returned by the first call.

### Syntax

```
void * clMapHostPipeIntelFPGA (cl_mem pipe,
                              cl_map_flags map_flags,
                              size_t requested_size,
                              size_t * mapped_size,
                              cl_int * errcode_ret);
```

### Return Values

The `clMapHostPipeIntelFPGA` function returns a valid non-NULL address and `CL_SUCCESS` is returned in `errcode_ref` if the host pipe is mapped successfully. Otherwise, it returns `NULL` and `errcode_ret` is set with one of the following values:

**Table 6. Error Values Set in `errcode_ret`**

Description	Error Values
This error is returned when the: <ul style="list-style-type: none"> <li>• Pipe is not created with the <code>clCreatePipe</code> function.</li> <li>• <code>CL_MEM_HOST_WRITE_ONLY</code> or <code>CL_MEM_HOST_READ_ONLY</code> flag is not used when creating a pipe.</li> <li>• Pipe has been released.</li> </ul>	<code>CL_INVALID_MEM_OBJECT</code>
Pipe is not bound to a kernel argument using the <code>clSetKernelArg</code> function.	<code>CL_INVALID_KERNEL</code>
The <code>requested_size</code> attribute is not a multiple of the packet size that is specified to the <code>clCreatePipe</code> function.	<code>CL_INVALID_VALUE</code>
The <code>mapped_size</code> attribute is <code>NULL</code> .	
Runtime cannot provide buffer space for the host to read or write. This may occur when a new mapping is requested before data is transferred to or from the device from a previous <code>clUnmapHostPipeIntelFPGA</code> function call, so attempting to map again may succeed.	<code>CL_OUT_OF_RESOURCES</code>
Runtime is unable to allocate the host addressable memory.	<code>CL_OUT_OF_HOST_MEMORY</code>

### `clUnmapHostPipeIntelFPGA` Function

This function signals to the runtime that the host no longer uses `size_to_unmap` bytes of a host-addressable buffer that the `clMapHostPipeIntelFPGA` function has returned previously. In the case of a writeable host pipe, calling the `clUnmapHostPipeIntelFPGA` function allows the unmapped data to become available to the kernel. If the `size_to_unmap` value is smaller than the `mapped_size` value specified by the `clMapHostPipeIntelFPGA` function, then multiple `clUnmapHostPipeIntelFPGA` function calls are necessary to unmap full capacity of the buffer. You may include multiple `clUnmapHostPipeIntelFPGA` function calls to unmap successive bytes in the buffer returned by a `clMapHostPipeIntelFPGA` function call, up to the `mapped_size` value defined by the `clMapHostPipeIntelFPGA` call.

### Syntax

```
cl_int clUnmapHostPipeIntelFPGA (cl_mem pipe,
                                void * mapped_ptr,
                                size_t size_to_unmap,
                                size_t * unmapped_size);
```

### Return Values

The `clUnmapHostPipeIntelFPGA` function returns `CL_SUCCESS` on successful unmapping. Otherwise, it returns the following error:

**Table 7. Return Values for `clUnmapHostPipeIntelFPGA` Function**

Return Value	Description
CL_INVALID_VALUE	The <code>mapped_ptr</code> attribute is not a valid pointer returned by a <code>clMapHostPipeIntelFPGA</code> function call associated with the pipe, or if the <code>mapped_ptr</code> has already been fully unmapped.
	The <code>requested_size</code> attribute is not a multiple of the packet size that was specified to the <code>clCreatePipe</code> function.
	The <code>requested_size</code> attribute is larger than the remaining unmapped bytes in the <code>mapped_ptr</code> buffer.
	The <code>unmapped_size</code> attribute is <code>NULL</code> .

### Recommendations for Using Host Pipes Efficiently

To use host pipes efficiently, Intel recommends the following:

- Use map or unmap with multiple packets rather than single packet with read or write.
- Use two threads to simultaneously map and unmap.
- Ensure kernel can read or write data to host pipe in every clock cycle.

### Recommendations for Avoiding Hangs, Stalls or Deadlocks

To avoid hangs, stalls, or deadlocks situation, Intel recommends the following:

- Packets may get lost when enqueueing kernels from different `cl_programs`. Use a single `cl_program` object in the host application.
- Ensure that the same amount of data is written and read on kernel and host.

### 5.5.6.3. Creating a Host Accessible Pipe

The `clCreatePipe` function, defined in the *Section 5.4.1* of the *OpenCL 2.2 API Specification*, contains a `flags` parameter. The legal values of `flags` for `clCreatePipe` function are `CL_MEM_READ_WRITE` and `CL_MEM_HOST_NO_ACCESS`. If the value passed to `flags` is 0, then the specification defines that both of these flags are implicitly passed as the default.

To enable host access (reading or writing) to pipes, the `cl_intel_fpga_host_pipe` extension legalizes the following two flags values to `clCreatePipe`:

- `CL_MEM_HOST_READ_ONLY`
- `CL_MEM_HOST_WRITE_ONLY`

When one of these flags is passed to the `clCreatePipe` function, the corresponding `cl_mem` object can be passed as the first argument to `clReadPipeIntelFPGA` and `clWritePipeIntelFPGA` functions. Throughout the remainder of the `cl_intel_fpga_host_pipe` extension, such a pipe is referred to as a host pipe.



**Warning:** It is illegal to specify both `CL_MEM_HOST_READ_ONLY` and `CL_MEM_HOST_WRITE_ONLY` on the same pipe, or to mix either of those values with `CL_MEM_READ_WRITE` or `CL_MEM_HOST_NO_ACCESS`, or both. Invalid flags combinations are detected by the OpenCL runtime, and causes `clCreatePipe` to return the `CL_INVALID_VALUE` error.

#### 5.5.6.4. Example Use of the `cl_intel_fpga_host_pipe` Extension

The following are the example kernel and host codes of the `cl_intel_fpga_host_pipe` extension:

##### Kernel Code

```
#pragma OPENCL EXTENSION cl_intel_fpga_host_pipe : enable

kernel void reader(__attribute__((intel_host_accessible))
                  __read_only pipe ulong4 host_in) {
    ulong4 val;
    if (read_pipe(host_in, &val)) {
        ....
    }
    ....
}

kernel void writer(__attribute__((intel_host_accessible))
                  __write_only pipe ulong4 device_out) {
    ulong4 val;
    ....
    if (write_pipe(device_out, &val)) {
        ....
    }
}
```

##### Host Code

```
....

cl_kernel read_kern = clCreateKernel(program, "reader", NULL);
cl_kernel write_kern = clCreateKernel(program, "writer", NULL);

cl_mem read_pipe = clCreatePipe(context, CL_MEM_HOST_READ_ONLY,
                                sizeof( cl_ulong4 ), 128,
                                // Number of packets that can be buffered
                                NULL, &error);

cl_mem write_pipe = clCreatePipe(context, CL_MEM_HOST_WRITE_ONLY,
                                 sizeof( cl_ulong4 ), 64,
                                 // Number of packets that can be buffered
                                 NULL, &error);

// Bind pipes to kernels
clSetKernelArg(read_kern, 0, sizeof(cl_mem), (void *)&write_pipe);
clSetKernelArg(write_kern, 0, sizeof(cl_mem), (void *)&read_pipe);

// Enqueue kernels
....

cl_ulong4 val
if (!clReadPipeIntelFPGA (read_pipe, &val)) {
    cl_int result = clWritePipeIntelFPGA (write_pipe, &val);
    // Check write success/failure and handle
    ....
}

....
```

## 5.6. Implementing Arbitrary Precision Integers

Use the Intel FPGA SDK for OpenCL arbitrary precision integer extension to define integers with a custom bit-width. You can define integer custom bit-widths up to and including 64 bits.

To use the arbitrary precision integer extension, include the following line in your list of header files in your kernel code:

```
#include "ihc_apint.h"
```

When you compile a kernel that includes the `ihc_apint.h` header file, you must include the `-I $INTELFPGAOCSDKROOT/include/kernel_headers` option with the `aoc` command. For example:

```
aoc <other command options> -I $INTELFPGAOCSDKROOT/include/kernel_headers  
<my_kernel_file>
```

The header defines signed and unsigned arbitrary precision integer type definitions of the form `int $d$ _t` and `uint $d$ _t`, where  $d$  is the bit width of the type.

For example, you can declare a 10-bit signed and unsigned arbitrary precision integers as follows:

```
int10_t x_signed;  
uint10_t x_unsigned;
```

You can declare arbitrary precision integers with widths up to 64 bits.

If you do operations where the bit width of the result is larger than the bit widths of the arguments, you must explicitly cast one of the arguments to the resulting bit width.

For example, if you had the following operation, the result overflows the declared size of the integer:

```
int10_t a;  
int10_t b;  
int20_t res;  
  
res = a * b;
```

In the example, the compiler attempts to instantiate a multiplier that multiplies two 10-bit integers and put the results into another 10-bit integer. The result is then sign extended or zero extended up to 20-bits.

To prevent the overflow, explicitly cast one of the arguments to the resulting bit width as follows:

```
res = ((int20_t)a) * b
```

**Remember:** When you compile a program for x86-64 platforms, the bit widths for arbitrary precision integers are rounded up to either 32 bits or 64 bits. When you compile a kernel for an FPGA platform, the bit widths are not rounded up and the arbitrary precision integers remain at their declared bit width.

As a result, an operation that appears to work correctly in an x86-64 program can overflow and lose precision when you compile that same operation in an FPGA kernel. The additional precision provided by bit-width rounding on x86-64 platforms masks possible overflow and precision-loss problems you might encounter when you compile your FPGA kernel.

## 5.7. Using Predefined Preprocessor Macros in Conditional Compilation

You may take advantage of predefined preprocessor macros that allow you to conditionally compile portions of your kernel code.

To introduce Intel FPGA SDK for OpenCL Offline Compiler version-specific code and optimizations, structure your kernel program in the following manner:

```
#if INTELFPGA_CL >= 191
    // use new features added in 19.1
#else
    // do things the old way
#endif
```

Where `INTELFPGA_CL` is the Intel predefined preprocessor macro set to the Intel FPGA SDK for OpenCL Offline Compiler version. This macro enables you to version your code based on the compiler version.

### Related Information

[Defining Preprocessor Macros to Specify Kernel Parameters \(-D<macro\\_name>\)](#) on page 125

## 5.8. Declaring `__constant` Address Space Qualifiers

There are several limitations and workarounds you must consider when you include `__constant` address space qualifiers in your kernel.

### Function Scope `__constant` Variables

The Intel FPGA SDK for OpenCL Offline Compiler does not support function scope `__constant` variables. Replace function scope `__constant` variables with file scope constant variables. You can also replace function scope `__constant` variables with `__constant` buffers that the host passes to the kernel.

### File Scope `__constant` Variables

If the host always passes the same constant data to your kernel, consider declaring that data as a constant preinitialized file scope array within the kernel file. Declaration of a constant preinitialized file scope array creates a ROM directly in the hardware to store the data. This ROM is available to all work-items in the NDRange.

For example:

```
__constant int my_array[8] = {0x0, 0x1, 0x2, 0x3, 0x4, 0x5, 0x6, 0x7};

__kernel void my_kernel (__global int * my_buffer)
{
    size_t gid = get_global_id(0);
    my_buffer[gid] += my_array[gid % 8];
}
```

In this case, the offline compiler sets the values for `my_array` in a ROM because the file scope constant data does not change between kernel invocations.

### Pointers to `__constant` Parameters from the Host

You can replace file scope constant data with a pointer to a `__constant` parameter in your kernel code if the data is not fixed across kernel invocations. You must then modify your host application in the following manner:

1. Create `cl_mem` memory objects associated with the pointers in global memory.
2. Load constant data into `cl_mem` objects with `clEnqueueWriteBuffer` prior to kernel execution.
3. Pass the `cl_mem` objects to the kernel as arguments with the `clSetKernelArg` function.

For simplicity, if a constant variable is of a complex type, use a `typedef` argument, as shown in the table below:

**Table 8. Replacing File Scope `__constant` Variable with Pointer to `__constant` Parameter**

If your source code is structured as follows:	Rewrite your code to resemble the following syntax:
<pre>__constant int Payoff[2][2] = {{ 1, 3}, {5, 3}}; __kernel void original(__global int * A) {     *A = Payoff[1][2];     // and so on }</pre>	<pre>__kernel void modified(__global int * A, __constant Payoff_type * PayoffPtr ) {     *A = (PayoffPtr)[1][2];     // and so on }</pre>

**Attention:** Use the same type definition in both your host application and your kernel.

## 5.9. Including Structure Data Types as Arguments in OpenCL Kernels

Pass structure parameters (`struct`) in OpenCL kernels either by value or as a pointer to a structure.

**Attention:** The `__global struct` argument points to a buffer that must be created in the host program to store the structure. To prevent pointer aliasing, include a `restrict` qualifier in the declaration of the pointer to the structure.

### 5.9.1. Matching Data Layouts of Host and Kernel Structure Data Types

If you use structure data types (`struct`) as arguments in OpenCL kernels, match the member data types and align the data members between the host application and the kernel code.

To match member data types, use the `cl_` version of the data type in your host application that corresponds to the data type in the kernel code. The `cl_` version of the data type is available in the `opencl.h` header file. For example, if you have a data member of type `float4` in your kernel code, the corresponding data member you declare in the host application is `cl_float4`.

Align the structures and align the `struct` data members between the host and kernel applications. Manage the alignments carefully because of the variability among different host compilers.

For example, if you have `float4` OpenCL data types in the struct, the alignments of these data items must satisfy the OpenCL specification (that is, 16-byte alignment for `float4`).

The following rules apply when the Intel FPGA SDK for OpenCL Offline Compiler compiles your OpenCL kernels:

1. Alignment of built-in scalar and vector types follow the rules outlined in Section 6.1.5 of the *OpenCL Specification version 1.0*.

The offline compiler usually aligns a data type based on its size. However, the compiler aligns a value of a three-element vector the same way it aligns a four-element vector.

2. An array has the same alignment as one of its elements.
3. A `struct` (or a `union`) has the same alignment as the maximum alignment necessary for any of its data members.

Consider the following example:

```
struct my_struct
{
    char data[3];
    float4 f4;
    int index;
};
```

The offline compiler aligns the `struct` elements above at 16-byte boundaries because of the `float4` data type. As a result, both `data` and `index` also have 16-byte alignment boundaries.

4. The offline compiler does not reorder data members of a `struct`.
5. Normally, the offline compiler inserts a minimum amount of data structure padding between data members of a `struct` to satisfy the alignment requirements for each data member.
  - a. In your OpenCL kernel code, you may specify data packing (that is, no insertion of data structure padding) by applying the `packed` attribute to the `struct` declaration. If you impose data packing, ensure that the alignment of data members satisfies the OpenCL alignment requirements. The Intel FPGA SDK for OpenCL does not enforce these alignment requirements. Ensure that your host compiler respects the kernel attribute and sets the appropriate alignments.
  - b. In your OpenCL kernel code, you may specify the amount of data structure padding by applying the `aligned(N)` attribute to a data member, where  $N$  is the amount of padding. The SDK does not enforce these alignment requirements. Ensure that your host compiler respects the kernel attribute and sets the appropriate alignments.

For Windows systems, some versions of the Microsoft Visual Studio compiler pack structure data types by default. If you do not want to apply data packing, specify an amount of data structure padding as shown below:

```
struct my_struct
{
    __declspec(align(16)) char data[3];

    /*Note that cl_float4 is the only known float4 definition on the
    host*/
    __declspec(align(16)) cl_float4 f4;

    __declspec(align(16)) int index;
};
```

*Tip:* An alternative way of adding data structure padding is to insert dummy struct members of type char or array of char.

#### Related Information

- [Modifying Host Program for Structure Parameter Conversion](#) on page 104
- [OpenCL Specification version 1.0](#)

### 5.9.2. Disabling Insertion of Data Structure Padding

You may instruct the Intel FPGA SDK for OpenCL Offline Compiler to disable automatic padding insertion between members of a struct data structure.

- To disable automatic padding insertion, insert the `packed` attribute prior to the kernel source code for a struct data structure.

For example:

```
struct __attribute__((packed)) Context
{
    float param1;
    float param2;
    int param3;
    uint param4;
};
__kernel void algorithm(__global float * restrict A,
                       __global struct Context * restrict c)
{
    if ( c->param3 )
    {
        // Dereference through a pointer and so on
    }
}
```

For more information, refer to the *Align a Struct with or without Padding* section of the *Intel FPGA SDK for OpenCL Best Practices Guide*.

#### Related Information

[Align a Struct with or without Padding](#)

### 5.9.3. Specifying the Alignment of a Struct

You may instruct the Intel FPGA SDK for OpenCL Offline Compiler to set a specific alignment of a struct data structure.

- To specify the struct alignment, insert the `aligned(N)` attribute prior to the kernel source code for a struct data structure.

For example:

```
struct __attribute__((aligned(2))) Context
{
    float param1;
    float param2;
    int param3;
    uint param4;
};
__kernel void algorithm(__global float * A,
                       __global struct Context * restrict c)
{
    if ( c->param3 )
    {
        // Dereference through a pointer and so on
    }
}
```

For more information, refer to the *Align a Struct with or without Padding* section of the *Intel FPGA SDK for OpenCL Best Practices Guide*.

#### Related Information

[Align a Struct with or without Padding](#)

## 5.10. Inferring a Register

The Intel FPGA SDK for OpenCL Offline Compiler can implement data that is in the private address space in registers or in block RAMs. In general, the offline compiler chooses registers if the access to a variable is fixed and does not require any dynamic indexes. Accessing an array with a variable index usually forces the array into block RAMs. Implementing private data as registers is beneficial for data accesses that should occur in a single cycle (for example, feedback in a single work-item loop).

The offline compiler infers private arrays as registers either as single values or in a piecewise fashion. Piecewise implementation results in very efficient hardware; however, the offline compiler must be able to determine data accesses statically. To facilitate piecewise implementation, hardcode the access points into the array. You can also facilitate register inference by unrolling loops that access the array.

If array accesses are not inferable statically, the offline compiler might infer the array as registers. However, the offline compiler limits the size of these arrays to 64 bytes in length for single work-item kernels. There is effectively no size limit for kernels with multiple work-items.

Consider the following code example:

```
int array[SIZE];
for (int j = 0; j < N; ++j)
{
    for (int i = 0; i < SIZE - 1; ++i)
    {
        array[i] = array[i + 1];
    }
}
```

The indexing into `array[i]` is not inferable statically because the loop is not unrolled. If the size of `array[SIZE]` is less than or equal to 64 bytes for single work-item kernels, the offline compiler implements `array[SIZE]` into registers as a single value. If the size of `array[SIZE]` is greater than 64 bytes for single work-item kernels, the offline compiler implements the entire array in block RAMs. For multiple work-item kernels, the offline compiler implements `array[SIZE]` into registers as a single value provided that its size is less than 1 kilobyte (KB).

### 5.10.1. Inferring a Shift Register

The shift register design pattern is a very important design pattern for efficient implementation of many applications on the FPGA. However, the implementation of a shift register design pattern might seem counterintuitive at first.

Consider the following code example:

```
channel int in, out;

#define SIZE 512
//Shift register size must be statically determinable

__kernel void foo()
{
    int shift_reg[SIZE];
    //The key is that the array size is a compile time constant

    // Initialization loop
    #pragma unroll
    for (int i=0; i < SIZE; i++)
    {
        //All elements of the array should be initialized to the same value
        shift_reg[i] = 0;
    }

    while(1)
    {
        // Fully unrolling the shifting loop produces constant accesses
        #pragma unroll
        for (int j=0; j < SIZE-1; j++)
        {
            shift_reg[j] = shift_reg[j + 1];
        }
        shift_reg[SIZE - 1] = read_channel_intel(in);

        // Using fixed access points of the shift register
        int res = (shift_reg[0] + shift_reg[1]) / 2;

        // 'out' channel will have running average of the input channel
        write_channel_intel(out, res);
    }
}
```

In each clock cycle, the kernel shifts a new value into the array. By placing this shift register into a block RAM, the Intel FPGA SDK for OpenCL Offline Compiler can efficiently handle multiple access points into the array. The shift register design pattern is ideal for implementing filters (for example, image filters like a Sobel filter or time-delay filters like a finite impulse response (FIR) filter).



When implementing a shift register in your kernel code, keep in mind the following key points:

1. Unroll the shifting loop so that it can access every element of the array.
2. All access points must have constant data accesses. For example, if you write a calculation in nested loops using multiple access points, unroll these loops to establish the constant access points.
3. Initialize all elements of the array to the same value. Alternatively, you may leave the elements uninitialized if you do not require a specific initial value.
4. If some accesses to a large array are not inferable statically, they force the offline compiler to create inefficient hardware. If these accesses are necessary, use `__local` memory instead of `__private` memory.
5. Do not shift a large shift register conditionally. The shifting must occur in very loop iteration that contains the shifting code to avoid creating inefficient hardware.

### 5.11. Enabling Double Precision Floating-Point Operations

The Intel FPGA SDK for OpenCL offers preliminary support for all double precision floating-point functions.

Before declaring any double precision floating-point data type in your OpenCL kernel, include the following `OPENCL_EXTENSION` pragma in your kernel code:

```
#pragma OPENCL_EXTENSION cl_khr_fp64 : enable
```

### 5.12. Single-Cycle Floating-Point Accumulator for Single Work-Item Kernels

Single work-item kernels that perform accumulation in a loop can take advantage of the single-cycle floating-point accumulator feature of the Intel FPGA SDK for OpenCL Offline Compiler. The offline compiler searches for these kernel instances and attempts to map an accumulation that executes in a loop into the accumulator structure.

The offline compiler supports an accumulator that adds or subtracts a value. To use this feature, describe the accumulation in a way that allows the offline compiler to infer the accumulator.

- Attention:**
- The accumulator is only available on Intel Arria 10 devices.
  - The accumulator must be part of a loop.
  - The accumulator must have an initial value of 0.
  - The accumulator cannot be conditional.

Below are examples of a description that results in the correct inference of the accumulator by the offline compiler.

```
channel float4 RANDOM_STREAM;

__kernel void acc_test(__global float *a, int k) {
    // Simplest example of an accumulator.
    // In this loop, the accumulator acc is incremented by 5.
    int i;
    float acc = 0.0f;
    for (i = 0; i < k; i++) {
        acc+=5;
    }
}
```

```

    }
    a[0] = acc;
}

__kernel void acc_test2(__global float *a, int k) {
    // Extended example showing that an accumulator can be
    // conditionally incremented. The key here is to describe the increment
    // as conditional, not the accumulation itself.
    int i;
    float acc = 0.0f;
    for (i = 0; i < k; i++) {
        acc += ((i < 30) ? 5 : 0);
    }
    a[0] = acc;
}

__kernel void acc_test3(__global float *a, int k) {
    // A more complex case where the accumulator is fed
    // by a dot product.
    int i;
    float acc = 0.0f;
    for (i = 0; i < k; i++) {
        float4 v = read_channel_intel(RANDOM_STREAM);
        float x1 = v.x;
        float x2 = v.y;
        float y1 = v.z;
        float y2 = v.w;

        acc += (x1*y1+x2*y2);
    }
    a[0] = acc;
}

__kernel void loader(__global float *a, int k) {
    int i;
    float4 my_val = 0;
    for(i = 0; i < k; i++) {
        if ((i%4) == 0)
            write_channel_intel(RANDOM_STREAM, my_val);
        if ((i%4) == 0) my_val.x = a[i];
        if ((i%4) == 1) my_val.y = a[i];
        if ((i%4) == 2) my_val.z = a[i];
        if ((i%4) == 3) my_val.w = a[i];
    }
}

```

### 5.12.1. Programming Strategies for Inferring the Accumulator

To take advantage of the single cycle floating-point accumulator feature, you can modify the accumulator description in your kernel code to improve efficiency or work around programming restrictions.

#### Describing an Accumulator Using Multiple Loops

Consider a case where you want to describe an accumulator using multiple loops, with some of the loops being unrolled:

```

float acc = 0.0f;
for (i = 0; i < k; i++) {
    #pragma unroll
    for(j=0; j < 16; j++)
        acc += (x[i+j]*y[i+j]);
}

```

In this situation, it is important to compile the kernel with the `-ffp-reassociate` Intel FPGA SDK for OpenCL Offline Compiler command option to enable the offline compiler to rearrange the operations in a way that exposes the accumulation. If you do not compile the kernel with `-ffp-reassociate`, the resulting accumulator structure will have a high initiation interval (II). II is the number of cycles between launching successive loop iterations. The higher the II value, the longer the accumulator structure must wait before it can process the next loop iteration.

### Modifying a Multi-Loop Accumulator Description

In cases where you cannot compile an accumulator description using the `-ffp-reassociate` offline compiler command option rewrite the code to expose the accumulation.

For the code example above, rewrite it in the following manner:

```
float acc = 0.0f;
for (i = 0; i < k; i++) {
    float my_dot = 0.0f;
    #pragma unroll
    for(j=0;j < 16; j++)
        my_dot += (x[i+j]*y[i+j]);
    acc += my_dot;
}
```

### Modifying an Accumulator Description Containing a Variable or Non-Zero Initial Value

Consider a situation where you might want to apply an offset to a description of an accumulator that begins with a non-zero value:

```
float acc = array[0];
for (i = 0; i < k; i++) {
    acc += x[i];
}
```

Because the accumulator hardware does not support variable or non-zero initial values in a description, you must rewrite the description.

```
float acc = 0.0f;
for (i = 0; i < k; i++) {
    acc += x[i];
}
acc += array[0];
```

Rewriting the description in the above manner enables the kernel to use an accumulator in a loop. The loop structure is then followed by an increment of `array[0]`.

## 5.13. Integer Promotion Rules

The rules of integer promotion applied when you use `intX_t` data types are different from standard C/C++ rules. Your kernel design should account for these differing rules.

- If both operands are of standard integer type (for example `char` or `short`), integers are promoted following the C/C++ standard. That is, the operation is carried out in the data type and size of the largest operand, with at least 32 bits. The expression returns the result in that larger data type.
- If both operands are `intX_t` data types, operations are carried out in the largest `intX_t` data type even if that data type is smaller than 32 bits. The expression returns the result in that type.
- If the expression has one standard data type and one `intX_t` data type, the rules for `intX_t` data type promotion apply. The resulting expression type is always an `intX_t` data type. For example, if the largest data type is a standard integer type `short`, the resulting data type is an `int16_t`.
- In C/C++, literals are by default an `int` data type, so when you use a literal without any casting, the expression type is always at least 32 bits. For example, if you have code as shown in the following snippet, the comparison is carried out in 32 bits:

```
int5_t ap;
...
if (ap < 4) {
...

```

- If operands are of different signage and the unsigned type is at least as large as the other type, the operation is carried out as an unsigned operation. Otherwise, the unsigned operand is converted to a signed value.

For example, if you have code as shown in the following snippet, `-1` expands to a 32-bit negative number (`0xffffffff`) while the `uint3_t` ends up as the positive 32-bit number `7` (`0x00000007`), which are not equal.

```
uint3_t x = 7;
if (x != -1) {
    // FAIL
}
```

## 6. Designing Your Host Application

---

Intel offers guidelines on host requirements and procedures on structuring the host application. If applicable, implement these design strategies when you create or modify a host application for your OpenCL kernels.

[Host Programming Requirements](#) on page 93

[Allocating OpenCL Buffers for Manual Partitioning of Global Memory](#) on page 94

[Triggering Collection Profiling Data During Kernel Execution](#) on page 99

[Accessing Custom Platform-Specific Functions](#) on page 103

[Modifying Host Program for Structure Parameter Conversion](#) on page 104

[Managing Host Application](#) on page 105

[Allocating Shared Memory for OpenCL Kernels Targeting SoCs](#) on page 115

[Sharing Multiple Devices Across Multiple Host Programs](#) on page 117

### 6.1. Host Programming Requirements

When designing your OpenCL host application for use with the Intel FPGA SDK for OpenCL, ensure that the application satisfies the following host programming requirements.

#### 6.1.1. Host Machine Memory Requirements

The machine that runs the host application must have enough host memory to support several components simultaneously.

The host machine must support the following components:

- The host application and operating system.
- The working set for the host application.
- The maximum amount of OpenCL memory buffers that can be allocated at once. Every device-side `cl_mem` buffer is associated with a corresponding storage area in the host process. Therefore, the amount of host memory necessary might be as large as the amount of external memory supported by the FPGA.

#### 6.1.2. Host Binary Requirement

When compiling the host application, target one of these architectures: x86-64 (64-bit) or ARM® 32-bit ARMV7-A for SoCs. The Intel FPGA SDK for OpenCL host runtime does not support x86-32 (32-bit) binaries.

### 6.1.3. Multiple Host Threads

The Intel FPGA SDK for OpenCL host library is thread-safe.

All OpenCL APIs are thread safe except the `clSetKernelArg` function.

It is safe to call `clSetKernelArg` from any host thread or in a reentrant way as long as concurrent calls to `clSetKernelArg` operate on different `cl_kernel` objects.

#### Related Information

[Multi-Threaded Host Application](#)

### 6.1.4. Out-of-order Command Queues

The OpenCL host runtime command queues support out-of-order command execution. Out-of-order command queues can be used to execute commands both sequentially (if event dependencies are specified), or concurrently (if no event dependencies are specified).

### 6.1.5. Requirement for Multiple Command Queues to Execute Kernels Concurrently

To execute kernels within the same OpenCL program object concurrently, instantiate a separate in-order command queue for each kernel you want to run concurrently. This can also be achieved by using a single out-of-order command queue and enqueueing independent kernels (no dependencies between each other).

A single in-order command queue can only dispatch a single operation for execution. Subsequent operations are not dispatched until the previous operation is fully complete. Thus, to execute kernels within the same OpenCL program object concurrently, instantiate a separate command queue for each kernel you want to run concurrently.

Similarly, multiple in-order command queues are also required to concurrently execute different transfers (`clEnqueueReadBuffer` or `clEnqueueWriteBuffer`), including executing transfers concurrently with kernels. For example, to execute kernel A concurrently with kernel B as well as concurrently with a `clEnqueueWriteBuffer` transfer, you must create three command queues and enqueue each of the operations in a separate queue. Achieving this in a steady state leads to maximum utilization of the FPGA device.

Out-of-order command queues may also be used to launch buffer writes, reads, and kernel execution concurrently. Enqueueing events that do not have any dependencies on one another into an out-of-order queue schedules them to execute concurrently, with the exception that these events are not blocked by their other dependencies.

## 6.2. Allocating OpenCL Buffers for Manual Partitioning of Global Memory

Manual partitioning of global memory buffers allows you to control memory accesses across buffers to maximize the memory bandwidth. You can partition buffers across interfaces of the same memory type or across interfaces of different memory types.

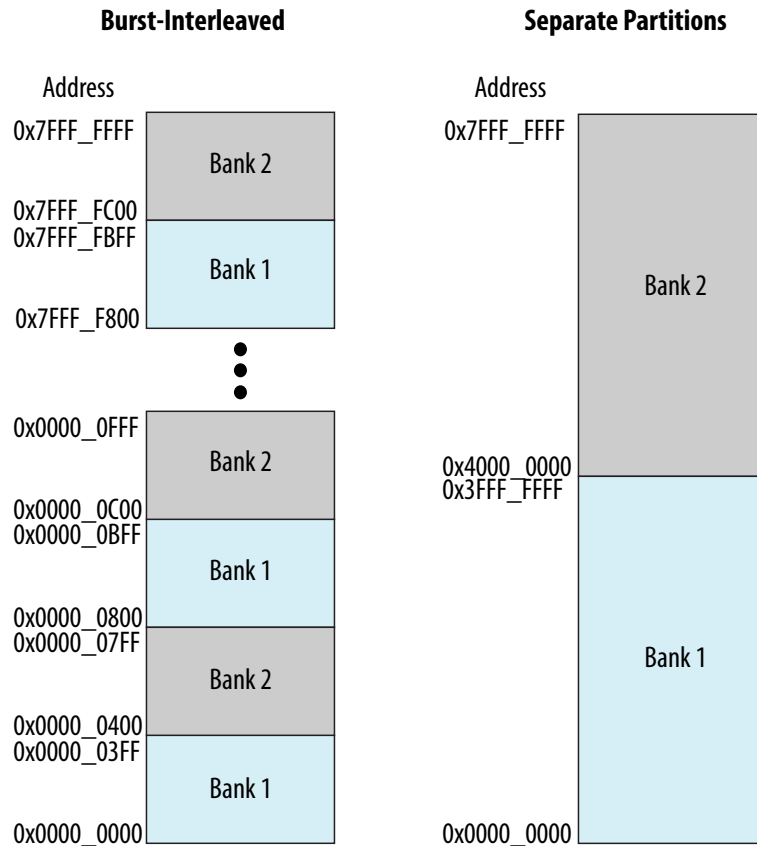
### 6.2.1. Partitioning Buffers Across Multiple Interfaces of the Same Memory Type

Before you partition the memory across multiple interfaces of the same memory type, you must first disable burst-interleaving during OpenCL kernel compilation. Then, in the host application, you must specify the memory bank to which you allocate the OpenCL buffer.

*Tip:* For oneAPI SYCL-specific instructions, refer to [Global Memory Accesses Optimization](#) topic in the *FPGA Optimization Guide for Intel oneAPI Toolkits*.

By default, the Intel FPGA SDK for OpenCL Offline Compiler configures each global memory type in a burst-interleaved fashion. Usually, the burst-interleaving configuration leads to the best load balancing between the memory banks. However, there might be situations where it is more efficient to partition the memory into non-interleaved regions.

The figure below illustrates the differences between burst-interleaved and non-interleaved memory partitions.



To manually partition some or all of the available global memory types, perform the following tasks:

1. Compile your OpenCL kernel using the `-no-interleaving=<global_memory_type>` flag to configure the memory bank(s) of the specified memory type as separate addresses.

For more information about the use of the `-no-interleaving=<global_memory_type>` flag, refer to the *Disabling Burst-Interleaving of Global Memory (-no-interleaving=<global\_memory\_type>)* section.

2. Create an OpenCL buffer in your host application, and allocate the buffer to one of the banks using the `CL_CHANNEL` flags.
  - Specify `CL_CHANNEL_1_INTELFPGA` to allocate the buffer to the lowest available memory region.
  - Specify `CL_CHANNEL_2_INTELFPGA` to allocation memory to the second bank (if available).

**Attention:** Allocate each buffer to a single memory bank only. If the second bank is not available at runtime, the memory is allocated to the first bank. If no global memory is available, the `clCreateBuffer` call fails with the error message `CL_MEM_OBJECT_ALLOCATION_FAILURE`.

### Related Information

[Disabling Burst-Interleaving of Global Memory \(-no-interleaving=<global\\_memory\\_type>\)](#) on page 129

## 6.2.2. Partitioning Buffers Across Different Memory Types (Heterogeneous Memory)

The board support package for your FPGA board can assemble a global memory space consisting of different memory technologies (for example, DRAM or SRAM). The board support package designates one such memory, which might consist of multiple interfaces, as the default memory. All buffers reside there.

*Tip:* For oneAPI SYCL-specific instructions, refer to [Global Memory Accesses Optimization](#) topic in the *FPGA Optimization Guide for Intel oneAPI Toolkits*.

To use the heterogeneous memory, modify the code in your `.cl` file as follows:

1. Determine the names of the global memory types available on your FPGA board in one of the following ways:
  - Refer to the board vendor's documentation for more information.
  - Find the names in the `board_spec.xml` file of your board Custom Platform. For each global memory type, the name is the unique string assigned to the `name` attribute of the `global_mem` element.
2. To instruct the host to allocate a buffer to a specific global memory type, insert the `buffer_location("<memory_type>")` attribute, where `<memory_type>` is the name of the global memory type provided by your board vendor.

For example:

```
__kernel void foo(__global __attribute__((buffer_location("DDR"))) int *x,
                 __global __attribute__((buffer_location("QDR"))) int *y)
```

If you do not specify the `buffer_location` attribute, the host allocates the buffer to the default memory type automatically. To determine the default memory type, consult the documentation provided by your board vendor. Alternatively, in the `board_spec.xml` file of your Custom Platform, search for the memory type that is defined first or has the attribute `default=1` assigned to it.



Intel recommends that you define the `buffer_location` attribute in a preprocessor macro for ease of reuse, as follows:

```
#define QDR\  
__global __attribute__((buffer_location("QDR")))  
  
#define DDR\  
__global __attribute__((buffer_location("DDR")))  
__kernel void foo (QDR uint * data, DDR uint * lup)  
{  
    //statements  
}
```

**Attention:** If you assign a kernel argument to a non-default memory (for example, `QDR uint * data` and `DDR uint * lup` from the code above), you cannot declare that argument using the `constant` keyword. In addition, you cannot perform atomic operations with pointers derived from that argument.

By default, the host allocates buffers into the main memory when you load kernels into the OpenCL runtime via the `clCreateProgramWithBinary` function. During kernel invocation, the host automatically relocates heterogeneous memory buffers that are bound to kernel arguments to the main memory.

3. To avoid the initial allocation of heterogeneous memory buffers in the main memory, include the `CL_MEM_HETEROGENEOUS_INTELFPGA` flag when you call the `clCreateBuffer` function. Also, bind the `cl_mem` buffer to the argument that used the `buffer_location` attribute using `clSetKernelArg` before doing any reads or writes from that buffer, as follows:

```
mem = clCreateBuffer(context, flags|CL_MEM_HETEROGENEOUS_INTELFPGA,  
                    memSize, NULL, &errNum);  
  
clSetKernelArg(kernel, 0, sizeof(cl_mem), &mem);  
clEnqueueWriteBuffer(queue, mem, CL_FALSE, 0, N, 0, NULL, &write_event);  
clEnqueueNDRangeKernel(queue, kernel, 1, NULL, global_work_size, NULL,  
                        0, NULL, &kernel_event);
```

For example, the following `clCreateBuffer` call allocates memory into the lowest available memory region of a nondefault memory bank:

```
mem = clCreateBuffer(context, (CL_MEM_HETEROGENEOUS_INTELFPGA|  
CL_CHANNEL_1_INTELFPGA),  
                    memSize, NULL, &errNum);
```

**Note:** Host programs using `CL_MEM_HETEROGENEOUS_INTELFPGA` and `CL_CHANNEL_*_INTELFPGA` flags must include the `CL/cl_ext_intelfpga.h` header file.

The `clCreateBuffer` call allocates memory into a certain global memory type based on what you specify in the kernel argument. If a memory (`cl_mem`) object residing in a memory type is set as a kernel argument that corresponds to a different memory technology, the host moves the memory object automatically when it queues the kernel. Do not pass a buffer as kernel arguments that associate it with multiple memory technologies.

For more information about optimizing heterogeneous global memory accesses, refer to the *Heterogeneous Memory Buffers* and the *Manual Partitioning of Global Memory* sections of the *Intel FPGA SDK for OpenCL Best Practices Guide*.

### Related Information

- [Manual Partitioning of Global Memory](#)
- [Heterogeneous Memory Buffers](#)

## 6.2.3. Creating a Pipe Object in Your Host Application

To implement OpenCL pipes in your kernel, you must create Intel FPGA SDK for OpenCL-specific pipe objects in your host application.

An SDK-specific pipe object is not a true OpenCL pipe object as described in the OpenCL Specification version 2.0. This implementation allows you to migrate away from Intel FPGA products with a conformant solution. The SDK-specific pipe object is a memory object (`cl_mem`); however, the host does not allocate any memory for the pipe itself.

The following `clCreatePipe` host API creates a pipe object:

```
cl_mem clCreatePipe(cl_context context,
                   cl_mem_flags flags,
                   cl_uint pipe_packet_size,
                   cl_uint pipe_max_packets,
                   const cl_pipe_properties *properties,
                   cl_int *errcode_ret)
```

For more information about the `clCreatePipe` host API function, refer to section 5.4.1 of the *OpenCL Specification version 2.0*.

Below is an example syntax of the `clCreatePipe` host API function:

```
cl_int status;
cl_mem c0_pipe = clCreatePipe(context,
                              0,
                              sizeof(int),
                              1,
                              NULL,
                              &status);
status = clSetKernelArg(kernel, 1, sizeof(cl_mem), &c0_pipe);
```

**Caution:** The SDK does not support dynamic channel assignment at runtime. The SDK statically links the pipes during compilation based on matching names.

### Related Information

[OpenCL Specification version 2.0 \(API\)](#)

## 6.2.4. Enabling All Global Memory

By default, the OpenCL runtime reserves a small portion of the memory at the start of the device's memory to facilitate `NULL`-pointer checking. This ensures that valid global memory pointers always have a non-zero value and allows a zero value to be interpreted as a `NULL` pointer.

To enable 100% of the device's memory to be used, including the reserved memory, set the `CL_ALLOW_GLOBAL_MEM_AT_NULL_ADDRESS_INTELFPGA` environment variable to 1.

**Note:** If the `CL_ALLOW_GLOBAL_MEM_AT_NULL_ADDRESS_INTELFPGA` environment variable is set to 1, a valid global memory pointer can be 0 when observed by the kernel. In other words, you cannot reliably check inside the kernel that a NULL global memory pointer is invalid.

### 6.3. Triggering Collection Profiling Data During Kernel Execution

The Intel FPGA dynamic profiler for OpenCL can be used to collect performance information from the hardware when the design is executed. For instructions about how to add the profiler to your hardware design and how to view the collected data, refer to Chapter 5 of *Intel FPGA SDK for OpenCL Best Practices Guide*.

In cases where kernel execution finishes after the host application completes and temporal profiling is disabled, you can query the FPGA explicitly to collect profile data during kernel execution.

**Note:** All information in this section are written with the assumption that temporal profiling is not enabled. For more information about temporal profiling, refer to [Temporal Performance Collection](#) in the *Intel FPGA SDK for OpenCL Best Practices Guide*.

**Tip:** For oneAPI SYCL-specific instructions, refer to [Intel FPGA Dynamic Profiler for DPC++](#) topic in the *FPGA Optimization Guide for Intel oneAPI Toolkits*.

When you profile your OpenCL kernel during compilation, a `profile.mon` file is generated automatically. The profile data is then written to `profile.mon` after kernel execution completes on the FPGA. However, if kernel execution completes after the host application terminates, no profiling information for that kernel invocation is available in the `profile.mon` file. In this case, you can modify your host code to acquire profiling information during kernel execution.

**Important:** Collecting profile data during kernel execution may add overhead to kernel executions by increasing the latency in your kernel.

- To query the FPGA to collect profile data while the kernel is running, call the following host library call:

```
extern CL_API_ENTRY cl_int CL_API_CALL
clGetProfileInfoIntelFPGA(cl_event);
```

where `cl_event` is the kernel event. The kernel event you pass to this host library call must be the same one you pass to the `clEnqueueNDRangeKernel` call.

- Important:** — If kernel execution completes before the invocation of `clGetProfileInfoIntelFPGA`, the function returns an event error message.
- Host programs that use `clGetProfileInfoIntelFPGA` and `clGetProfileDataDeviceIntelFPGA` function calls must include the `CL/cl_ext_intelfpga.h` header file.

**Caution:** Invoking the `clGetProfileInfoIntelFPGA` function during kernel execution disables the profile counters momentarily so that the Intel FPGA dynamic profiler for OpenCL can collect data from the FPGA. As a result, you lose some profiling information during this interruption. If you call this function at very short intervals, the profile data might not accurately reflect the actual performance behavior of the kernel.

Consider the following example host code:

```
int main()
{
    ...
    clEnqueueNDRangeKernel(queue, kernel, ..., NULL);
    ...
    clEnqueueNDRangeKernel(queue, kernel, .., NULL);
    ...
}
```

This host application runs on the assumption that a kernel launches twice and then completes. In the `profile.mon` file, there are two sets of profile data, one for each kernel invocation. To collect profile data while the kernel is running, modify the host code in the following manner:

```
int main()
{
    ...
    clEnqueueNDRangeKernel(queue, kernel, ..., &event);

    //Get the profile data before the kernel completes
    clGetProfileInfoIntelFPGA(event);

    //Wait until the kernel completes
    clFinish(queue);

    ...
    clEnqueueNDRangeKernel(queue, kernel, ..., NULL);
    ...
}
```

The call to `clGetProfileInfoIntelFPGA` adds a new entry in the `profile.mon` file.

### Related Information

[Profile Your Kernel to Identify Performance Bottlenecks](#)

### 6.3.1. Profiling Autorun Kernels

**Note:** All information in this section are written with the assumption that temporal profiling is not enabled. For more information about temporal profiling, refer to [Temporal Performance Collection](#) in the *Intel FPGA SDK for OpenCL Best Practices Guide*.

Unlike enqueued kernels that automatically generate profiler data on completion (if the compiler flag is set), autorun kernels never complete. Hence, you must explicitly indicate when to profile kernels by calling the `clGetProfileDataDeviceIntelFPGA` host library call. All profiler data is output to a `profile.mon` file. Data collected by the host library call is a snapshot of the autorun profile data.

Following is the code snippet of the `clGetProfileDataDeviceIntelFPGA` host library call:

```
cl_int clGetProfileDataDeviceIntelFPGA (cl_device_id device_id,
                                       cl_program program,
                                       cl_bool read_enqueue_kernels,
                                       cl_bool read_auto_enqueued,
                                       cl_bool clear_counters_after_readback,
                                       size_t param_value_size,
                                       void *param_value,
                                       size_t *param_value_size_ret,
                                       cl_int *errcode_ret);
```

where,

- `read_enqueue_kernels` parameter profiles enqueued kernels. In this release, this parameter has no effect.
- `read_auto_enqueued` parameter profiles autorun kernels.
- Following are the placeholder parameters for the future releases:
  - `clear_counters_after_readback`
  - `param_value_size`
  - `param_value`
  - `param_value_size_ret`
  - `errcode_ret`

**Notice:** Only autorun kernels are supported by this host library call. You can enter `TRUE` for the `read_enqueue_kernels` parameter, but the boolean is ignored. This does not mean that enqueued kernels are not profiled. If the compiler `profile` flag is set to include enqueued kernels, the profile data is captured normally at the end of execution. The only difference is that the `clGetProfileDataDeviceIntelFPGA` host library call does not profile enqueued kernels in addition to the profiling already done automatically for the enqueued kernels.

The `clGetProfileDataDeviceIntelFPGA` host library call returns `CL_SUCCESS` on success. Else, it returns one of the following errors:

- `CL_INVALID_DEVICE` if the device is not a valid device.
- `CL_INVALID_PROGRAM` if the program is not a valid program.

**Caution:** The `clGetProfileDataDeviceIntelFPGA` host library call does not trigger a programming operation of the provided program on the provided device. If the program is not already programmed to the device at the time of the host library call, then the host library call returns `CL_INVALID_PROGRAM` error.

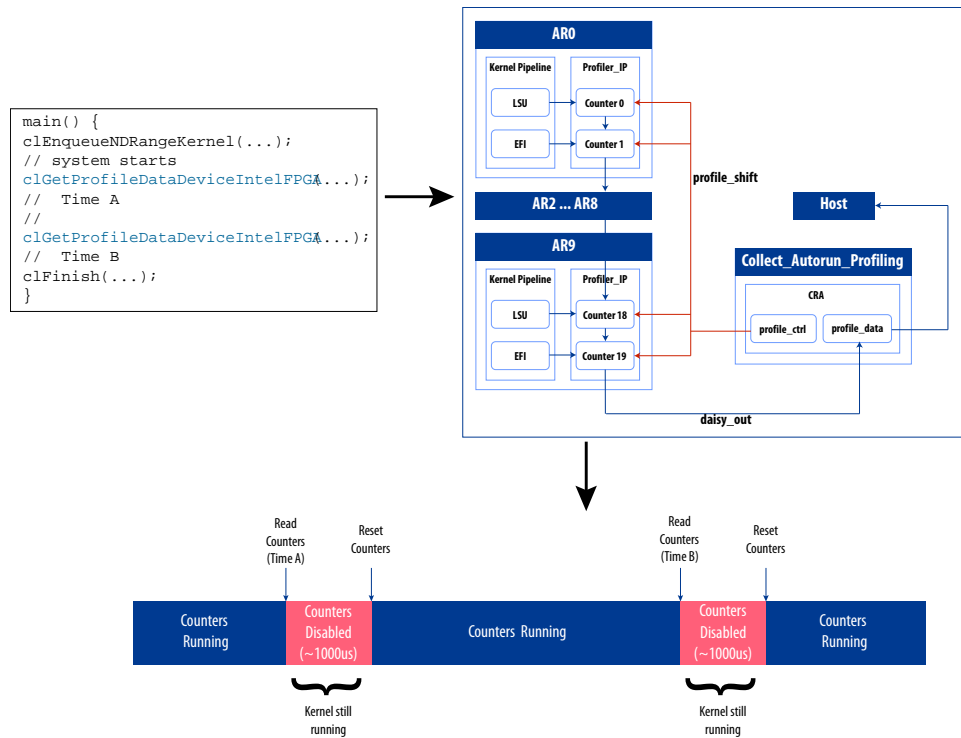
**Table 9. `clGetProfileDataDeviceIntelFPGA` Host Library Call Parameter Combinations**

	<code>read_auto_enqueued</code>
Profile only enqueued kernels <i>Note:</i> Automatically outputs profile information once the execution is completed.	False
Profile only autorun kernels	True
Profile both enqueued and autorun kernels	True

### 6.3.1.1. Multiple Autorun Profiling Calls

Because autorun kernels run continuously, the host application can include multiple `clGetProfileDataDeviceIntelFPGA` host library calls to profile autorun kernels at certain points within the execution or in specific time ranges. Every time the host application calls the `clGetProfileDataDeviceIntelFPGA` host library call, it reads the profile counters and then resets them to zero. Calling `clGetProfileDataDeviceIntelFPGA` multiple times allows the host application to profile autorun kernels over time ranges.

**Figure 13. Multiple Autorun Profile Captures Flow**



### 6.3.2. Profile Data Acquisition

Profile data acquisition from running kernels is paused during read back operations.

**Attention:** Although data acquisition is paused, kernels themselves are still running. Therefore, during read back operations, no kernel data is recorded.

Pausing data acquisition is not synchronized exactly across all kernels. The skew between halting profile data acquisition across kernels is dependent on the communication link with the device, driver overhead, and congestion on communication buses. Exact synchronized snapshotting of profile data between kernels should not be relied upon.

### 6.4. Accessing Custom Platform-Specific Functions

You have the option to include in your application user-accessible functions that are available in your Custom Platform. However, when you link your host application to the FPGA Client Driver (FCD), you cannot directly reference these Custom Platform-specific functions. To reference Custom Platform-specific user-accessible functions while linking to the FCD, include the `clGetBoardExtensionFunctionAddressIntelFPGA` extension in your host application.

The `clGetBoardExtensionFunctionAddressIntelFPGA` extension specifies an API that retrieves a pointer to a user-accessible function from the Custom Platform.

**Attention:** For Linux systems, the `clGetBoardExtensionFunctionAddressIntelFPGA` function works with or without FCD. For Windows systems, the function only works in conjunction with FCD. Consult with your board vendor to determine if FCD is supported in your Custom Platform.

Definitions of the extension interfaces are available in the `INTELFPGAOCLSDKROOT/host/include/CL/cl_ext_intelfpga.h` file.

- To obtain a pointer to a user-accessible function in your Custom Platform, call the following function in your host application:

```
void* clGetBoardExtensionFunctionAddressIntelFPGA (
    const char* function_name,
    cl_device_id device
);
```

Where:

*function\_name* is the name of the user-accessible function that your Custom Platform vendor provides,

and

*device* is the device ID returned by the `clGetDeviceIDs` function.

After locating the user-accessible function, the `clGetBoardExtensionFunctionAddressIntelFPGA` function returns a pointer to the user-accessible function. If the function does not exist in the Custom Platform, `clGetBoardExtensionFunctionAddressIntelFPGA` returns `NULL`.

**Attention:** To access the `clGetBoardExtensionFunctionAddressIntelFPGA` API via the Installable Client Driver (ICD), ensure that the ICD extension API `clGetExtensionFunctionAddressIntelFPGA` retrieves the pointer to the `clGetBoardExtensionFunctionAddressIntelFPGA` API first.

The following code example shows how you can access the Custom Platform-specific function via ICD:

```
clGetBoardExtensionFunctionAddressIntelFPGA_fn
clGetBoardExtensionFunctionAddressIntelFPGA =
    (clGetBoardExtensionFunctionAddressIntelFPGA_fn)
    clGetExtensionFunctionAddressForPlatform
    (platform, "clGetBoardExtensionFunctionAddressIntelFPGA");

if (clGetBoardExtensionFunctionAddressIntelFPGA == NULL){
    printf ("Failed to get
           clGetBoardExtensionFunctionAddressIntelFPGA\n");
}

void * board_extension_function_ptr =
clGetBoardExtensionFunctionAddressIntelFPGA("function_name",device_id);
```

### Related Information

- [Linking Your Host Application to the Khronos ICD Loader Library](#) on page 107
- [OpenCL Installable Client Driver \(ICD\) Loader](#)

## 6.5. Modifying Host Program for Structure Parameter Conversion

If you convert any structure parameters to pointers-to-constant structures in your OpenCL kernel, you must modify your host application accordingly.

Perform the following changes to your host application:

- Allocate a `cl_mem` buffer to store the structure contents.



**Attention:** You need a separate `cl_mem` buffer for every kernel that uses a different structure value.

2. Set the structure kernel argument with a pointer to the structure buffer, not with a pointer to the structure contents.
3. Populate the structure buffer contents before queuing the kernel. Perform one of the following steps to ensure that the structure buffer is populated before the kernel launches:
  - Queue the structure buffer on the same command queue as the kernel queue.
  - Synchronize separate kernel queues and structure buffer queues with an event.
4. When your application no longer needs to call a kernel that uses the structure buffer, release the `cl_mem` buffer.

#### Related Information

- [Including Structure Data Types as Arguments in OpenCL Kernels](#) on page 84
- [Matching Data Layouts of Host and Kernel Structure Data Types](#) on page 84

## 6.6. Managing Host Application

The Intel FPGA SDK for OpenCL includes utility commands you can invoke to obtain information on flags and libraries necessary for compiling and linking your host application.

**Attention:** To cross-compile your host application to an SoC FPGA board, include the `--arm` option in your utility command.

**Caution:** For Linux systems, if you debug your host application using the GNU Project Debugger (GDB), invoke the following command prior to running the host application:

```
handle SIG44 nostop
```

Without this command, the GDB debugging process terminates with the following error message:

```
Program received signal SIG44, Real-time event 44.
```

### 6.6.1. Displaying Example Makefile Fragments (example-makefile or makefile)

To display example Makefile fragments for compiling and linking a host application against host runtime libraries available with the Intel FPGA SDK for OpenCL, invoke the `example-makefile` or `makefile` utility command.

- At a command prompt, invoke the `aocl example-makefile` or `aocl makefile` utility command.

The software displays an output similar to the following:

The following are example Makefile fragments for compiling and linking a host program against the host runtime libraries included with the Intel® FPGA SDK for OpenCL™.

Example GNU makefile on Linux, with GCC toolchain:

```
AOCL_COMPILE_CONFIG=$(shell aocl compile-config)
AOCL_LINK_CONFIG=$(shell aocl link-config)

host_prog : host_prog.o
    g++ -o host_prog host_prog.o $(AOCL_LINK_CONFIG)

host_prog.o : host_prog.cpp
    g++ -c host_prog.cpp $(AOCL_COMPILE_CONFIG)
```

Example GNU makefile on Windows, with Microsoft Visual C++ command line compiler:

```
AOCL_COMPILE_CONFIG=$(shell aocl compile-config)
AOCL_LINK_CONFIG=$(shell aocl link-config)

host_prog.exe : host_prog.obj
    link -nologo /OUT:host_prog.exe host_prog.obj $(AOCL_LINK_CONFIG)

host_prog.obj : host_prog.cpp
    cl /MD /Fohost_prog.obj -c host_prog.cpp $(AOCL_COMPILE_CONFIG)
```

Example GNU makefile cross-compiling to ARM SoC from Linux or Windows, with Linaro GCC cross-compiler toolchain:

```
CROSS-COMPILER=arm-linux-gnueabi-
AOCL_COMPILE_CONFIG=$(shell aocl compile-config --arm)
AOCL_LINK_CONFIG=$(shell aocl link-config --arm)

host_prog : host_prog.o
    $(CROSS-COMPILER)g++ -o host_prog host_prog.o $(AOCL_LINK_CONFIG)

host_prog.o : host_prog.cpp
    $(CROSS-COMPILER)g++ -c host_prog.cpp $(AOCL_COMPILE_CONFIG)
```

## 6.6.2. Compiling and Linking Your Host Application

The OpenCL host application uses standard OpenCL runtime APIs to manage device configuration, data buffers, kernel launches, and synchronization. The host application also contains functions such as file I/O, or portions of the source code that do not run on an accelerator device. The Intel FPGA SDK for OpenCL includes utility commands you can invoke to obtain information on C header files describing the OpenCL APIs, and board-specific MMD and host runtime libraries with which you must link your host application.

**Important:** For Windows systems, you must add the `/MD` flag to link the host runtime libraries against the multithreaded dynamic link library (DLL) version of the Microsoft C Runtime library. You must also compile your host application with the `/MD` compilation flag, or use the `/NODEFAULTLIB` linker option to override the selection of runtime library.

*Remember:* Include the path to the `INTELFPGAOCSDKROOT/host/<OS_platform>/bin` folder in your library search path when you run your host application.

[Linking Your Host Application to the Khronos ICD Loader Library](#) on page 107

[Displaying Flags for Compiling Host Application \(compile-config\)](#) on page 109

[Displaying Paths to OpenCL Host Runtime and MMD Libraries \(ldflags\)](#) on page 109

[Listing OpenCL Host Runtime and MMD Libraries \(ldlibs\)](#) on page 109

[Displaying Information on OpenCL Host Runtime and MMD Libraries \(link-config or linkflags\)](#) on page 109

### 6.6.2.1. Linking Your Host Application to the Khronos ICD Loader Library

The Intel FPGA SDK for OpenCL supports the OpenCL ICD extension from the Khronos Group. The OpenCL ICD extension allows you to have multiple OpenCL implementations on your system. With the OpenCL ICD Loader Library, you may choose from a list of installed platforms and execute OpenCL API calls that are specific to your OpenCL implementation of choice.

In addition to the SDK's host runtime libraries, Intel supplies a version of the ICD Loader Library that supports the OpenCL Specification version 1.0 and the implemented APIs from the OpenCL Specification versions 1.1, 1.2, and 2.0. To use an ICD library from another vendor, consult the vendor's documentation on how to link to their ICD library.

Before linking your OpenCL host application to the ICD Loader Library, you must also set up the FCD for loading the board MMD libraries. If you have not set up the FCD yet, refer to [Managing an FPGA Board](#) on page 19 for more information.

Ensure that you have set up both ICD and FCD correctly. You can verify this by using the `aocl diagnose -icd-only` utility, which populates the corresponding ICD/FCDs and verifies if the libraries are registered in the system.

- If the output of the `aocl diagnose` utility displays `ICD diagnostics PASSED`, then when you build your host application, the host application automatically gets linked with the ICD Loader Libraries.
- If the `aocl diagnose` utility fails to detect ICD, follow these steps to verify the ICD setup:
  - For Windows system, open the `regedit` with administrator privilege and go to the Windows registry key **HKEY\_LOCAL\_MACHINE\SOFTWARE\Khronos\OpenCL\Vendors**. The value of **Name** should be `altera_icd.dll`. You can find this dynamic library file in `<INTELFPGAOCCLSDKROOT>/host/windows64/bin`. The **Type** should be **DWORD**, and the **Data** should be **00000000**. For example, the registry key should resemble the following:
 

```
HKEY_LOCAL_MACHINE\SOFTWARE\Khronos\OpenCL\Vendors]  
"alteracl_icd.dll"=dword:00000000
```
  - For Linux system, ensure that the file `/etc/OpenCL/vendors/Altera.icd` exists in the system and contains the text `libalteracl.so`.
- If the `aocl diagnose` utility fails to detect FCD, follow these steps for checking the FCD setup:
  - For Windows system, check the libraries in the registry key **HKEY\_LOCAL\_MACHINE\SOFTWARE\Intel\OpenCL\Boards** to determine if you have chosen to install for all users. Otherwise, check the registry key **HKEY\_CURRENT\_USER\SOFTWARE\Intel\OpenCL\Boards**. The value of **Name** should be `<path_to_the_mmd_library>`, and the **Data** should be a **DWORD** that is set to 0. For example, the registry key should resemble the following:
 

```
[HKEY_LOCAL_MACHINE\SOFTWARE\Intel\OpenCL\Boards]  
"c:\board_vendor_a\my_board_mmd.dll"=dword:00000000
```
  - For Linux system, ensure that the file `/opt/Intel/OpenCL/Boards/my_board.fcd` exists in the system and contains the name of the vendor-specific libraries (for example, `/data/board_vendor_a/libmy_board_mmd.so`).
 

**Attention:** • If your board vendor provides multiple libraries, you must double check that they are in a particular order. Consult with your board vendor to determine the correct order to load the libraries. You must list the libraries in the registry in their loading order.

    - For Intel Arria 10 SoC boards, when you build the SD flash card image for your Custom Platform, create an `Altera.icd` file containing the text `libalteracl.so`. Store the `Altera.icd` file in the `/etc/OpenCL/vendors` directory of your Custom Platform. Refer to [Building the Software and SD Card Image for the Intel Arria 10 SoC Development Kit Reference Platform](#) for more information.

### 6.6.2.2. Displaying Flags for Compiling Host Application (compile-config)

To display a list of flags necessary for compiling a host application, invoke the `compile-config` utility command.

1. At a command prompt, invoke the `aocl compile-config` utility command. The software displays the path to the folder or directory in which the OpenCL API header files reside. For example:
  - For Windows systems, the path is `-I%INTELFPGAOCCLSDKROOT%/host/include`
  - For Linux systems, the path is `-I$INTELFPGAOCCLSDKROOT/host/include` where `INTELFPGAOCCLSDKROOT` points to the location of the software installation.
2. Add this path to your C preprocessor.

**Attention:** In your host source, include the `opencl.h` OpenCL header file, located in the `INTELFPGAOCCLSDKROOT/host/include/CL` folder or directory.

### 6.6.2.3. Displaying Paths to OpenCL Host Runtime and MMD Libraries (ldflags)

To display the paths necessary for linking a host application to the OpenCL host runtime and MMD libraries, invoke the `ldflags` utility command.

- At a command prompt, invoke the `aocl ldflags` utility command. The software displays the paths for linking your host application with the OpenCL host runtime libraries that provide OpenCL platform and runtime APIs. The OpenCL host runtime libraries are available in the `INTELFPGAOCCLSDKROOT/host/<OS_platform>/lib` directory.

### 6.6.2.4. Listing OpenCL Host Runtime and MMD Libraries (ldlibs)

To display the names of the OpenCL host runtime and MMD libraries necessary for linking a host application, invoke the `ldlibs` utility command.

- At a command prompt, invoke the `aocl ldlibs` utility command. The software lists the OpenCL host runtime libraries residing in the `INTELFPGAOCCLSDKROOT/host/<OS_platform>/lib` directory. It also lists the Custom Platform-specific MMD libraries residing in the `/<board_family_name>/<OS_platform>/lib` directory of your Custom Platform.

*Note:* If you set up FCD correctly, the software does not list the MMD libraries.

- For Windows systems, the output is `OpenCL.lib`
- For Linux systems, the output is `-lOpenCL`.

### 6.6.2.5. Displaying Information on OpenCL Host Runtime and MMD Libraries (link-config or linkflags)

To display a list of flags necessary for linking a host application with OpenCL host runtime and MMD libraries, invoke the `link-config` or `linkflags` utility command.

This utility command combines the functions of the `ldflags` and `ldlibs` utility commands.

- At a command prompt, invoke the `aocl link-config` or `aocl linkflags` command.
  - For Windows systems, the output is `/libpath:%INTELFPGAOCCLSDKROOT%/host/windows64/lib OpenCL.lib`.
  - For Linux systems, the output is `-L/$INTELFPGAOCCLSDKROOT/host/[linux64|arm32]/lib -lOpenCL`

### 6.6.3. Using OpenCL ICD Extension APIs

When using ICD, if you are calling an OpenCL host API function that is not officially mentioned in the OpenCL specification, you must first call the `clGetExtensionFunctionAddress` function to obtain a function pointer to that host API function and then call the host API function using that function pointer. This applies to all non-standard calls that include `IntelFPGA` in their function name.

Consider the following example code snippet:

```
extern CL_API_ENTRY cl_int CL_API_CALL
clGetProfileDataDeviceIntelFPGA(
    cl_device_id          /*device_id*/,
    cl_program           /*program*/,
    cl_bool              /*read_enqueue_kernels*/,
    cl_bool              /*read_auto_enqueued*/,
    cl_bool              /*clear_counters_after_readback*/,
    size_t               /*param_value_size*/,
    void *               /*param_value*/,
    size_t *             /*param_value_size_ret*/,

    cl_int *             /*errcode_ret*/ );
```

Replace the following function call:

```
cl_int status = clGetProfileDataDeviceIntelFPGA (device, program, false,
                                                true, false, 0, NULL,
                                                NULL, NULL);
```

with code using the following syntax to define and load the function pointer:

```
typedef cl_int (*clGetProfileDataDevice_fn) (cl_device_id, cl_program,
                                             cl_bool, cl_bool, cl_bool,
                                             size_t, void *,
                                             size_t *, cl_int *);
clGetProfileDataDevice_fn get_profile_data_ptr = (clGetProfileDataDevice_fn)
clGetExtensionFunctionAddressForPlatform (platform,
"clGetProfileDataDeviceIntelFPGA");
```

and use the function pointer as the function call:

```
cl_int status = (get_profile_data_ptr) (device, program, false, true,
                                       false, 0, NULL, NULL, NULL);
```

### 6.6.4. Programming an FPGA via the Host

The Intel FPGA SDK for OpenCL Offline Compiler is an offline compiler that compiles kernels independently of the host application. To load the kernels into the OpenCL runtime, include the `clCreateProgramWithBinary` function in your host application.

**Caution:** If your host system consists of multiple processors, only one processor can access the FPGA at a given time. Consider an example where there are two host applications, corresponding to two processors, attempting to launch kernels onto the same FPGA at the same time. The second host application receives an error message indicating that the device is busy. The second host application cannot run until the first host application releases the OpenCL context.

1. Compile your OpenCL kernel with the offline compiler to create the `.aocx` file.
2. Include the `clCreateProgramWithBinary` function in your host application to create the `cl_program` OpenCL program objects from the `.aocx` file.
3. Include the `clBuildProgram` function in your host application to create the program executable for the specified device.

Below is an example host code on using `clCreateProgramWithBinary` to program an FPGA device:

```
size_t lengths[1];
unsigned char* binaries[1] = {NULL};
cl_int status[1];
cl_int error;
cl_program program;
const char options[] = "";

FILE *fp = fopen("program.aocx", "rb");
fseek(fp, 0, SEEK_END);
lengths[0] = ftell(fp);
binaries[0] = (unsigned char*)malloc(sizeof(unsigned char)*lengths[0]);
rewind(fp);
fread(binaries[0], lengths[0], 1, fp);
fclose(fp);

program = clCreateProgramWithBinary(context,
                                   1,
                                   device_list,
                                   lengths,
                                   (const unsigned char **)binaries,
                                   status,
                                   &error);
clBuildProgram(program, 1, device_list, options, NULL, NULL);
```

If the `clBuildProgram` function executes successfully, it returns `CL_SUCCESS`.

4. Create kernel objects from the program executable using the `clCreateKernelsInProgram` or `clCreateKernel` function.
5. Include the kernel execution function to instruct the host runtime to execute the scheduled kernel(s) on the FPGA.
  - To enqueue a command to execute an NDRange kernel, use `clEnqueueNDRangeKernel`.
  - To enqueue a single work-item kernel, use `clEnqueueTask`.

**Attention:** Intel recommends that you release an event object when it is not in use. The SDK keeps an event object live until you explicitly instruct it to release the event object. Keeping an unused event object live causes unnecessary memory usage.

To release an event object, call the `clReleaseEvent` function.

You can load multiple FPGA programs into memory, which the host then uses to reprogram the FPGA as required.

For more information about these OpenCL host runtime API calls, refer to the *OpenCL Specification version 1.0*.

### Related Information

[OpenCL Specification version 1.0](#)

#### 6.6.4.1. Programming Multiple FPGA Devices

If you install multiple FPGA devices in your system, you can direct the host runtime to program a specific FPGA device by modifying your host code.

*Important:* Linking your host application to FCD allows you to target multiple FPGA devices from different Custom Platforms. However, this feature has limited support for Custom Platforms that are compatible with SDK versions prior to 16.1.

You can present up to 128 FPGA devices to your system in the following manner:

- Multiple FPGA accelerator boards, each consisting of a single FPGA.
- Multiple FPGAs on a single accelerator board that connects to the host system via a PCIe switch.
- Combinations of the above.

The host runtime can load kernels onto each and every one of the FPGA devices. The FPGA devices can then operate in a parallel fashion.

### Related Information

[Accessing Custom Platform-Specific Functions](#) on page 103

##### 6.6.4.1.1. Probing the OpenCL FPGA Devices

The host must identify the number of OpenCL FPGA devices installed into the system.

1. To query a list of FPGA devices installed in your machine, invoke the `aocl diagnose` command.
2. To direct the host to identify the number of OpenCL FPGA devices, add the following lines of code to your host application:

```
//Get the platform
ciErrNum = clGetPlatformID(&cpPlatform);

//Get the devices
ciErrNum = clGetDeviceIDs(cpPlatform,
                          CL_DEVICE_TYPE_ALL,
                          0,
                          NULL,
                          &ciDeviceCount);
cdDevices = (cl_device_id * )malloc(ciDeviceCount * sizeof(cl_device_id));
ciErrNum = clGetDeviceIDs(cpPlatform,
                          CL_DEVICE_TYPE_ALL,
                          ciDeviceCount,
                          cdDevices,
                          NULL);
```

For example, on a system with two OpenCL FPGA devices, `ciDeviceCount` has a value of 2, and `cdDevices` contains a list of two device IDs (`cl_device_id`).



### Related Information

[Querying the Device Name of Your FPGA Board \(diagnose\)](#) on page 22

#### 6.6.4.1.2. Querying Device Information

You can direct the host to query information on your OpenCL FPGA devices.

- To direct the host to output a list of OpenCL FPGA devices installed into your system, add the following lines of code to your host application:

```
char buf[1024];
for (unsigned i = 0; i < ciDeviceCount; i++)
{
    clGetDeviceInfo(cdDevices[i], CL_DEVICE_NAME, 1023, buf, 0);
    printf("Device %d: '%s'\n", i, buf);
}
```

When you query the device information, the host lists your FPGA devices in the following manner: `Device <N>: <board_name>: <name_of_FPGA_board>`

Where:

- `<N>` is the device number.
- `<board_name>` is the board designation you use to target your FPGA device when you invoke the `aoc` command.
- `<name_of_FPGA_board>` is the advertised name of the FPGA board.

For example, if you have two identical FPGA boards on your system, the host generates an output that resembles the following:

```
Device 0: board_1: Stratix V FPGA Board
Device 1: board_1: Stratix V FPGA Board
```

**Note:** The `clGetDeviceInfo` function returns the board type (for example, `board_1`) that the Intel FPGA SDK for OpenCL Offline Compiler lists on-screen when you invoke the `aoc -list-boards` command. If your accelerator board contains more than one FPGA, each device is treated as a "board" and is given a unique name.

### Related Information

[Listing the Available FPGA Boards and Custom Platforms \(-list-boards and -list-board-packages\)](#) on page 17

#### 6.6.4.1.3. Loading Kernels for Multiple FPGA Devices

If your system contains multiple FPGA devices, you can create specific `cl_program` objects for each FPGA and load them into the OpenCL runtime.

The following host code demonstrates the usage of the `clCreateProgramWithBinary` and `createMultiDeviceProgram` functions to program multiple FPGA devices:

```
cl_program createMultiDeviceProgram(cl_context context,
                                    const cl_device_id *device_list,
                                    cl_uint num_devices,
                                    const char *aocx_name);

// Utility function for loading file into Binary String
//
unsigned char* load_file(const char* filename, size_t *size_ret)
```

```

{
    FILE *fp = fopen(aocx_name, "rb");
    fseek(fp, 0, SEEK_END);
    size_t len = ftell(fp);
    char *result = (unsigned char*)malloc(sizeof(unsigned char)*len);
    rewind(fp);
    fread(result, len, 1, fp);
    fclose(fp);
    *size_ret = len;
    return result;
}

//Create a Program that is compiled for the devices in the "device_list"
//
cl_program createMultiDeviceProgram(cl_context context,
                                   const cl_device_id *device_list,
                                   cl_uint num_devices,
                                   const char *aocx_name)
{
    printf("creating multi device program %s for %d devices\n",
           aocx_name, num_devices);
    const unsigned char **binaries =
        (const unsigned char**)malloc(num_devices*sizeof(unsigned char*));
    size_t *lengths=(size_t*)malloc(num_devices*sizeof(size_t));
    cl_int err;

    for(cl_uint i=0; i<num_devices; i++)
    {
        binaries[i] = load_file(aocx_name, &lengths[i]);
        if (!binaries[i])
        {
            printf("couldn't load %s\n", aocx_name);
            exit(-1);
        }
    }

    cl_program p = clCreateProgramWithBinary(context,
                                             num_devices,
                                             device_list,
                                             lengths,
                                             binaries,
                                             NULL,
                                             &err);

    free(lengths);
    free(binaries);

    if (err != CL_SUCCESS)
    {
        printf("Program Create Error\n");
    }
    return p;
}

// main program
main ()
{
    // Normal OpenCL setup
}
program = createMultiDeviceProgram(context,
                                   device_list,
                                   num_devices,
                                   "program.aocx");
clBuildProgram(program, num_devices, device_list, options, NULL, NULL);

```

### 6.6.5. Termination of the Runtime Environment and Error Recovery

In the event that the host application terminates unexpectedly, you must restart the runtime environment and reprogram the FPGA.

The runtime environment is a library that is compiled as part of the host application. When the host application terminates, the runtime environment also terminates along with any tracking activity that it performs. If you restart the host application, a new runtime environment and its associated tracking activities reinitializes. The initialization functions reset the kernel's hardware state.

In some cases, unexpected termination of the host application causes the configuration of certain hardware (for example, PCIe hard IP) to be incomplete. To restore the configuration of these hardware, the host needs to reprogram the FPGA.

If you use a Custom Platform that implements customized hardware blocks, be aware that restarting the host application and resetting these blocks might have design implications:

- When the host application calls the `clGetPlatformIDs` function, all kernels and channels are reset for all available devices.
- When the host application calls the `clGetPlatformIDs` function, it resets FIFO buffers and channels as it resets the device.
- The host application initializes memory buffers via the `clCreateBuffer` and `clEnqueueWriteBuffer` function calls. You cannot access the contents of buffers from a previous host execution within a new host execution.

### 6.7. Allocating Shared Memory for OpenCL Kernels Targeting SoCs

Intel recommends that OpenCL kernels that run on Intel SoCs access shared memory instead of the FPGA DDR memory. FPGA DDR memory is accessible to kernels with very high bandwidths. However, read and write operations from the ARM CPU to FPGA DDR memory are very slow because they do not use direct memory access (DMA). Reserve FPGA DDR memory only for passing temporary data between kernels or within a single kernel for testing purposes.

*Note:*

- Mark the shared buffers between kernels as volatile to ensure that buffer modification by one kernel is visible to the other kernel.
- To access shared memory, you only need to modify the host code. Modifications to the kernel code are unnecessary.
- You cannot use the library function `malloc` or the operator `new` to allocate physically shared memory. Also, the `CL_MEM_USE_HOST_PTR` flag does not work with shared memory.

In DDR memory, shared memory must be physically contiguous. The FPGA cannot consume virtually contiguous memory without a scatter-gather direct memory access (SG-DMA) controller core. The `malloc` function and the `new` operator are for accessing memory that is virtually contiguous.

- CPU caching is disabled for the shared memory.
- When you use shared memory, one copy of the data is used for both the host and the kernel. When this memory is used, OpenCL memory calls are done as zero-copy transfers for buffer reads, buffer writers, maps, and unmaps.

The ARM CPU and the FPGA can access the shared memory simultaneously. You do not need to include the `clEnqueueReadBuffer` and `clEnqueueWriteBuffer` calls in your host code to make data visible to either the FPGA or the CPU.

- To allocate and access shared memory, structure your host code in a similar manner as the following example:

```
cl_mem src = clCreateBuffer(..., CL_MEM_ALLOC_HOST_PTR, size, ...);
int *src_ptr = (int*)clEnqueueMapBuffer (..., src, size, ...);
*src_ptr = input_value; //host writes to ptr directly
clSetKernelArg (..., src);
clEnqueueNDRangeKernel(...);
clFinish();
printf ("Result = %d\n", *dst_ptr); //result is available immediately
clEnqueueUnmapMemObject(..., src, src_ptr, ...);
clReleaseMemObject(src); // actually frees physical memory
```

You can include the `CONFIG_CMA_SIZE_MBYTES` kernel configuration option to control the maximum total amount of shared memory available for allocation. In practice, the total amount of allocated shared memory is smaller than the value of `CONFIG_CMA_SIZE_MBYTES`.

- Important:**
- If your target board has multiple DDR memory banks, the `clCreateBuffer(..., CL_MEM_READ_WRITE, ...)` function allocates memory to the nonshared DDR memory banks. However, if the FPGA has access to a single DDR bank that is shared memory, then `clCreateBuffer(..., CL_MEM_READ_WRITE, ...)` allocates to shared memory, similar to using the `CL_MEM_ALLOC_HOST_PTR` flag.
  - The shared memory that you request with the `clCreateBuffer(..., CL_MEM_ALLOC_HOST_PTR, size, ...)` function is allocated in the Linux OpenCL kernel driver, and it relies on the contiguous memory allocator (CMA) feature of the Linux kernel. For detailed information on enabling and configuring the CMA, refer to the *Recompiling the Linux Kernel for the Intel Arria 10 SoC Development Kit* and *Compiling and Installing the OpenCL Linux Kernel Driver* sections of the *Intel FPGA SDK for OpenCL Intel Arria 10 SoC Development Kit Reference Platform Porting Guide*.
- To transfer data from shared hard processor system (HPS) DDR to FPGA DDR efficiently, include a kernel that performs the `memcpy` function, as shown below.

```
__attribute__((num_simd_work_items(8)))
mem_stream(__global uint * src, __global uint * dst)
{
```

```
size_t gid = get_global_id(0);
dst[gid] = src[gid];
}
```

**Attention:** Allocate the `src` pointer in the HPS DDR as shared memory using the `CL_MEM_ALLOC_HOST_PTR` flag.

- If the host allocates constant memory to shared HPS DDR system and then modifies it after kernel execution, the modifications might not take effect. As a result, subsequent kernel executions might use outdated data. To prevent kernel execution from using outdated constant memory, perform one of the following tasks:
  - a. Do not modify constant memory after its initialization.
  - b. Create multiple constant memory buffers if you require multiple `__constant` data sets.
  - c. If available, allocate constant memory to the FPGA DDR on your accelerator board.

#### Related Information

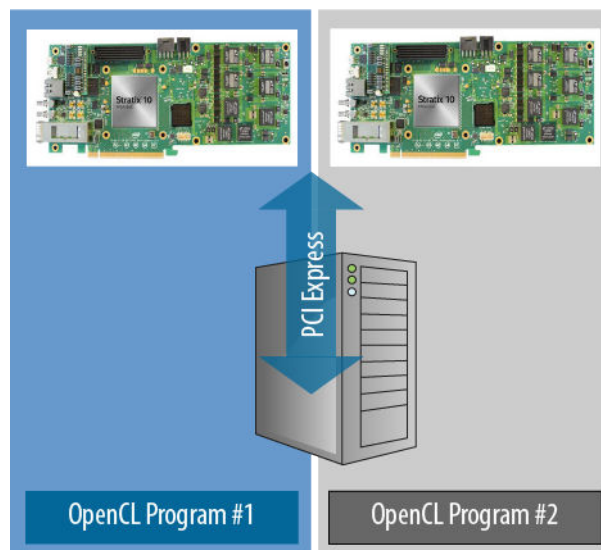
- [Recompiling the Linux Kernel for the Intel Arria 10 SoC Development Kit](#)
- [Compiling and Installing the OpenCL Linux Kernel Driver](#)

## 6.8. Sharing Multiple Devices Across Multiple Host Programs

In a system with multiple FPGA devices, each device appears as a separate `cl_device_id` object in the OpenCL host API. You can query various device properties using the `clGetDeviceInfo` function. Based on the properties, you can select devices you want to use in your program.

You can then pass these devices to `clCreateContext` function where the devices are locked by that OpenCL program until it either calls `clReleaseContext` function or terminates.

**Figure 14. Sharing Multiple Devices Across Multiple Host Programs**



Multiple processes or multiple trusted users can arbitrate between devices in a multi-device system using this locking mechanism. In case users have decided ahead of time which device (by name) each person uses, they can use `clGetDeviceInfo` function to select the `cl_device_id` with the name that a user is assigned to. To arbitrate more dynamically in case each of the  $N$  users want  $D_i$  devices, use the following scheme:

1. Each user queries the `clGetDeviceIDs` function to obtain a list of devices.
2. Each user chooses  $D_i$  devices (ideally randomly to minimize collisions) and passes those to the `clCreateContext` function.

It is possible that during step 2, another user may have already called `clCreateContext` function with that same device, in which case, the `clCreateContext` function call fails. The user should then repeat steps 2 (optionally changing the device selection) until it succeeds.

Consider the following example code snippet:

```
do {
    for(i = num_devices - 1; i >= 0; i--) {
        context = clCreateContext(0, 1, &(device_ids[i]), NULL, NULL, &status);

        if (status != CL_SUCCESS) {
            printf("Failed to get context with %d (error: %d), waiting\n", i, status);
            sleep(1);
        } else {
            device_id = device_ids[i];
            break;
        }
    }
} while (status != CL_SUCCESS);
// Exit this loop only when we've succeeded in creating a context
```

**Note:**

The only limitation is that only one process can own and operate on a device at a time. All other processes must wait for that process to release the context for obtaining access to the device.

## 7. Compiling Your OpenCL Kernel

---

The Intel FPGA SDK for OpenCL offers a list of compiler options that allows you to customize the kernel compilation process. An Intel FPGA SDK for OpenCL Offline Compiler command consists of the `aoc` command, compiler option(s) and settings, and kernel filenames. You can invoke an `aoc` command to direct the compiler to target a specific FPGA board, generate reports, or implement optimization techniques.

Before you compile an OpenCL kernel, verify that the `QUARTUS_ROOTDIR_OVERRIDE` environment variable points to the the Intel Quartus Prime Pro Edition software.

If these environment variables do not have the correct settings, follow the instructions in the *Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables* section of the *Intel FPGA SDK for OpenCL Pro Edition Getting Started Guide* to modify the settings.

### Compiling Encrypted Source:

When you compile an encrypted `.cl` file that is provided to you, you can compile only that file with the `aoc` command. You cannot compile multiple encrypted `.cl` files at the same time with the `aoc` command. You cannot compile kernel source files that you have encrypted yourself.

### Related Information

- [Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables \(Windows\)](#)
- [Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables \(Linux\)](#)

### 7.1. Compiling Your Kernel to Create Hardware Configuration File

You can compile an OpenCL kernel and create the hardware configuration file (that is, the `.aocx` file) in a single step.

Intel recommends that you use this one-step compilation strategy under the following circumstances:

- After you optimize your kernel via the Intel FPGA SDK for OpenCL design flow, and you are now ready to create the `.aocx` file for deployment onto the FPGA.
- You have one or more simple kernels that do not require any optimization.

To compile the kernel and generate the `.aocx` file in one step, invoke the `aoc` `<your_kernel_filename1>.cl [<your_kernel_filename2>.cl ...]` command.

Where `[<your_kernel_filename2>.cl ...]` are the optional space-delimited file names of kernels that you can compile in addition to `<your_kernel_filename1>.cl`.

The Intel FPGA SDK for OpenCL Offline Compiler groups the `.cl` files into a temporary file. It then compiles this file to generate the `.aocx` file.

**Note:** If you run the `aoc` command to compile a `.aoco` file (that is, `aoc <your_kernel_filename>.aoco`), the offline compiler generates both a `.aocr` file and a `.aocx` file.

## 7.2. Compiling Your Kernel without Building Hardware (-c)

To direct the Intel FPGA SDK for OpenCL Offline Compiler to compile your OpenCL kernel and generate the output of the OpenCL parser without creating a hardware configuration file, include the `-c` option in your `aoc` command.

**Note:** The `-c` flag is not supported in the incremental compilation flow.

- At a command prompt, invoke the `aoc -c <your_kernel_filename1>.cl [<your_kernel_filename2>.cl ...]` command.

Where [`<your_kernel_filename2>.cl ...`] are the optional space-delimited file names of kernels that you can compile in addition to `<your_kernel_filename1>.cl`.

When you invoke the `aoc` command with the `-c` flag, the offline compiler compiles the kernel(s) and creates the following files and directories:

- A `.aoco` file for each `.cl` kernel source file. The offline compiler creates the `.aoco` file(s) in a matter of seconds to minutes.

## 7.3. Compiling and Linking Your Kernels or Object Files without Building Hardware (-rtl)

To direct the Intel FPGA SDK for OpenCL Offline Compiler to compile your OpenCL kernels (`.cl`), generate an intermediate `.aoco` object file for each kernel, and then link them together to create a `.aocr` file without creating a hardware configuration file, include the `-rtl` option in your `aoc` command.



- To compile one or more kernel source files, at a command prompt, invoke the `aoc -rtl <your_kernel_filename1>.cl [<your_kernel_filename2>.cl ...]` command.

Where [`<your_kernel_filename2>.cl ...`] are the optional space-delimited file names of kernels that you can compile in addition to `<your_kernel_filename1>.cl`.

When you invoke the `aoc` command with the `-rtl` flag, the offline compiler compiles the kernels and creates the following files and directories:

- An intermediate `.aoco` file for each `.cl` kernel source file. It is not presented unless you specify the `-save-temps` `aoc` command option. The offline compiler then links them and generates a `.aocr` file. It takes the offline compiler a matter of seconds to minutes to create a `.aoco` file or the `.aocr` file.
  - A `<your_kernel_filename>` folder or subdirectory. It contains intermediate files that the SDK uses to build the hardware configuration file necessary for FPGA programming.
- To compile one or more `.aoco` object files, at a command prompt, invoke the `aoc -rtl <your_kernel_filename>.aoco [<your_kernel_filename2>.aoco ...]` command.

Where [`<your_kernel_filename2>.aoco ...`] are the optional space-delimited file names of object files that you can compile in addition to `<your_kernel_filename1>.aoco`.

When you invoke the `aoc` command with the `-rtl` flag, the offline compiler creates the following files and directories:

- The offline compiler links all the `.aoco` files and generates a `.aocr` file.
- A `<your_kernel_filename>` folder or subdirectory. It contains intermediate files that the SDK uses to build the hardware configuration file necessary for FPGA programming.

## 7.4. Specifying the Location of Header Files (-I=<directory>)

To add a directory to the list of directories that the Intel FPGA SDK for OpenCL Offline Compiler searches for header files during kernel compilation, include the `-I=<directory>` option in your `aoc` command.

If the header files are in the same directory as your kernel, you do not need to include the `-I=<directory>` option in your `aoc` command. The offline compiler automatically searches the current folder or directory for header files.

- At a command prompt, invoke the `aoc -I=<directory> <your_kernel_filename>.cl` command.

**Caution:** For Windows systems, ensure that your include path does not contain any trailing slashes. The offline compiler considers a trailing forward slash (/) or backward slash (\) as illegal.

The offline compiler generates an error message if you invoke the `aoc` command in the following manner:

```
aoc -I=<drive>\<folder>\<subfolder>\ <your_kernel_filename>.cl  
or
```

```
aoc -I=<drive>/<folder>/<subfolder>/ <your_kernel_filename>.cl
```

The correct way to specify the include path is as follows:

```
aoc -I=<drive>\<folder>\<subfolder> <your_kernel_filename>.cl  
or
```

```
aoc -I=<drive>/<folder>/<subfolder> <your_kernel_filename>.cl
```

## 7.5. Specifying the Name of an Intel FPGA SDK for OpenCL Offline Compiler Output File (-o <filename>)

To specify the name of a `.aocr` file or a `.aocx` file, include the `-o <filename>` option in your `aoc` command.

- If you implement the multistep compilation flow, specify the names of the output files in the following manner:
  - a. To specify the name of the `.aoco` file that the offline compiler creates during an intermediate compilation step, invoke the `aoc -rtl -o <your_object_filename>.aocr <your_kernel_filename>.cl -save-temps` command.
  - b. To specify the name of the `.aocx` file that the offline compiler creates during the final compilation step, invoke the `aoc -o <your_executable_filename>.aocx <your_object_filename>.aocr` command.
- If you implement the one-step compilation flow, specify the name of the `.aocx` file by invoking the `aoc -o <your_executable_filename>.aocx <your_kernel_filename>.cl` command.

## 7.6. Compiling a Kernel for a Specific FPGA Board and Custom Platform (-board=<board\_name>) and (-board-package=<board\_package\_path>)

To compile your OpenCL kernel for a specific FPGA board, include the `-board=<board_name>` option in the `aoc` command.

**Warning:** During intermediate compilation with the `-rtl` flag, if you compile a `kernel1.cl` file for board X (that is, `-board=X`) and a `kernel2.cl` file for board Y (that is, `-board=Y`), the Intel FPGA SDK for OpenCL Offline Compiler issues an error and exits when you run the `aoc -rtl kernel1.aoco kernel2.aoco` command.

When you compile your kernel by including the `-board=<board_name>` option in the `aoc` command, the Intel FPGA SDK for OpenCL Offline Compiler defines the preprocessor macro `AOCL_BOARD_<board_name>` to be 1, which allows you to compile device-optimized code in your kernel.

- To obtain the names of the available FPGA boards in your Custom Platform, invoke the `aoc -list-boards` command.

For example, the offline compiler generates the following output:

```
Board List:  
FPGA_board_1
```

where `FPGA_board_1` is the `<board_name>`.

You can also list out all the available FPGA boards from a specific Custom Platform. Include the `-board-package=<custom_platform_path>` option in the `aoc` command. At the command prompt, invoke the following command:

```
aoc -board-package=<custom_platform_path> -list-boards=<board_name>
```

The Intel FPGA SDK for OpenCL Offline Compiler lists the available boards within the specific Custom Platform.

- To compile your OpenCL kernel for `FPGA_board_1`, invoke the `aoc -board=FPGA_board_1 <your_kernel_filename>.cl` command. The offline compiler defines the preprocessor macro `AOCL_BOARD_FPGA_board_1` to be 1 and compiles kernel code that targets `FPGA_board_1`.
- If there are multiple Custom Platforms (board packages) installed, you can compile your kernel with the board variant from a specific Custom Platform by including `-board-package=<custom_platform_path>` option with `-board=<board_name>`. At the command prompt, invoke the following command:

```
aoc -board-package=<custom_platform_path> -board=<board_name>
```

The Intel FPGA SDK for OpenCL Offline Compiler compiles the kernel with the board specified in the `<custom_platform_path>`.

- To list Custom Platforms available in the system, include the `-list-board-packages` option in the `aoc` command. At a command prompt, invoke the `aoc -list-board-packages` command. The Intel FPGA SDK for OpenCL Offline Compiler generates an output that resembles the following:

```
Installed board packages:  
<board_package_1>  
...
```

Where `<board_package_N>` is the board package of the Custom Platform installed in your system or shipped within the Intel FPGA SDK for OpenCL.

**Tip:** To readily identify compiled kernel files that target a specific FPGA board, Intel recommends that you rename the kernel binaries by including the `-o` option in the `aoc` command.

- To target your kernel to `FPGA_board_1` in the one-step compilation flow, invoke the following command:

```
aoc -board=FPGA_board_1 <your_kernel_filename>.cl -o  
<your_executable_filename>_FPGA_board_1.aocx
```

- To target your kernel to `FPGA_board_1` in the multistep compilation flow, perform the following tasks:

1. Invoke the following command to generate the `.aoco` file:

```
aoc -rtl -board=FPGA_board_1 <your_kernel_filename>.cl -o  
<my_object_filename>_FPGA_board_1.aocr -save-temps
```

2. Invoke the following command to generate the `.aocx` file:

```
aoc -board=FPGA_board_1 <your_object_filename>_FPGA_board_1.aocr -o  
<your_executable_filename>_FPGA_board_1.aocx
```

- If you have an accelerator board consisting of two FPGAs, each FPGA device has an equivalent "board" name (for example, `board_fpga_1` and `board_fpga_2`). To target a `kernel_1.cl` to `board_fpga_1` and a `kernel_2.cl` to `board_fpga_2`, invoke the following commands:

```
aoc -board=board_fpga1 kernel_1.cl  
aoc -board=board_fpga2 kernel_2.cl
```

### Related Information

[Specifying the Name of an Intel FPGA SDK for OpenCL Offline Compiler Output File \(`-o <filename>`\)](#) on page 122

## 7.7. Resolving Hardware Generation Fitting Errors during Kernel Compilation (`-high-effort`)

Sometimes, OpenCL kernel compilation fails during the hardware generation stage because the design fails to meet fitting constraints. In this case, recompile the kernel using the `-high-effort` option of the `aoc` command.

When kernel compilation fails because of a fitting constraint problem, the Intel FPGA SDK for OpenCL Offline Compiler displays the following error message:

```
Error: Kernel fit error, recommend using -high-effort.  
Error: Cannot fit kernel(s) on device
```

- To overcome this problem, recompile your kernel by invoking the following command:

```
aoc -high-effort <your_kernel_filename>.cl
```

After you invoke the command, the offline compiler displays the following message:

```
High-effort hardware generation selected, compile time may increase  
significantly.
```

The offline compiler makes three attempts to recompile your kernel and generate hardware. Modify your kernel if compilation still fails after the `-high-effort` attempt.

## 7.8. Specifying Schedule $f_{\max}$ Target for Kernels (`-clock=<clock_target>`)

The schedule  $f_{\max}$  target determines the pipelining effort the scheduler attempts during the scheduling process.

You can use one or both of the following options to specify the kernel specific  $f_{\max}$  target:

- By using the `__attribute__((scheduler_target_fmax_mhz(__x)))` source-level attribute.
- By directing the Intel FPGA SDK for OpenCL Offline Compiler to globally compile all kernels with `-clock=<clock target in Hz/KHz/MHz/GHz or s/ms/us/ns/ps>` option in the `aoc` command.

If you use both the command line option and source level attribute, the kernel attribute takes the priority. Consider the following example:

```
kernel void k1(){  
    ...  
}  
__attribute__((scheduler_target_fmax_mhz(200)))  
kernel void k2(){  
    ...  
}
```

In you direct the offline compiler to compile the above code with `-clock=300MHz` in the `aoc` command, the compiler schedules kernel `k1` at 300 MHz and kernel `k2` at 200 MHz.

**Important:** The schedule target  $f_{\max}$  determines the pipelining effort during compilation. Run Intel Quartus Prime Pro Edition compiler to get the actual  $f_{\max}$  value.

## 7.9. Defining Preprocessor Macros to Specify Kernel Parameters (`-D<macro_name>`)

The Intel FPGA SDK for OpenCL Offline Compiler supports preprocessor macros that allow you to pass macro definitions and compile code on a conditional basis.

- To pass a preprocessor macro definition to the offline compiler, invoke the `aoc -D<macro_name> <kernel_filename>.cl` command.
- To override the existing value of a defined preprocessor macro, invoke the `aoc -D<macro_name>=<value> <kernel_filename>.cl` command.

Consider the following code snippet for the kernel `sum`:

```
#ifndef UNROLL_FACTOR  
    #define UNROLL_FACTOR 1  
#endif  
  
__kernel void sum (__global const int * restrict x,  
                  __global int * restrict sum)
```

```

{
  int accum = 0;

  #pragma unroll UNROLL_FACTOR
  for(size_t i = 0; i < 4; i++)
  {
    accum += x[i + get_global_id(0) * 4];
  }
  sum[get_global_id(0)] = accum;
}

```

To override the UNROLL\_FACTOR of 1 and set it to 4, invoke the `aoc -DUNROLL_FACTOR=4 sum.cl` command. Invoking this command is equivalent to replacing the line `#define UNROLL_FACTOR 1` with `#define UNROLL_FACTOR 4` in the `sum` kernel source code.

- To use preprocessor macros to control how the offline compiler optimizes your kernel without modifying your kernel source code, invoke the `aoc -o <hardware_filename>.aocx -D<macro_name>=<value> <kernel_filename>.cl`

Where:

`-o` is the offline compiler option you use to specify the name of the `.aocx` file that the offline compiler generates.

`<hardware_filename>` is the name of the `.aocx` file that the offline compiler generates using the preprocessor macro value you specify.

*Tip:* To preserve the results from both compilations on your file system, compile your kernels as separate binaries by using the `-o` flag of the `aoc` command.

For example, if you want to compile the same kernel multiple times with required work-group sizes of 64 and 128, you can define a `WORK_GROUP_SIZE` preprocessor macro for the kernel attribute `reqd_work_group_size`, as shown below:

```

__attribute__((reqd_work_group_size(WORK_GROUP_SIZE,1,1)))
__kernel void myKernel(...)
for (size_t i = 0; i < 1024; i++)
{
    // statements
}

```

Compile the kernel multiple times by typing the following commands:

```
aoc -o myKernel_64.aocx -DWORK_GROUP_SIZE=64 myKernel.cl
```

```
aoc -o myKernel_128.aocx -DWORK_GROUP_SIZE=128 myKernel.cl
```

## 7.10. Generating Compilation Progress Report (-v)

To direct the Intel FPGA SDK for OpenCL Offline Compiler to report on the progress of a compilation, include the `-v` option in your `aoc` command.

- To direct the offline compiler to report on the progress of a full compilation, invoke the `aoc -v <your_kernel_filename>.cl` command.

The offline compiler generates a compilation progress report similar to the following example:

```
aoc: Environment checks are completed successfully.
You are now compiling the full flow!!
aoc: Selected target board a10gx
aoc: Running OpenCL parser...
aoc: OpenCL parser completed successfully.
aoc: Compiling...
aoc: Linking with IP library ...
aoc: First stage compilation completed successfully.
aoc: Setting up project for CvP revision flow...
aoc: Hardware generation completed successfully.
```

- To direct the offline compiler to report on the progress of an intermediate compilation step that does not build hardware, invoke the `aoc -rtl -v <your_kernel_filename>.cl` command.

The offline compiler generates a compilation progress report similar to the following example:

```
aoc: Environment checks are completed successfully.
aoc: Selected target board a10gx
aoc: Running OpenCL parser...
aoc: OpenCL parser completed successfully.
aoc: Compiling...
aoc: Linking with IP library ...
aoc: First stage compilation completed successfully.
aoc: To compile this project, run "aoc <your_kernel_filename>.aoco"
```

- To direct the offline compiler to report on the progress of a compilation for emulation, invoke the `aoc -march=emulator -v <your_kernel_filename>.cl` command.

The offline compiler generates a compilation progress report similar to the following example:

```
aoc: Environment checks are completed successfully.
You are now compiling the full flow!!
aoc: Selected target board a10gx
aoc: Running OpenCL parser...ex
aoc: OpenCL parser completed successfully.
aoc: Compiling for Emulation ...
aoc: Emulator Compilation completed successfully.
Emulator flow is successful.
```

### Related Information

- [Compiling and Linking Your Kernels or Object Files without Building Hardware \(-rtl\)](#) on page 120
- [Emulating and Debugging Your OpenCL Kernel](#) on page 140

## 7.11. Displaying the Estimated Resource Usage Summary On-Screen (-report)

By default, the Intel FPGA SDK for OpenCL Offline Compiler estimates hardware resource usage during compilation. The offline compiler factors in the usage of external interfaces such as PCIe, memory controller, and DMA engine in its calculations. During kernel compilation, the offline compiler generates an estimated resource usage summary in the `<your_kernel_filename>.log` file within the `<your_kernel_filename>` directory. To review the estimated resource usage summary on-screen, include the `-report` option in the `aoc` command.

You can review the estimated resource usage summary without performing a full compilation. To review the summary on-screen prior to generating the hardware configuration file, include the `-rtl` option in your `aoc` command.

- At a command prompt, invoke the `aoc -rtl <your_kernel_filename>.cl -report` command. The offline compiler generates an output similar to the following example:

```
+-----+
; Estimated Resource Usage Summary ;
+-----+-----+-----+
; Resource + Usage ;
+-----+-----+-----+
; Logic utilization ; 35% ;
; ALUTs ; 22% ;
; Dedicated logic registers ; 15% ;
; Memory blocks ; 29% ;
; DSP blocks ; 0% ;
+-----+-----+-----+
```

### Related Information

[Compiling and Linking Your Kernels or Object Files without Building Hardware \(-rtl\)](#) on page 120

## 7.12. Suppressing Warning Messages from the Intel FPGA SDK for OpenCL Offline Compiler (-W)

To suppress all warning messages, include the `-w` option in your `aoc` command.

- At a command prompt, invoke the `aoc -W <your_kernel_filename>.cl` command.

## 7.13. Converting Warning Messages from the Intel FPGA SDK for OpenCL Offline Compiler into Error Messages (-Werror)

To convert all warning messages into error messages, include the `-Werror` option in your `aoc` command.

- At a command prompt, invoke the `aoc -Werror <your_kernel_filename>.cl` command.



## 7.14. Removing Debug Data from Compiler Reports and Source Code from the .aocx File (-g0)

By default, the Intel FPGA SDK for OpenCL Offline Compiler includes source information in compiler reports and embeds the source code into the .aocx binary when it compiles the .cl or .aoco file. Include the -g0 option in the aoc command to remove source information from the compiler reports and to remove source code and customer IP information from the .aocx file.

- To remove source information in reports and remove source code and customer IP information from the .aocx file, invoke the `aoc -g0 <your_kernel_filename>.cl` command.

## 7.15. Disabling Burst-Interleaving of Global Memory (-no-interleaving=<global\_memory\_type>)

The Intel FPGA SDK for OpenCL Offline Compiler cannot burst-interleave global memory across different memory types. You can disable burst-interleaving for all global memory banks of the same type and manage them manually by including the -no-interleaving=<global\_memory\_type> option in your aoc command. Manual partitioning of memory buffers overrides the default burst-interleaved configuration of global memory.

**Caution:** The -no-interleaving option requires a global memory type parameter. If you do not specify a memory type, the offline compiler issues an error message.

- To direct the offline compiler to disable burst-interleaving for the default global memory, invoke the `aoc <your_kernel_filename>.cl -no-interleaving=default` command.  
Your accelerator board might include multiple global memory types. To identify the default global memory type, refer to board vendor's documentation for your Custom Platform.
- For a heterogeneous memory system, to direct the offline compiler to disable burst-interleaving of a specific global memory type, perform the following tasks:
  - a. Consult the `board_spec.xml` file of your Custom Platform for the names of the available global memory types (for example, DDR and quad data rate (QDR)).
  - b. To disable burst-interleaving for one of the memory types (for example, DDR), invoke the `aoc <your_kernel_filename>.cl -no-interleaving=DDR` command.  
The offline compiler enables manual partitioning for the DDR memory bank, and configures the other memory bank in a burst-interleaved fashion.
  - c. To disable burst-interleaving for more than one type of global memory buffers, include a -no-interleaving=<global\_memory\_type> option for each global memory type.  
For example, to disable burst-interleaving for both DDR and QDR, invoke the `aoc <your_kernel_filename>.cl -no-interleaving=DDR -no-interleaving=QDR` command.

**Caution:** Do not pass a buffer as kernel arguments that associate it with multiple memory technologies.

## 7.16. Forcing Ring Interconnect for Global Memory (-global-ring)

The Intel FPGA SDK for OpenCL Offline Compiler attempts to choose an optimal global memory interconnect topology based on various characteristics of the design.

To override the compiler's choice and force a ring topology, use the `-global-ring` option in your `aoc` command. This can improve your kernel  $f_{max}$ .

Example: `aoc -global-ring <your_kernel_filename>.cl`

## 7.17. Forcing a Single Store Ring to Reduce Area at the Expense of Write Throughput to Global Memory (-force-single-store-ring)

To narrow the global memory interconnect in order to save area while limiting write-only throughput to one bank's worth, use the `-force-single-store-ring` option in your `aoc` command.

When the Intel FPGA SDK for OpenCL Offline Compiler implements a ring topology for the global memory interconnect (either by automatic choice or by forcing the ring through `-global-ring`), it widens the interconnect by default to allow more writes to occur in parallel. This allows for the saturation of global memory throughput using write-only traffic. The `-force-single-store-ring` option allows you to save area if you do not require that much write bandwidth.

Example: `aoc -force-single-store-ring <your_kernel_filename>.cl`

## 7.18. Forcing Fewer Read Data Reorder Units to Reduce Area at the Expense of Read Throughput to Global Memory (-num-reorder)

To narrow the global memory interconnect in order to save area while reducing read-only throughput, use the `-num-reorder=N` option in your `aoc` command, where  $N$  is the number of bank's worth of read bandwidth you desire.

When the Intel FPGA SDK for OpenCL Offline Compiler implements a ring topology for the global memory interconnect (either by automatic choice or by forcing the ring through `-global-ring`), it widens the interconnect by default to allow more reads to occur in parallel. This allows for the saturation of global memory throughput using read-only traffic. For example, if on a two-bank BSP you require only one bank's worth of read bandwidth, set `-num-reorder=1`.

Example: `aoc -num-reorder=1 <your_kernel_filename>.cl`

## 7.19. Configuring Constant Memory Cache Size (-const-cache-bytes=<N>)

Include the `-const-cache-bytes=<N>` flag in your `aoc` command to direct the Intel FPGA SDK for OpenCL Offline Compiler to configure the constant memory cache size (rounded up to the closest power of 2).

The default constant cache size is 16 kB.

- To configure the constant memory cache size, invoke the `aoc -const-cache-bytes=<N> <your_kernel_filename>.cl` command, where <N> is the cache size in bytes.

For example, to configure a 32 kB cache during compilation of the OpenCL kernel `myKernel.cl`, invoke the `aoc -const-cache-bytes=32768 myKernel.cl` command.

*Note:* This argument has no effect if none of the kernels uses the `__constant` address space.

## 7.20. Relaxing the Order of Floating-Point Operations (`-ffp-reassociate`)

Include the `-ffp-reassociate` option in your `aoc` command to direct the Intel FPGA SDK for OpenCL Offline Compiler to relax the order of arithmetic floating-point operations using a balanced tree hardware implementation.

This flag turns on the same optimizations as the deprecated `-fp-relaxed` flag for all instructions, unless denoted otherwise by the `fp reassociate` pragma. For information about the `fp reassociate` pragma, refer to [Floating Point Optimizations \(fp contract and fp reassociate Pragma\)](#) on page 38.

*Note:* Implementing a balanced tree structure leads to more efficient hardware at the expense of numerical variation in results. To implement this optimization control, your program must be able to tolerate small variations in the floating-point results.

To direct the offline compiler to execute a balanced tree hardware implementation, invoke the `aoc -ffp-reassociate <your_kernel_filename>.cl` command.

## 7.21. Reducing Floating-Point Rounding Operations (`-ffp-contract=fast`)

Include the `-ffp-contract=fast` option in your `aoc` command to direct the Intel FPGA SDK for OpenCL Offline Compiler to remove intermediary floating-point rounding operations and conversions whenever possible, and to carry additional bits to maintain precision.

This flag turns on the same optimizations as the deprecated `-fpc` flag for all instructions, unless denoted otherwise by the `fp contract` pragma. For information about the `fp contract` pragma, refer to [Floating Point Optimizations \(fp contract and fp reassociate Pragma\)](#) on page 38.

*Note:* Implementing this optimization control also changes the rounding mode. It rounds towards zero only at the end of a chain of floating-point arithmetic operations (that is, multiplications, additions, and subtractions).

- To direct the offline compiler to reduce the number of rounding operations, invoke the `aoc -ffp-contract=fast <your_kernel_filename>.cl` command.

## 7.22. Speeding Up Your OpenCL Compilation (-fast-compile)

To save 40-90% of compilation time and quickly create the `.aocx` file of your kernel, include the `-fast-compile` Intel FPGA SDK for OpenCL Offline Compiler command option in your `aoc` command.

Intel recommends that you use the `-fast-compile` option for internal development only.

The `-fast-compile` feature achieves significant savings in compilation time by lowering optimization efforts.

At the command prompt, invoke the `aoc -rtl <your_kernel_filename>.cl -fast-compile` command.

- Warning:**
- Enabling the `-fast-compile` feature might cause some performance issues such as:
    - Higher resource use
    - Lower  $f_{max}$  and as a result lower application performance
    - Lower power efficiency
  - When compiling your kernel using the `-fast-compile` flag, you might see functional failures due to timing violations in your design. In such cases, either avoid using the `-fast-compile` flag or try compiling your kernel with different seeds.

- Attention:**
- You can only use the `-fast-compile` compiler option to compile OpenCL designs targeting Intel Arria 10 and newer devices.
  - After you finalize a design, compile your OpenCL kernel without the `-fast-compile` option over multiple seeds to obtain the best performance.
  - Regardless of whether the `-fast-compile` feature is enabled, the initial compilation of any OpenCL system on a new board and with a new version of Intel FPGA SDK for OpenCL Pro Edition takes an additional 45 to 60 minutes to complete. The additional time is used to cache some parts of the compilation for future compilations (this behavior does *not* affect kernel performance). To create this cache, define the environment variable `$AOCL_TMP_DIR` to a writable directory that you can share. By default, this cache is stored in `/var/tmp/aocl/$USER` on Linux and `%USERPROFILE%\AppData\Local\aocl` on Windows. You can share this writable directory by setting it to a shared network location.  
After you create the cache, you do not need to create it again for the current version of the Intel FPGA SDK for OpenCL and the current targeted board.

## 7.23. Compiling Your Kernel Incrementally (-incremental)

To compile the modifications you make to your OpenCL design incrementally, include the `-incremental` Intel FPGA SDK for OpenCL Offline Compiler command option in your `aoc` command.

If you have a large, multi-kernel system and you only want to modify a single kernel, the Intel FPGA SDK for OpenCL Offline Compiler can reuse the results from a previous compilation, and only synthesize, place, and route the kernel(s) that you have modified. Leveraging this incremental compilation feature allows you to dramatically reduce compilation time.

**Attention:** Enable the incremental compilation feature for internal development of your OpenCL design only. For best circuit performance on your final design, run a full compilation.

Example incremental compilation flow:

```
aoc -incremental <your_kernel_filename>.cl  
/*****Update kernels in your OpenCL design*****/  
aoc -incremental -fast-compile <your_kernel_filename>.cl
```

1. Perform an initial setup compilation in a clean directory, with the incremental mode enabled, by invoking the `aoc -incremental <your_kernel_filename>.cl` command.

*Note:* You must enable the `-incremental` flag when performing the setup compilation.

This setup compilation does not reuse any result from a previous compilation. When performing a setup compilation, do not include the `-fast-compile` offline compiler command option in the `aoc` command because it increases the probability of encountering errors in future incremental compilations.

*Tip:* Intel recommends that you perform a fresh setup compilation whenever compilation time is not a concern because it reduces the probability of compilation failures in future incremental compilations. Performing many consecutive incremental compilations increases the probability of compilation failures. It also decreases the hardware performance and efficiency of the generated `.aocx` file.

2. Modify the kernels in your OpenCL design.  
Your design may contain multiple `.cl` files.
3. Perform an incremental compilation on your design. For optimal compilation speed, also include the `-fast-compile` flag in your `aoc` command:  

```
aoc -incremental -fast-compile <your_kernel_filename>.cl
```
4. Review the **Incremental compile** section of the `report.html` file to verify the changes that the offline compiler has detected.

The `report.html` file is in the `<your_kernel_filename>/reports` directory.

### 7.23.1. The Incremental Compile Report

Compiling your OpenCL design with the `-incremental` Intel FPGA SDK for OpenCL Offline Compiler command option instructs the offline compiler to include an **Incremental compile** section in the `report.html` file in the `<your_kernel_filename>` project directory.

The Incremental compile report provides the following metrics on your OpenCL design:

- The  $\langle\% \rangle$  of design not preserved metric at the bottom of the report provides a quick summary of the overall changes to your design. It is the best predictor of compilation time.

*Note:* The FPGA resources listed in the Incremental compile report are calculated based on the estimated area models that the Intel FPGA SDK for OpenCL Offline Compiler produces. The area numbers represent an estimate of the area usage in a standard (that is, non-incremental) compilation. You can use these numbers to gauge the area your design consumes in a standard compilation.

The FPGA resource information might not fully match the final area in the Intel Quartus Prime Pro Edition software compilation reports.

Figure 15. Incremental compile Report for a Setup Compilation

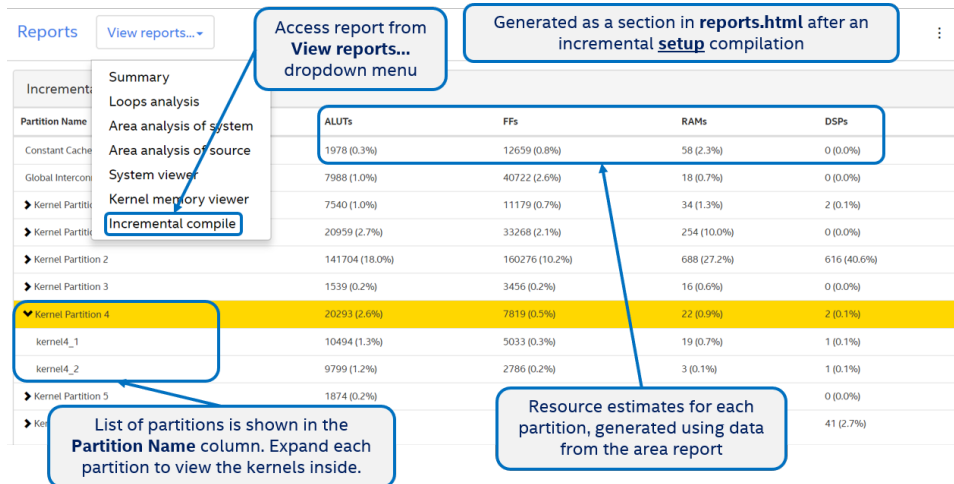
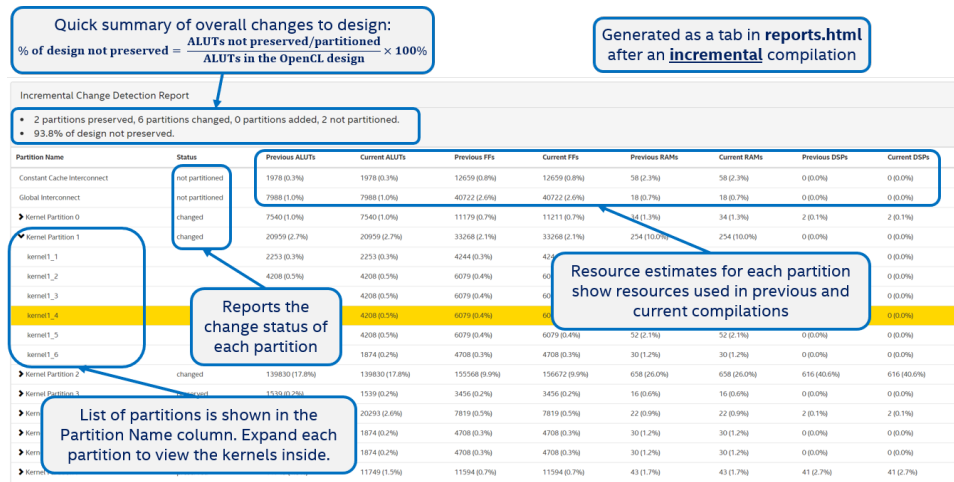


Figure 16. Incremental compile Report for an Incremental Compilation



## 7.23.2. Additional Command Options for Incremental Compilation

The Intel FPGA SDK for OpenCL incremental compilation feature includes optional functions that you can enable to customize the compilation of your OpenCL design.

### Grouping Multiple Kernels into Partitions (-incremental-grouping=<filename>)

By default, the Intel FPGA SDK for OpenCL Offline Compiler places each kernel in your design into a separate partition during incremental compilation. You have the option to group multiple kernels into a single partition by including the `-incremental-grouping=<partition_filename>` command option in your `aoc` command. In general, compilation speed is faster if your design contains fewer partitions.

Example: `aoc -incremental-grouping=<partition_filename> <your_kernel_filename>.cl`

**Note:** The offline compiler recompiles all kernels in a group even if you modify only one of the kernels. Intel recommends that you group kernels that you typically want to modify at the same time.

If your grouped kernels perform many load and store operations, you can speed up compilation further by also including the `-incremental=aggressive` option in your `aoc` command.

The partition file that you pass to the `-incremental-grouping` option is a plain text file. Each line in the file specifies a new partition containing a semi-colon (;)-delimited list of kernel names. For example, the following lines in a partition file specify three partitions, each containing four kernels:

```
reader0;reader1;reader2;reader3
accum0;accum1;accum2;accum3
writer0;writer1;writer2;writer3
```

### Compiling a Design in Aggressive Mode (-incremental=aggressive)

To increase the speed of an incremental compilation at the expense of area usage and throughput, include the `-incremental=aggressive` command option in your `aoc` command.

This feature is especially effective when the kernels in your design perform load and store operations to many buffers, or when you have grouped multiple kernels together using the `-incremental-grouping` command option.

Example: `aoc -incremental=aggressive -incremental-grouping=<partition_filename> <your_kernel_filename>.cl`

**Caution:**

- Enabling the aggressive mode might result in throughput degradations that are larger than what the Fmax degradation indicates.
- For each OpenCL design, avoid changing the compilation mode between incremental compilations. If you compile your design in aggressive mode, enable aggressive mode for all subsequent incremental compilations that you perform on this design. Each time you switch the incremental compilation mode, compilation takes longer to complete.

### Specifying a Custom Input Directory (-incremental-input-dir=<path\_to\_directory>)

During incremental compilation, the offline compiler creates a default `<your_kernel_filename>` project directory in the current working directory to store intermediate compilation files. To base your incremental compilation on a nondefault project directory, specify the directory by including the `-incremental-input-dir=<path_to_directory>` command option in your `aoc` command.

You must include the `-incremental-input-dir` option if you compile your design in one or both of the following scenarios:

- Run the `aoc` command from a different working directory than the previous compilation.
- Included the `-o <filename>` command option in the previous compilation.

Consider the following example where there is a `mykernel.cl` file in the initial working directory and another revision of the same `mykernel.cl` file in the `new_rev` subdirectory:

```
aoc -incremental mykernel.cl
cd new_rev
aoc -incremental -fast-compile mykernel.cl -incremental-input-dir=./mykernel
```

In this scenario, the offline compiler reuses the files in `mykernel` project directory from the first compilation as the basis for the second compilation. The offline compiler creates a `new_rev/mykernel` project directory for the second compilation without modifying any file in the original `mykernel` directory.

The `-incremental-input-dir` command option is useful if multiple developers share the same incremental setup compilation. Each developer can run subsequent incremental compilations in their own workspace without overwriting other developers' compilation results.

### Disabling Automatic Retry (-incremental-flow=no-retry)

By default, the offline compiler automatically retries a failed incremental compilation by performing a second compilation without preserving any partitions. This second compilation takes longer to complete because it recompiles the entire design.

To disable the offline compiler's automatic retry mechanism, include the `-incremental-flow=no-retry` command option in your `aoc` command. If you enable this feature, the offline compiler does not perform another incremental compilation after the first attempt fails. In addition, the offline compiler does not generate a `.aocx` file.

Enabling this feature allows you to implement your own failure mitigation strategies such as:

- Compiling multiple seeds in parallel to increase the probability of at least one compilation succeeding without retrying.
- Executing a **non-incremental** fast compilation instead of an incremental fast compilation (that is, `aoc -fast-compile <your_kernel_filename>.cl`).



### 7.23.3. Limitations of the Incremental Compilation Feature

The Intel FPGA SDK for OpenCL incremental compilation feature is only available to OpenCL designs targeting the Intel Arria 10 FPGAs.

In addition to device support, the incremental compilation has the following limitations:

- You experience area,  $F_{\max}$ , and power degradations when you enable the incremental compilation feature (`-incremental`) or the fast compilation feature (`-fast-compile`), or both.
- In congested designs, incremental compilations can experience severe (that is, 25% or more)  $F_{\max}$  reductions compared to the initial setup compilation. If the  $F_{\max}$  reduction is unacceptable, perform a non-incremental fast compilation instead to reduce the amount of  $F_{\max}$  degradation while preserving some of the savings in compilation time.
- The offline compiler does not detect changes in RTL libraries that you have included by invoking the `-l <library_name>.aoclib` offline compiler command option. After you modify an RTL library, you must perform a setup compilation again.

The offline compiler prints a warning message as a reminder to rerun the setup compilation.

- When compiling an OpenCL kernel containing calls to HLS tasks, incremental compile may trigger recompilation for unaffected kernels. However, this is not a functional bug. It may result in a more conservative incremental compile.

### 7.24. Compiling Your Kernel with Memory Error Correction Coding (-ecc)

**Attention:** Error correction coding (ECC) is an early Intel FPGA SDK for OpenCL feature that is at the preview stage. Full use of this feature, including the reporting of corrected errors and detected but uncorrected errors, requires an ECC-ready Custom Platform from your board vendor.

Include the `-ecc` option in your `aoc` command to direct the Intel FPGA SDK for OpenCL Offline Compiler to enable error correction coding on the kernel memories (that is, M20ks and MLABs).

The ECC implementation has single error correction and double error detection capabilities for each 32-bit word.

**Caution:** Enabling the ECC feature costs an area overhead both in the number of RAMs and ALMs, as well as causes degradation in the  $F_{\max}$  of the system.

- To direct the offline compiler to enable error correction coding hardware implementation, invoke the `aoc -ecc <your_kernel_filename>.cl` command.

## 7.25. Disabling Hardware Kernel Invocation Queue (-no-hardware-kernel-invocation-queue)

To direct the Intel FPGA SDK for OpenCL Offline Compiler to reduce kernel area use by removing kernel invocation queue in OpenCL kernel, include the `-no-hardware-kernel-invocation-queue` option in your `aoc` command.

Example: `aoc -no-hardware-kernel-invocation-queue <your_kernel_filename>.cl`

Using this option may result in longer kernel execution time as the kernel invocation queue allows OpenCL runtime environment to queue kernel launches in accelerator so that the accelerator can start execution on the next invocation as soon as previous invocation of the same kernel is complete.

**Attention:** Use the `-no-hardware-kernel-invocation-queue` option only when your kernel execution time is much greater than the system and OpenCL runtime environment overhead hidden by the kernel invocation queue (20-100us), or if the Intel FPGA SDK for OpenCL Offline Compiler has a difficulty fitting your kernel.

Refer to *Utilizing Hardware Kernel Invocation Queue* topic in *Intel FPGA SDK for OpenCL Pro Edition: Best Practices Guide* for more information about how to utilize the kernel invocation queue.

### Related Information

[Utilizing Hardware Kernel Invocation Queue](#)

## 7.26. Modifying the Handshaking Protocol (-hyper-optimized-handshaking)

To modify the handshaking protocol used in certain areas of your design, use the `-hyper-optimized-handshaking=<auto|off>` option in your `aoc` command.

The `-hyper-optimized-handshaking` option can be set to one of the following values:

**auto** The default behavior without the option specified. The compiler enables the optimization if it is possible to do so, else it sets it to `off`.

Use this value when you want to achieve a higher  $f_{\max}$ . When you enable the optimization, the Intel FPGA SDK for OpenCL Offline Compiler adds pipeline registers to the handshaking paths of the stallable nodes. As a result, you might observe higher  $f_{\max}$  at the cost of increased area and latency.

Example: `aoc -hyper-optimized-handshaking <your_kernel_filename>.cl`

**off** The compiler attempts to optimize for lower latency at the potential cost of lower  $f_{\max}$ . Disabling hyper-optimized handshaking might also decrease area. This is useful for smaller designs where you are willing to give up  $f_{\max}$  for typically lower latency and area.

```
Example: aoc -hyper-optimized-handshaking=off  
<your_kernel_filename>.cl
```

*Remember:* The `-hyper-optimized-handshaking` option applies only to designs targeting Intel Stratix 10. If you use this option on target devices other than Intel Stratix 10 devices, the compiler fails and produces an error.

## 7.27. Pipelining Loops in Non-task Kernels (-auto-pipeline)

To direct the Intel FPGA SDK for OpenCL Offline Compiler to compile your design and pipeline loops in non-task (NDRange) kernels, include the `-auto-pipeline` option in your `aoc` command. The host program invokes non-task kernels through the kernel execution function `clEnqueueNDRangeKernel`.

At a command prompt, invoke the `aoc -auto-pipeline <your_kernel_filename1>.cl [<your_kernel_filename2>.cl ...]` command, where `[<your_kernel_filename2>.cl ...]` are the optional space-delimited file names of kernels that you can compile in addition to `<your_kernel_filename1>.cl`.

Without the `-auto-pipeline` option, the compiler does not pipeline the loops in non-task kernels, but instead executes different work items in parallel. With the `-auto-pipeline` option, the compiler attempts to pipeline the loops in your design, but the pipelining is not guaranteed.

*Note:* The `-auto-pipeline` option might improve or degrade performance depending on the memory access pattern in your design.

- If the auto-pipelining is successful, the [Loop Analysis](#) report displays the message `ND-Range as task`. The compiler-generated loops appear marked as `Compiler generated NDRange-as-task loop` in the report.
- If the compiler chooses not to auto-pipeline the loops, the Loop Analysis report displays a message for the kernel. The reasons for not auto-pipelining a loop can be one of the following:
  - A [barrier](#) in the function is not at the top-level function scope.
  - Kernel uses a local or private memory.
  - Kernel uses a volatile or atomic memory or channels.
  - Kernel uses the `num_compute_units(X,Y,Z)` attribute.

*Tip:* If you do not want the compiler to pipeline some infrequently used loops while allowing other loops to be auto-pipelined, use the [disable\\_loop\\_pipelining](#) pragma on specific loops when using the `-auto-pipeline` option. This pragma disables the loop pipelining.

## 8. Emulating and Debugging Your OpenCL Kernel

---

The Intel FPGA SDK for OpenCL Emulator assesses the functionality of your kernel.

The Intel FPGA SDK for OpenCL Emulator generates a `.aocx` file that executes on x86-64 Windows or Linux host. This feature allows you to emulate the functionality of your kernel and iterate on your design without executing it on the actual FPGA each time. For Linux platform, you can also use the Emulator to perform functional debug.

*Important:* Emulation does not support cross-compilation to ARM processor. To run emulation on a design that targets an ARM SoC device, emulate on a non-SoC board (for example, `<board_package_path>/s10_ref` or `<board_package_path>/a10_ref`). When you are satisfied with the emulation results, you may target your design on an SoC board for subsequent optimization steps.

The emulator supports 64-bit Windows and Linux operating systems. On Linux systems, the GNU C Library (glibc) version 2.15 or later is required. .

*Note:* You cannot use the execution time of an emulated design to estimate its execution time on an FPGA. Furthermore, running an emulated design is not a substitute for natively running a functionally equivalent C/C++ implementation on an x86-64 host.

1. [Setting up the Emulator](#) on page 140
2. [Modifying Channels Kernel Code for Emulation](#) on page 141
3. [Compiling a Kernel for Emulation \(-march=emulator\)](#) on page 144
4. [Emulating Your OpenCL Kernel](#) on page 144
5. [Debugging Your OpenCL Kernel on Linux](#) on page 146
6. [Limitations of the Intel FPGA SDK for OpenCL Emulator](#) on page 147
7. [Discrepancies in Hardware and Emulator Results](#) on page 148
8. [Emulator Environment Variables](#) on page 150
9. [Extensions Supported by the Emulator](#) on page 151
10. [Emulator Known Issues](#) on page 152

### 8.1. Setting up the Emulator

If you installed the Intel FPGA SDK for OpenCL Pro edition with administrator privileges, no additional setup is needed. If you did not install the Intel FPGA SDK for OpenCL with administrator privileges, you must perform some additional steps to enable the emulator.

If you did not install the Intel FPGA SDK for OpenCL with administrator privileges, manually set up the emulator installable client driver (ICD) entry:

- **Linux:** Ensure that the file `/etc/OpenCL/vendors/Intel_FPGA_SSG_Emulator.icd` matches the file found in the directory that the environment variable `INTELFPGAOCSDKROOT` specifies. The `INTELFPGAOCSDKROOT` environment variable points to the location of the SDK installation.  
  
If the files do not match, or if it is missing from `/etc/OpenCL/vendors`, copy the `Intel_FPGA_SSG_Emulator.icd` file from the location specified by the `INTELFPGAOCSDKROOT` environment variable to the `/etc/OpenCL/vendors` directory.
- **Windows:** Ensure that the registry key `HKEY_LOCAL_MACHINE\SOFTWARE\Khronos\OpenCL\Vendors` contains the following value:  
`[HKEY_LOCAL_MACHINE\SOFTWARE\Khronos\OpenCL\Vendors]`  
`"intelocl64_emu.dll"=dword:00000000`

**Note:** The emulator in Intel FPGA SDK for OpenCL Pro Edition is built with GCC 7.2.0 as part of the offline compiler. When executing the host program for an emulated OpenCL device, the version of `libstdc++.so` must be at least that of GCC 7.2.0. In other words, the `LD_LIBRARY_PATH` environment variable must ensure that the correct version of `libstdc++.so` is found.

If the correct version of `libstdc++.so` is not found, the call to `clGetPlatformIDs` function fails to load the FPGA emulator platform and returns `CL_PLATFORM_NOT_FOUND_KHR` (error code -1001). Depending on which version of `libstdc++.so` is found, the call to `clGetPlatformIDs` may succeed, but a later call to the `clCreateContext` function may fail with `CL_DEVICE_NOT_AVAILABLE` (error code -2).

If the `LD_LIBRARY_PATH` does not point to a compatible `libstdc++.so`, use the following syntax to invoke the host program:

```
env LD_LIBRARY_PATH=<path to compatible libstdc++.so>:$LD_LIBRARY_PATH <host>  
[host arguments]
```

### Related Information

- [Linking Your Host Application to the Khronos ICD Loader Library](#) on page 107
- [OpenCL Installable Client Driver \(ICD\) Loader](#)
- [OpenCL Vector Addition Design Example](#)

## 8.2. Modifying Channels Kernel Code for Emulation

The Emulator emulates kernel-to-kernel channels. It does not support the emulation of I/O channels that interface with input or output features of your FPGA board. To emulate applications with a channel that reads or writes to an I/O channel, modify your kernel to add a read or write channel that replaces the I/O channel, and make the source code that uses it conditional.

The Intel FPGA SDK for OpenCL does not set the `EMULATOR` macro definition. You must set it manually either from the command line or in the source code.

Consider the following kernel example:

```
channel_ulong4 inchannel __attribute__((io("eth0_in")));

__kernel void send (int size) {
    for (unsigned i = 0; i < size; i++) {
        ulong4 data = read_channel_intel(inchannel);
        //statements
    }
}
```

To enable the Emulator to emulate a kernel with a channel that interfaces with an I/O channel, perform the following tasks:

1. Modify the kernel code in one of the following manner:

- Add a matching `write_channel_intel` call such as the one shown below.

```
#ifdef EMULATOR

__kernel void io_in (__global char * restrict arr, int size) {
    for (unsigned i = 0; i < size; i++) {
        ulong4 data = arr[i]; //arr[i] being an alternate data source
        write_channel_intel(inchannel, data);
    }
}
#endif
```

- Replace the I/O channel access with a memory access, as shown below:

```
__kernel void send (int size) {
    for (unsigned i = 0; i < size; i++) {
        #ifndef EMULATOR

            ulong4 data = read_channel_intel(inchannel);

        #else
            ulong4 data = arr[i]; //arr[i] being an alternate data source
        #endif
        //statements
    }
}
```

2. Modify the host application to create and start this conditional kernel during emulation.

### Related Information

[Implementing I/O Channels Using the io Channels Attribute](#) on page 52

## 8.2.1. Emulating a Kernel that Passes Pipes or Channels by Value

The Intel FPGA SDK for OpenCL Emulator supports a kernel that passes pipes or channels by value.

You may emulate a kernel that passes a channel or pipe by value, as shown in the following example:

```
channel_uint my_ch;

void my_function (channel_uint ch,
                 __global_uint * dst, int i)
{
    dst[i] = read_channel_intel(ch);
}
```

```
__kernel void
consumer (__global uint * restrict dst)
{
    for (int i=0;i<5;i++)
    {
        my_function(my_ch, dst, i );
    }
}
```

## 8.2.2. Emulating Channel Depth

When you compile your OpenCL kernel for emulation, the default channel depth is different from the default channel depth generated when your kernel is compiled for hardware. You can change this behavior when you compile your kernel for emulation with the `CL_CONFIG_CHANNEL_DEPTH_EMULATION_MODE` environment variable.

This `CL_CONFIG_CHANNEL_DEPTH_EMULATION_MODE` environment variable can take the following values:

`ignoredepth` All channels are given a channel depth chosen to provide the fastest execution time for your kernel emulation. Any explicitly set channel depth attribute is ignored.

This value is used by default if `CL_CONFIG_CHANNEL_DEPTH_EMULATION_MODE` environment variable is not set.

`default` Channels with an explicit depth attribute have their specified depth. Channels without a specified depth are given a default channel depth that is chosen to provide the fastest execution time for your kernel emulation.

`strict` All channel depths in the emulation are given a depth that matches the depth given for the FPGA compilation.

*Important:* For channels, the `CL_CONFIG_CHANNEL_DEPTH_EMULATION_MODE` environment variable should be set before you compile a kernel. For pipes, it should be set before you run the host program.

## 8.2.3. Emulating Applications with a Channel That Reads or Writes to an I/O Channel

The Intel FPGA SDK for OpenCL Emulator emulates kernel-to-kernel channels. It does not support the emulation of I/O channels that interface with input or output features of your FPGA board.

To emulate applications with a channel that reads or writes to an I/O channel, perform the following:

1. Modify your kernel to add a read or write channel that replaces the I/O channel.
2. Make the source code that uses the read or write channel conditional.

However, it is possible to emulate unmodified kernel code using the following procedures:

- **For input I/O channels**

1. Store input data to be transferred to the channel in a file with a name matching the `io` attribute in the channel declaration. Consider the following example:

```
channel_ulong4_inchannel __attribute__((io("eth0_in")));
```

2. Create a file named `eth0_in`.
3. Store the test input data in the `eth0_in` file.

- **For output I/O channels**

Output data is automatically written into a file with name equal to the channel `io` attribute.

### 8.3. Compiling a Kernel for Emulation (-march=emulator)

To compile an OpenCL kernel for emulation, include the `-march=emulator` option in your `aoc` command.

- Before you perform kernel emulation, perform the following tasks:
  - Verify that the environment variable `QUARTUS_ROOTDIR_OVERRIDE` points to Intel Quartus Prime Pro Edition software installation folder.
- Verify that the `LD_LIBRARY_PATH` environment variable setting includes all the paths described in the *Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables* section in the *Intel FPGA SDK for OpenCL Pro Edition Getting Started Guide*.
- To create kernel programs that are executable on x86-64 host systems, invoke the following command:

```
aoc -march=emulator <your_kernel_filename>.cl
```

- For Linux systems, the Intel FPGA SDK for OpenCL Offline Compiler offers symbolic debug support for the debugger.

The offline compiler debug support allows you to pinpoint the origins of functional errors in your kernel source code.

#### Related Information

- [Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables \(Windows\)](#)
- [Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables \(Linux\)](#)
- [Debugging Your OpenCL Kernel on Linux](#) on page 146

### 8.4. Emulating Your OpenCL Kernel

To emulate your OpenCL kernel, run the emulation `.aocx` file on the platform on which you built your kernel. The OpenCL Emulator uses a different OpenCL platform than when targeting FPGA hardware.



To emulate your kernel, perform the following steps:

1. Modify your host program to select the emulator OpenCL platform.

Select the emulation OpenCL platform in your host program by selecting platform with the following name:

```
Intel(R) FPGA Emulation Platform for OpenCL(TM)
```

2. Build a host application and link your host application to the Khronos ICD Loader Library. For more information, see *Linking Your Host Application to the Khronos ICD Loader Library*.
3. If necessary, move the `<your_kernel_filename>.aocx` file to a location where the host can find it easily, preferably the current working directory.
4. To run the host application for emulation:
  - For Windows:
    - a. Define the number of emulated devices by invoking the `set CL_CONFIG_CPU_EMULATE_DEVICES=<number_of_devices>` command.
    - b. Run the host application.
    - c. Invoke `set CL_CONFIG_CPU_EMULATE_DEVICES=` to unset the variable.
  - For Linux, invoke the `env CL_CONFIG_CPU_EMULATE_DEVICES=<number_of_devices> <host_application_filename>` command.

This command specifies the number of identical emulation devices that the Emulator needs to provide.

**Remember:** The emulation OpenCL platform (Intel(R) FPGA Emulation Platform for OpenCL(TM)) does not provide access to physical boards. Only the emulated devices are available.

**Tip:** If you want to use only one emulator device, you do not need to set the `CL_CONFIG_CPU_EMULATE_DEVICES` environment variable.

5. If you change your host or kernel program and you want to test it, only recompile the modified host or kernel program and then rerun emulation.

**Note:**

- The emulator in Intel FPGA SDK for OpenCL Pro Edition is built with GCC 7.2.0 as part of the offline compiler. When executing the host program for an emulated OpenCL device, the version of `libstdc++.so` must be at least that of GCC 7.2.0. In other words, the `LD_LIBRARY_PATH` environment variable must ensure that the correct version of `libstdc++.so` is found.

If the correct version of `libstdc++.so` is not found, the call to `clGetPlatformIDs` function fails to load the FPGA emulator platform and returns `CL_PLATFORM_NOT_FOUND_KHR` (error code -1001). Depending on which version of `libstdc++.so` is found, the call to `clGetPlatformIDs` may succeed, but a later call to the `clCreateContext` function may fail with `CL_DEVICE_NOT_AVAILABLE` (error code -2).

If `LD_LIBRARY_PATH` does not point to a sufficiently new `libstdc++.so`, use the following syntax to invoke the host program:

```
env LD_LIBRARY_PATH=<path to sufficiently new libstdc++.so>:$LD_LIBRARY_PATH  
<host> [host arguments]
```

- To enable debugging of kernel code, optimizations are disabled by default for the FPGA emulator. This can lead to sub-optimal execution speed when emulating kernel code.

You can pass the `-g0` flag to the `aoc compile` command to disable debugging and enable optimizations. This enables faster emulator execution.

**Related Information**

- [Linking Your Host Application to the Khronos ICD Loader Library](#) on page 107
- [Emulating Your OpenCL Kernel on Windows](#)

## 8.5. Debugging Your OpenCL Kernel on Linux

For Linux systems, you can direct the Intel FPGA SDK for OpenCL Emulator to run your OpenCL kernel in the debugger and debug it functionally as part of the host application. The debugging feature allows you to debug the host and the kernel seamlessly. You can step through your code, set breakpoints, and examine and set variables.

To compile your OpenCL kernel for debugging, perform the following steps:

1. To generate an emulator `.aocx` file for debugging, invoke the `aoc -march=emulator <your_kernel_filename>.cl` command.
2. Build the host application and link your host application to the Khronos ICD Loader Library. For details, see [Linking Your Host Application to the Khronos ICD Loader Library](#) on page 107.
3. Ensure that the `<your_kernel_filename>.aocx` file is in a location where the host expects to find it.
4. To run the application, invoke the command `env CL_CONFIG_CPU_EMULATE_DEVICES=<number_of_devices> gdb --args <your_host_program_name> [<host_program_arguments>]`, where `<number_of_devices>` is the number of identical emulation devices that the Emulator needs to provide.

*Note:* You can omit the `CL_CONFIG_CPU_EMULATE_DEVICES` environment variable if you only want to emulate a single device.

5. If you make changes to only your host application, recompile only the host before restarting the debugger. Similarly, if you make changes to only your kernel code, recompile only the kernel code before restarting the debugger.

To run your kernel in the debugger:

1. During program execution, the debugger cannot step from the host code to the kernel code. You must set a breakpoint before the actual kernel invocation by adding these lines:
  - a. `break <your_kernel>`

This line sets a breakpoint before the kernel.
  - b. `continue`

If you have not begun debugging your host, then type `start` instead.
2. The debugger does not recognize kernel names until the host actually loads the kernel functions. As a result, the debugger generates the following warning for the breakpoint you set before the execution of the first kernel:

```
Function "<your_kernel>" not defined.
```

```
Make breakpoint pending on future shared library load? (y or [n])
```

Answer `y`. After initial program execution, the debugger recognizes the function and variable names, and line number references for the duration of the session.

*Notice:* When debugging kernel code with the FPGA emulator, the debugger may display some source lines out of order for certain code patterns. Currently, there is no known workaround for this issue.

#### Related Information

[Linking Your Host Application to the Khronos ICD Loader Library](#) on page 107

## 8.6. Limitations of the Intel FPGA SDK for OpenCL Emulator

The Intel FPGA SDK for OpenCL Emulator feature has some limitations.

- Execution model

The Emulator supports the same compilation modes as the FPGA variant. As a result, you must call the `clCreateProgramBinary` function to create `cl_program` objects for emulation.
- Concurrent execution

Modeling of concurrent kernel executions has limitations. During execution, the Emulator does not actually run interacting work-items in parallel. Therefore, some concurrent execution behaviors, such as different kernels accessing global memory without a barrier for synchronization, might generate inconsistent emulation results between executions.
- Same address space execution

The Emulator executes the host runtime and the kernels in the same address space. Certain pointer or array usages in your host application might cause the kernel program to fail, and vice versa. Example usages include indexing external allocated memory and writing to random pointers. You may use memory leak detection tools such as Valgrind to analyze your program. However, the host might encounter a fatal error caused by out-of-bounds write operations in your kernel, and vice versa.
- Conditional channel operations

Emulation of channel behavior has limitations, especially for conditional channel operations where the kernel does not call the channel operation in every loop iteration. In these cases, the Emulator might execute channel operations in a different order than on the hardware.
- GCC version

Emulator host programs on Linux\* must be run with a version of `libstdc++.so` from GCC 7.2.0 or later. This can be achieved either by installing GCC 7.2.0 or later on your system, or setting the `LD_LIBRARY_PATH` environment variable such that a compatible `libstdc++.so` is identified.

## 8.7. Discrepancies in Hardware and Emulator Results

When you emulate a kernel, your OpenCL system might produce results different from that of the kernel compiled for hardware. You can further debug your kernel before you compile for hardware by running your kernel through simulation.

**Warning:** These discrepancies usually occur when the Intel FPGA SDK for OpenCL Emulator is unable to model some aspects of the hardware computation accurately, or when your program relies on an undefined behavior.

The most common reasons for differences in emulator and hardware results are as follows:

- Your OpenCL kernel code is using the `#pragma ivdep` directive. The Emulator does not model your OpenCL system when a true dependence is broken by a `pragma ivdep` directive. During a full hardware compilation, you observe this as an incorrect result.
- Your OpenCL kernel code is relying on uninitialized data. Examples of uninitialized data include uninitialized variables and uninitialized or partially initialized global buffers, local arrays, and private arrays.
- Your OpenCL kernel code behavior depends on the precise results of floating point operations. The Emulator uses floating point computation hardware of the CPU whereas the hardware run uses floating point cores implemented as FPGA cores. The use of `-ffp-reassociate aoc` option in your OpenCL kernel code might change the order of operations leading to further divergence in the floating point results.

*Note:* The OpenCL standard allows one or more least significant bits of floating point computations to differ between platforms, while still being considered *correct* on both such platforms.

- Your OpenCL kernel code behavior depends on the order of channel accesses in different kernels. The emulation of channel behavior has limitations, especially for conditional channel operations where the kernel does not call the channel operation in every loop iteration. In such cases, the Emulator might execute channel operations in an order different from that on the hardware.
- Your OpenCL kernel or host code is accessing global memory buffers out-of-bounds.

- Attention:**
- Uninitialized memory read and write behaviors are platform-dependent. Verify sizes of your global memory buffers when using all addresses within kernels, allocating `clCreateBuffer` function call, and transferring `clEnqueueReadBuffer` and `clEnqueueWriteBuffer` function calls.
  - You can use software memory leak detection tools, such as Valgrind, on the emulated version of your OpenCL system to analyze memory related problems. Absence of warnings from such tools does not mean the absence of problems. It only means that the tool could not detect any problem. In such a scenario, Intel recommends manual verification of your OpenCL kernel or host code.

- Your OpenCL kernel code is accessing local or private variables out-of-bounds. For example, accessing a local or private array out-of-bounds or accessing a private variable after it has gone out of scope.  
**Attention:** In software terms, these issues are referred to as stack corruption issues because accessing variables out-of-bounds usually affects unrelated variables located close to the variable being accessed on a software stack. Emulated OpenCL kernels are implemented as regular CPU functions, and have an actual stack that can be corrupted. When targeting hardware, no stack exists and hence, the stack corruption issues are guaranteed to manifest differently. You may use memory leak analyzer tools, such as Valgrind, when a stack corruption is suspected. However, stack related issues are usually difficult to identify. Intel recommends manual verification of your OpenCL kernel code to debug a stack related issue.
- Your OpenCL kernel code is using shifts that are larger than the type being shifted. For example, shifting a 64-bit integer by 65 bits. According to the OpenCL specification version 1.0, the behavior of such shifts is undefined.
- When you compile your OpenCL kernel for emulation, the default channel depth is different from the default channel depth generated when your kernel is compiled for hardware. This difference in channel depths might lead to scenarios where execution on the hardware hangs while kernel emulation works without any issue. Refer to [Emulating Channel Depth](#) on page 143 for information on how to fix the channel depth difference.
- In terms of ordering the printed lines, the output of the `printf` function might be ordered differently on the Emulator and hardware. This is because, in the hardware, `printf` data is stored in a global memory buffer and flushed from the buffer only when the kernel execution is complete, or when the buffer is full. In the Emulator, the `printf` function uses the x86 `stdout`.
- If you perform an unaligned load/store through upcasting of types, the FPGA and emulator might produce different results. A load/store of this type is undefined in the C99 specification.

For example, the following operation might produce unexpected results:

```
int tmp = *((int *) (my_ptr + 5));
```

### Related Information

[Debugging Your OpenCL Library Through Simulation \(Preview\)](#) on page 194

## 8.8. Emulator Environment Variables

Several environment variables are available to modify the behavior of the emulator.

### CL\_CONFIG\_CPU\_EMULATE\_DEVICES

Controls the number of identical emulator devices provided by the emulator platform. If not set, a single emulator device is available. Therefore, set this variable only if you want to emulate multiple devices.

### **CL\_CONFIG\_CPU\_TBB\_NUM\_WORKERS**

Indicates a maximum number of threads that can be used by the emulator. The default value is 32, and the maximum value is 255. Each thread can run a single kernel.

If the application requires several kernels to be executing simultaneously, the `CL_CONFIG_CPU_TBB_NUM_WORKERS` should be set appropriately (to the number of kernels used or a higher value).

### **CL\_CONFIG\_CPU\_FORCE\_LOCAL\_MEM\_SIZE**

Set the amount of available OpenCL local memory, with units, for example: 8MB, 256KB, or 1024B.

### **CL\_CONFIG\_CPU\_FORCE\_PRIVATE\_MEM\_SIZE**

Set the amount of available OpenCL private memory, with units, for example: 8MB, 256KB, or 1024B.

### **CL\_CONFIG\_CHANNEL\_DEPTH\_EMULATION\_MODE**

When you compile your OpenCL kernel for emulation, the channel depth is different from the channel depth generated when your kernel is compiled for hardware. You can change this behavior with the `CL_CONFIG_CHANNEL_DEPTH_EMULATION_MODE` environment variable. For details, see [Emulating Channel Depth](#) on page 143.

#### **Related Information**

[Emulating Channel Depth](#) on page 143

## **8.9. Extensions Supported by the Emulator**

The emulator offers varying levels of support for different OpenCL extensions.

The following OpenCL extensions are fully supported by the emulator:

- `cl_intel_fpga_host_pipe`
- `cl_khr_byte_addressable_store`
- `cl_khr_icd`
- `cles_khr_int64`

The emulator also supports the following OpenCL extensions:

- `cl_intel_channels`
- `cl_khr_local_int32_base_atomics`
- `cl_khr_local_int32_extended_atomics`
- `cl_khr_global_int32_base_atomics`
- `cl_khr_global_int32_extended_atomics`
- `cl_khr_fp64`
- `cl_khr_fp16`

## 8.10. Emulator Known Issues

A few known issues might affect your use of the emulator. Review these issues to avoid possible problems when using the emulator.

### Autorun Kernels

Autorun kernels shut down only after a host program exits, not after a `clReleaseProgram()` call.

### Compiler Diagnostics

Some compiler diagnostics are not yet implemented for the emulator.

### CL\_OUT\_OF\_RESOURCES Error Returned From `clEnqueueNDRangeKernel()`

This can occur when the kernel used more `__private` or `__local` memory than the emulator supports by default.

Try setting the `CL_CONFIG_CPU_FORCE_PRIVATE_MEM_SIZE` or the `CL_CONFIG_CPU_FORCE_LOCAL_MEM_SIZE` environment variables, as described in [Emulator Environment Variables](#) on page 150.

### CL\_INVALID\_VALUE Error Returned From `clCreateKernel()`

It is possible a call to the `clBuildProgram()` was missed.

This call is required by the OpenCL specification, even if a program is created from a binary. See section 5.4.2 of the OpenCL Specification version 1.0 for details.

### Related Information

- [Emulator Environment Variables](#) on page 150
- [The OpenCL Specification, Version 1.0](#)



## 9. Developing OpenCL Applications Using Third-party IDEs

Intel FPGA SDK for OpenCL helps to create, build, debug, and analyze OpenCL applications developed with the OpenCL API on Microsoft Windows\* and Linux\* OSes. It integrates with third-party integrated development environments (IDEs) on Linux\* and Windows\* to provide a seamless GUI experience for software development.

*Note:* If you are looking for instructions to work with Intel oneAPI toolkits for FPGA development, refer to [Intel® oneAPI DPC++ FPGA Workflows on Third-Party IDEs](#).

### 9.1. FPGA Workflows in Microsoft Visual Studio

Plugins for Microsoft Visual Studio\* 2017 and 2019 versions are automatically installed as part of the Intel FPGA SDK for OpenCL installation. For more information about the system requirements, refer to the [Prerequisites for the Intel FPGA SDK for OpenCL Pro Edition](#).

#### 9.1.1. Preparing the Visual Studio Environment

Before you run Visual Studio, ensure that you have performed these steps:

1. [Downloaded the Intel FPGA SDK for OpenCL Pro Edition software](#).
2. [Installed the Intel FPGA SDK for OpenCL Pro Edition software](#).
3. [Set the user environment variables](#).
4. [Installed an FPGA Board and set the `QUARTUS\_ROOTDIR\_OVERRIDE` user environment variable](#).

Visual Studio project displays an error if you have not set the environment. If you set the environment while Visual Studio is running, then reload Visual Studio.

Once you have completed the prerequisites, verify the environment by opening **Edit the system environment variables** > **Environment Variables** and search for environment through **System variables**.

#### Related Information

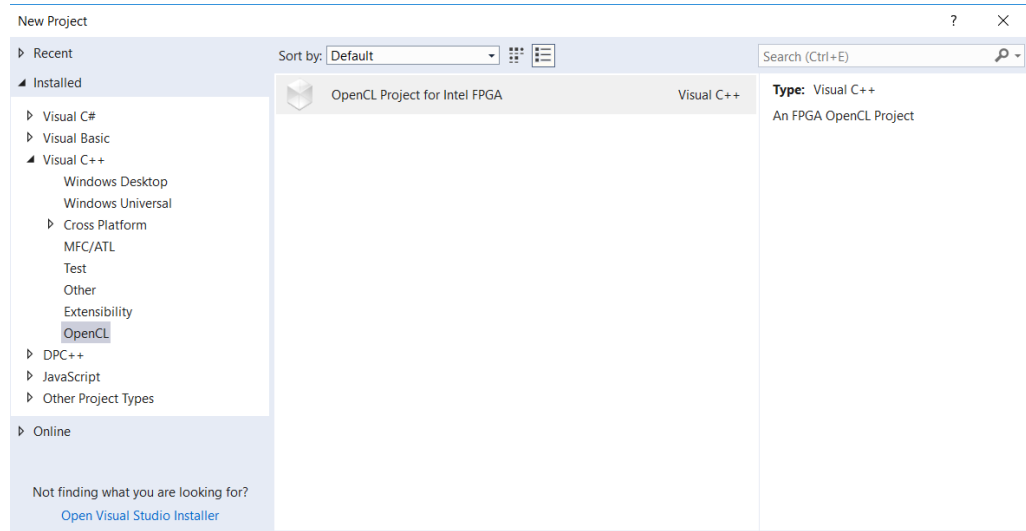
- [Getting Started with the Intel FPGA SDK for OpenCL Pro Edition for Windows](#)
- [Managing an FPGA Board](#) on page 19

#### 9.1.2. Creating an FPGA OpenCL Template

To create an FPGA OpenCL Template, perform these steps:

1. Navigate to the **File** > **New** > **Project** > **Visual C++** > **OpenCL** menu option.
2. Choose **OpenCL Project for Intel FPGA**. A new project is created. This project supports Intel FPGA by default.

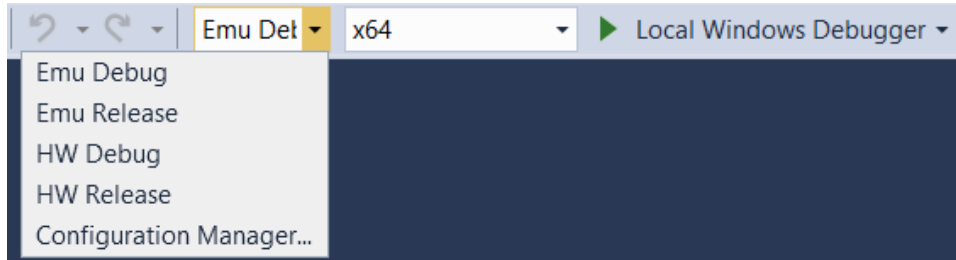
**Figure 17. Creating a New Project in Visual Studio**



### 9.1.3. Configuring the Build Targets

A default project supported by an FPGA has four different configurations, as shown in the following image:

**Figure 18. Build Targets Supported by an FPGA**



In [Build Targets Supported by an FPGA](#), Emu stands for FPGA Emulator configurations and HW stands for FPGA hardware configurations. These configurations allow you to configure different build targets the compiler supports.

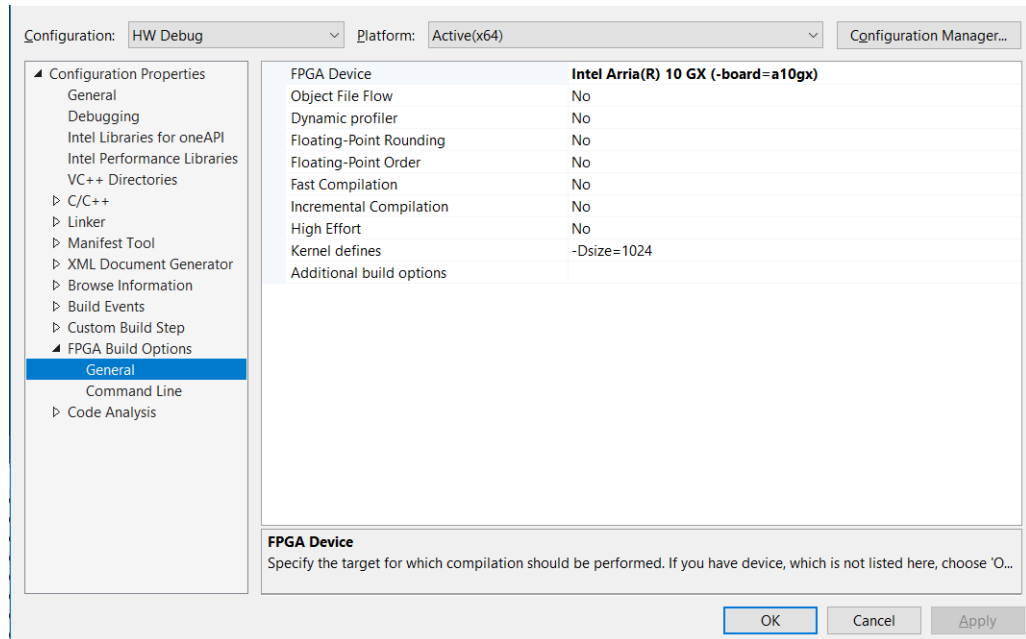
**Note:** The difference between the Debug and the Release target lies in the ability to debug the host code. The Debug target allows debugging of the host code only and not the kernel.

For more information about Emulation, refer to [Verifying Host Runtime Functionality via Emulation](#) topic in the *Intel FPGA SDK for OpenCL Getting Started Guide* and [Emulating and Debugging Your OpenCL Kernel](#) on page 140.

### 9.1.4. Configuring Build Options for a Project

To configure FPGA-specific build options, perform these steps:

1. Navigate to **Project > Properties > Configuration Properties > FPGA Build Options > General**. Build options are displayed.
2. Modify the build options as you desire.

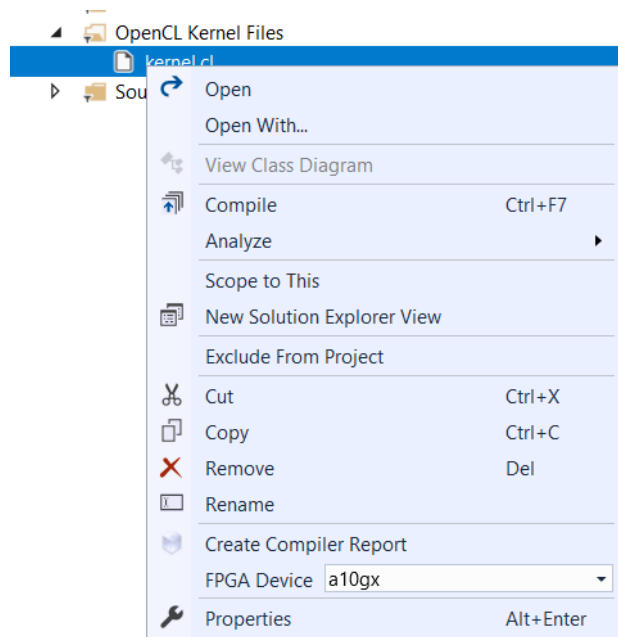


For more information about the build options, refer to [Compiling Your OpenCL Kernel](#) on page 119.

### 9.1.5. Generating the High-level Design Report

To generate the high-level design report, perform these steps:

1. Locate the `kernel.cl` file under the **OpenCL Kernel Files** folder.
2. Right-click on the `kernel.cl` file to view the context-sensitive menu.



3. In the **FPGA Device** context-sensitive menu option, select the FPGA board name you have installed from the drop-down list.

If you have installed a board but you do not see the board name in the drop-down list, then enter the board name manually.

*Note:* To obtain a list of the available devices, you can run the `aoc.exe -list-boards` command.

4. In the same context-sensitive menu, select the **Create Compiler Report** option to generate the high-level design report. After the report is generated, it opens in your default browser.

*Note:* If you select multiple OpenCL Kernel files, the **Create Compiler Report** option creates a combined report.

For more information about the high-level design report, refer to the [Reviewing Your Kernel's report.html File](#) in the *Intel FPGA SDK for OpenCL Best Practices Guide*

### 9.1.6. Building and Running the FPGA Template

Building and running the FPGA template with Visual Studio follows the same basic workflow as that of the common Visual Studio projects. For more information, refer to the [Build and run a C++ console app project](#) topic in the *Visual Studio C++ Tutorials* documentation.

## 9.2. FPGA Workflows in Eclipse

Extension for Eclipse is automatically installed as part of the Intel FPGA SDK for OpenCL installation. For more information about the system requirements, refer to [Prerequisites for the Intel FPGA SDK for OpenCL Pro Edition](#).

## 9.2.1. Preparing the Eclipse Environment

To apply transient environment variable settings, you must source `INTELFPGAOCSDKROOT/init_openc1.sh` script from the Intel FPGA SDK for OpenCL installation directory.

## 9.2.2. Creating a Simple FPGA application

In this tutorial, an Eclipse C/C++ Development Tooling (CDT) is used to create a simple Hello World application. This tutorial describes the process of creating a new OpenCL FPGA project (where the build is automatically managed by the CDT) and running the program.

To create a simple Hello World application using CDT, perform the following general steps:

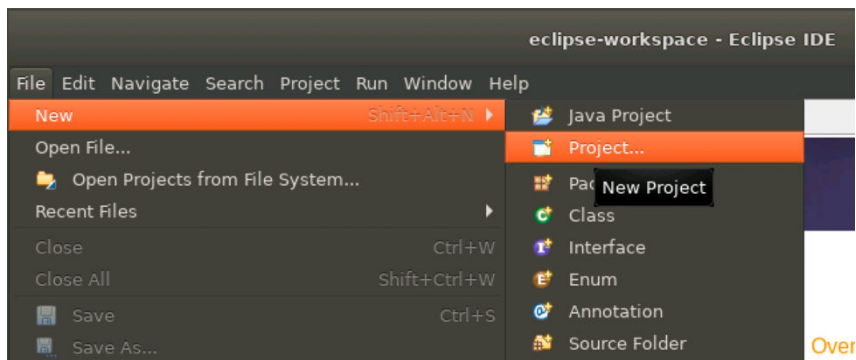
1. [Creating a Project](#) on page 157
2. [Reviewing the Code and Building the Project](#) on page 161
3. [Running the Application](#) on page 162

### 9.2.2.1. Creating a Project

Perform the following steps to create a project using Eclipse CDT:

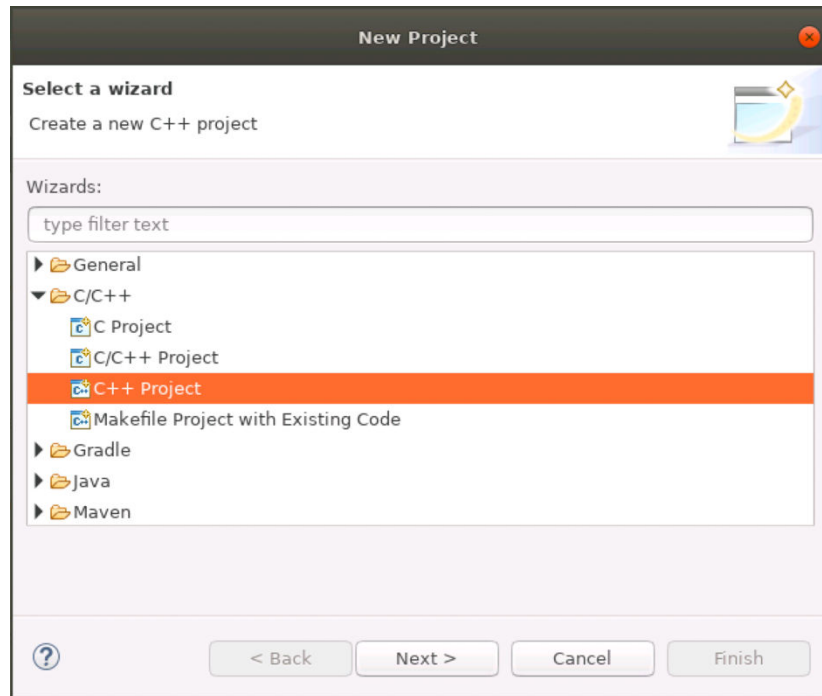
1. Select **File > New > Project** menu option.

**Figure 19. Eclipse CDT File Menu**



2. Select the type of project you want to create. For this tutorial, expand the **C/C++** folder, select **C++ Project**, and click the **Next** button.

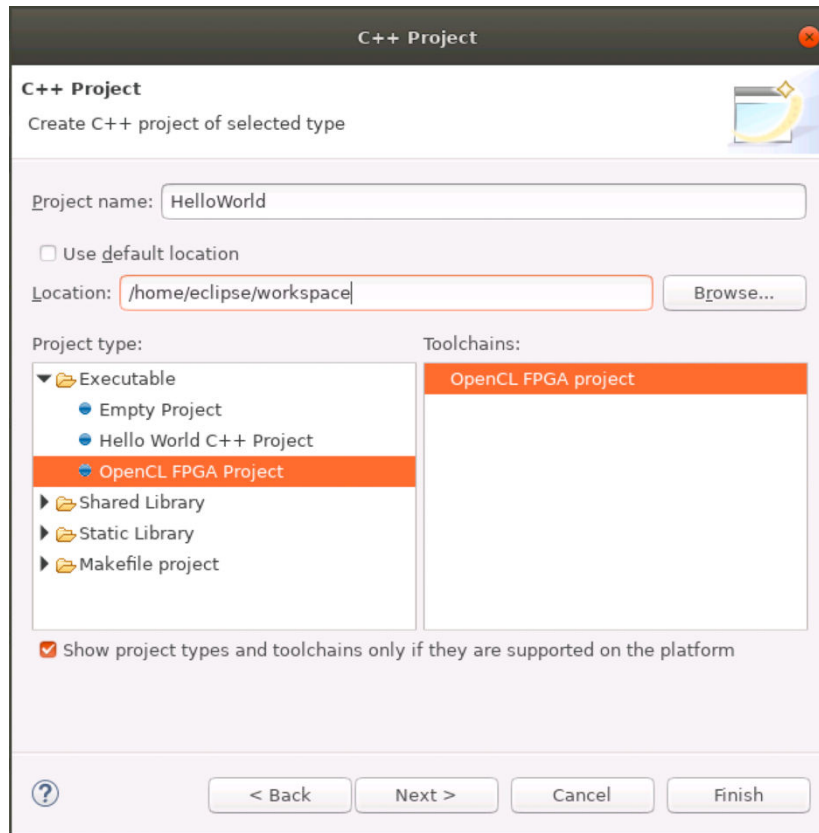
Figure 20. New Project Wizard



The **C++ Project** wizard opens. By default, the CDT filters the **Project Type** and **Toolchain** based on the language supported for the C++ Project wizard you selected for this tutorial.

3. In the **C++ Project** wizard, perform these steps:

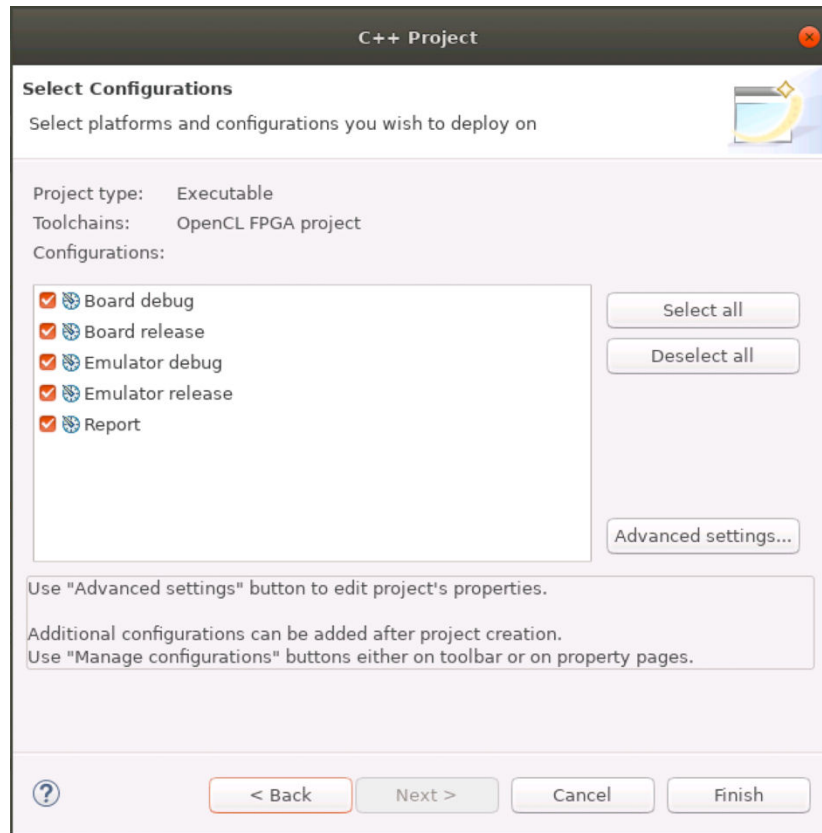
Figure 21. C++ Project Wizard



- In the **Project name** field, type a name for the project. For example, HelloWorld.
- In the **Project type** list, expand **Executable** folder and select **OpenCL FPGA Project**. This project type provides a simple Hello World application in OpenCL and the makefile is automatically created by the CDT.
- Under **Toolchains**, select **OpenCL FPGA project**.
- Click the **Next** button.

The **Select Configurations** dialog box displays a list of configurations based on the project type and toolchain selected earlier.

Figure 22. Select Configurations Dialog



*Note:* Board debug and Board release configuration options are equivalent to HW Debug and HW Release build target options in the Microsoft Visual Studio IDE. When you use these options, board compiles for hardware and executes on a board.

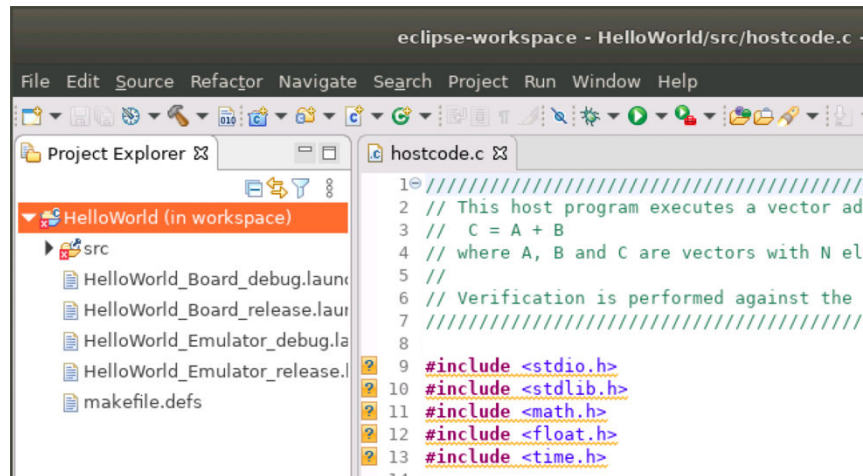
- Optional: If you want to change the default project settings, click the **Advanced settings** button. The **Project Properties** dialog launches for your new project, allowing you to change any of the project specific settings, such as includes paths, compiler options, and libraries.
- Click the **Finish** button.

*Note:* If the **C++ perspective** is not currently set as the default, you are prompted to determine if you want this project to be associated with the C/C++ perspective. Click **Yes**.

A project is created with the default settings and a full set of configurations based on the project type and toolchain you selected. You should now see the new project in the **Project Explorer** view.



Figure 23. New Project Displayed in the Project Explorer View



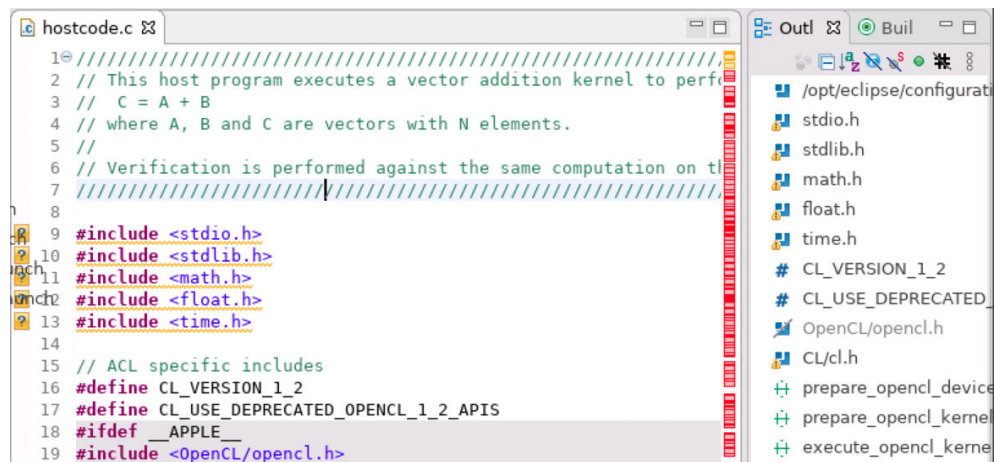
### 9.2.2.2. Reviewing the Code and Building the Project

Perform these steps to review your host code and build your project:

1. From the **Project Explorer** view, double-click the `.c` file created for your project (for example, `hostcode.c`).

You can find the `.c` file within the project `src` folder. It opens in a default editor and contains C++ template host code for the Hello World example project you selected earlier. In addition, the Outline view is also populated with objects created from your code.

Figure 24. Reviewing Code in the Project Explorer



*Tip:* You can specify a different editor and add or modify existing code templates by modifying your preferences in **Window > Preferences**.

2. Optional: Type additional code in this file and save the changes either by clicking **File > Save** or pressing **CTRL+S** keys combination.
3. Build your project using one of the following options:
  - Press **CTRL+B** keys combination.

- Select the project in the **Project Explorer** view and navigate to **Project > Build Project**.

*Note:* If a build generates any error or warning, you can view it in the **Problems** view. If you encounter difficulty, refer to the following topics in the Eclipse\* documentation:

- [Problem View](#)
- [Building C/C++ Projects](#)
- [Before you begin](#)

4. Review the build messages in the **Console** view. The project should build successfully.

You can observe that the **Outline** view has also been populated with objects created from your code. If you select an item from the **Outline** view, the corresponding text in the editor is highlighted.

### 9.2.2.3. Running the Application

To run the application, follow these steps:

1. Within the C/C++ Perspective, click **Project > Build Configurations > Set Active** and select the required configuration. For this tutorial, select the **Emulator debug** option.
2. Build your project either by clicking **Project > Build All** or by pressing **CTRL+B** keys combination.
3. Run the application either by clicking **Run > Run** or by pressing **CTRL+F11** keys combination.  
Now, you should have the Hello World application running in the **Console** view. The Console's title bar also displays which application is running.

### 9.2.3. Creating a Makefile Project

This tutorial describes the process of creating a new C++ project (includes a `makefile`) and debugging the program. You must create a project to contain your source code and related files. A project has an associated builder that can incrementally compile source files as they are changed.

To create a `makefile` project, perform the following steps:

1. Select **File > New > Project**.  
When you create a new project, you are required to specify the project type. This project type determines the toolchain, data, and tabs that the CDT uses or displays.
2. Select the type of project you want to create. For this tutorial, expand the **C/C++** folder and select the **C++ Project**. The C++ Project wizard launches.  
By default, the CDT filters the Project types and Toolchain list based on the language supported for the C++ Project wizard you selected for this tutorial.
3. In the **Project name** field, type `HelloWorld`.
4. Leave the **Use Default Location** option selected.

5. From the **Project types** list, expand the **Makefile** project and select **OpenCL FPGA Makefile Project**. This project allows you to enter the source file and the `makefile`.
6. From the **Toolchain** list, select the **FPGA makefile**.
7. Click **Next**.
8. Click **Finish**.

If a message box prompts you to change perspectives, click **Yes**.

Your new project displays in the **Project Explorer** view. The project is empty because you have not yet created files for your project. You may see an error since there is nothing to build yet for your project. You can now start writing the code for your HelloWorld program.

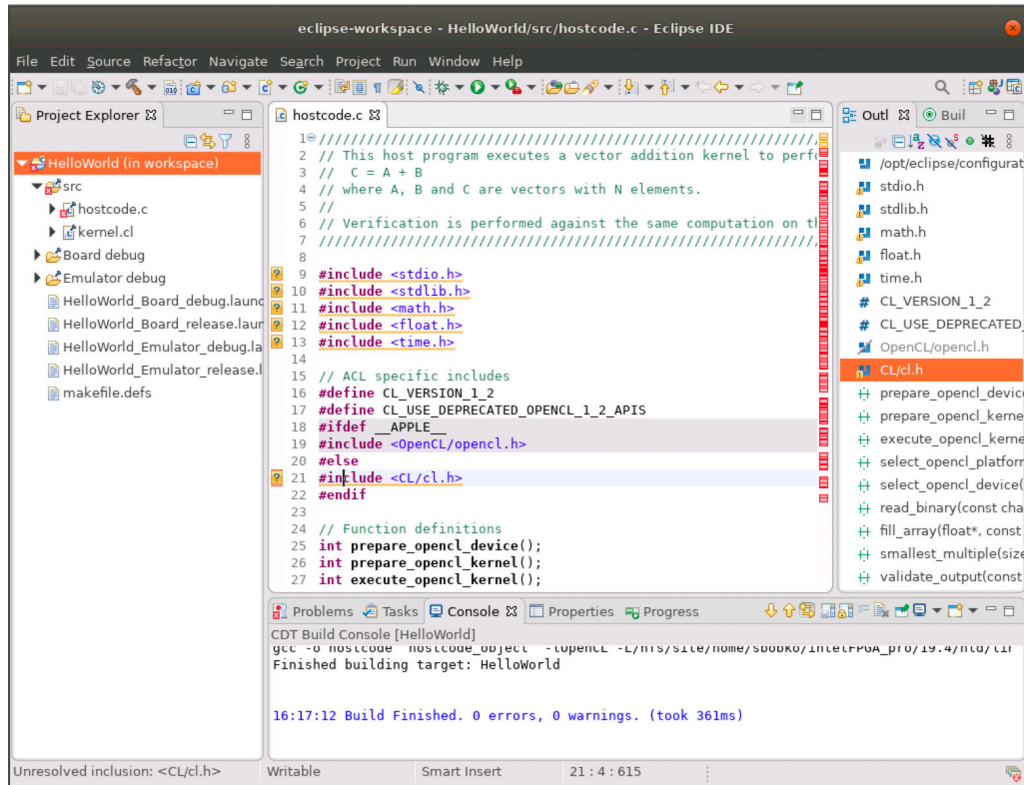
### 9.2.4. Building a Project

To build a project, perform these steps:

1. In the **Project Explorer** view, select your project. For this tutorial, you can select the **HelloWorld** project you created earlier.
2. To build the project, either click **Project > Build Project** or click the **Build** icon on the toolbar.  
In the **Console** view, you can view the output and results of the build command.
3. Click the **Console** view's tab to bring the view forward if it is not currently visible.

For some unknown reason, if you are unable to view the Console, then open it by selecting **Window > Show View > Console**.

Figure 25. Eclipse Workspace



- Tip:*
- The **Console** view displays the make output and build progress information.
  - The **Make Targets** view displays the makefile actions.
  - The **Problems** view displays compile warnings or errors.

### 9.3. Limitations

The following are the limitations when working with third-party IDEs:

- The Intel FPGA SDK for OpenCL plug-in for Visual Studio does not provide all of the build information in the **Output** window from the compiler by default. To view the full output from the compiler, go to **Tools > Options > Projects and Solutions > Build and Run** and set the output and log file verbosity to **Detailed**.
- The Intel FPGA SDK for OpenCL extension for Eclipse requires you to launch an Eclipse instance from the terminal session by sourcing the `init_opencl.sh` script.
- The C/C++ code editor in Eclipse does not include highlighting of several OpenCL-specific keywords. Additionally, it reports syntax errors for OpenCL functions in the host code.
- The compiler reports are generated for any build configuration in the Intel FPGA SDK for OpenCL extension for Eclipse. These reports are not available in the **Project Explorer** window by default. You must refresh the file tree manually.

## 10. Developing OpenCL Applications Using Intel Code Builder for OpenCL

---

The Intel Code Builder for OpenCL is a software development tool available as part of the Intel FPGA SDK for OpenCL. It enables development of OpenCL applications via well-known integrated development environments targeting the Intel FPGAs.

The Intel Code Builder for OpenCL provides a set of Microsoft Visual Studio and Eclipse plug-ins that enable capabilities for creating, building, debugging, and analyzing Windows and Linux applications accelerated with OpenCL.

### 10.1. Configuring the Intel Code Builder for OpenCL Offline Compiler Plug-in for Microsoft Visual Studio

To enable the Intel Code Builder for OpenCL offline compiler plug-in for Microsoft Visual Studio, perform the following steps:

1. In the Visual Studio software, select **Project** > **Properties**.
2. In the **Project** > **Properties** > **Code Builder** page, change the **Device** to your desired FPGA device.
3. In the **C/C++** > **General** property page, under **Additional Include Directories**, enter the full path to the directory where the OpenCL code header files are located (`$(INTELFPGAOCCLSDKROOT)\include`).
4. In the **Linker** > **General** property page, under **Additional Library Directories**, enter the full path to the directory where the OpenCL code run-time import library file is located. For example, for 64-bit application, add `$(INTELFPGAOCCLSDKROOT)\lib\x64`:
5. In the **Linker** > **Input** property page, under **Additional Dependencies**, enter the name of the OpenCL ICD import library file as `OpenCL.lib`.

### 10.2. Configuring the Intel Code Builder for OpenCL Offline Compiler Plug-in for Eclipse

To enable the Intel Code Builder for OpenCL offline compiler plug-in for Eclipse IDE, perform the following steps:

1. Copy the `CodeBuilder_<version>.jar` plug-in file from `$(INTELFPGAOCCLSDKROOT)/eclipse-plug-in` to `<ECLIPSE_ROOT_FOLDER>/dropins`.

**Attention:** In Linux, you must add `$INTELFPGAOCCLSDKROOT\bin` to the `LD_LIBRARY_PATH` environment variable.

2. Run the Eclipse IDE.
3. Select **Windows > Preferences**.
4. Switch to the Intel OpenCL dialog box.
5. Set the OpenCL binary directory to `$INTELFPGAOCCLSDKROOT/bin`.

Once the offline compiler is configured, you can use the **Code-Builder** menu to perform the following basic operations:

- Create a new session
- Open an existing session
- Save a session
- Build a session
- Compile a session
- Configure a session

For more information about the Intel Code Builder for OpenCL, refer to *Developer Guide for Intel SDK for OpenCL Applications*. For information about how to configure the Intel Code Builder for OpenCL for Microsoft Visual Studio, refer to *Intel Code Builder for OpenCL API for Microsoft Visual Studio*. For information about how to configure the Intel Code Builder for OpenCL for Eclipse, refer to *Intel Code Builder for OpenCL API for Eclipse*.

#### Related Information

- [Developer Guide for Intel SDK for OpenCL Applications 2017](#)
- [Intel Code Builder for OpenCL API Plug-in for Microsoft Visual Studio](#)
- [Intel Code Builder for OpenCL API Plug-in for Eclipse](#)

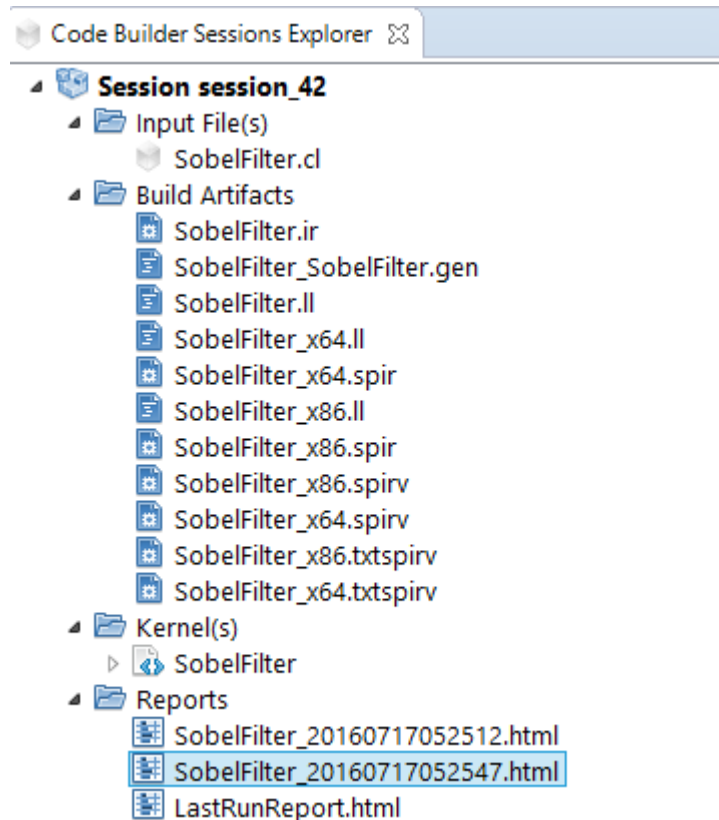
### 10.3. Creating a Session in the Intel Code Builder for OpenCL

Perform the following steps to create a session in the Intel Code Builder for OpenCL:

1. Select **Code-Builder > OpenCL Kernel Development > New Session**.
2. Specify the session name, path to the folder to store the session file and the content of the session (can be either an empty session or with a predefined OpenCL code).
3. Click **Done**.

Once the session is created, the new session appears in the **Code Builder Sessions Explorer** view.

Figure 26. Code Builder Sessions Explorer



*Note:* If you do not see the **Code Builder Session Explorer** view, select **Code-builder > OpenCL Kernel Development > Windows > Code Builder Session Explorer**.

## 10.4. Configuring a Session

A configuration is a set of analysis inputs such as assigned variables, number of iterations, global sizes and local sizes of a specific kernel, and so on. You can create a separate configuration for each set of inputs that you want to analyze.

You can configure a session by right-clicking the session in the **Code Builder Session Explorer** and selecting **Session Options**. Alternatively, you can also open the **Session Settings** dialog box by selecting **Code-Builder > OpenCL Kernel Development > Session Options**.

The **Session Settings** dialog box allows you to configure:

- Device options such as target machine, OpenCL platform, and OpenCL device.
- Build options such as offline compiler flags and build architecture.
- Build artifacts such as `.aocx` and `.aoco` files, and static reports.
- General options such as job architecture and network settings.

In the **Device Options** tab, ensure to select **Intel FPGA SDK for OpenCL** in the OpenCL platform drop-down list.

Under the **Build Options** tab, in the **OpenCL Build Options** section, enter the Intel FPGA SDK for OpenCL Offline Compiler flags manually.

**Attention:** If your kernel has channels, you must configure workflows. A workflow is a set of kernels, which can be executed sequentially. Workflow can be used to execute a workload with channels where you connect the input of one kernel with the output of the previous kernel (by assigning the same variable for both kernels).

For more information about configuring a session and variable management, refer to the *Developer Guide for Intel SDK for OpenCL Applications*.

#### Related Information

- [Configuring a Session in Microsoft Visual Studio](#)
- [Configurations and Settings in Eclipse](#)
- [Variable Management in Microsoft Visual Studio](#)
- [Variable Management in Eclipse](#)



## 11. Intel FPGA SDK for OpenCL Advanced Features

The Intel FPGA SDK for OpenCL provides advanced features you can use to control the following aspects of the design architecture and the Intel FPGA SDK for OpenCL Offline Compiler's behavior:

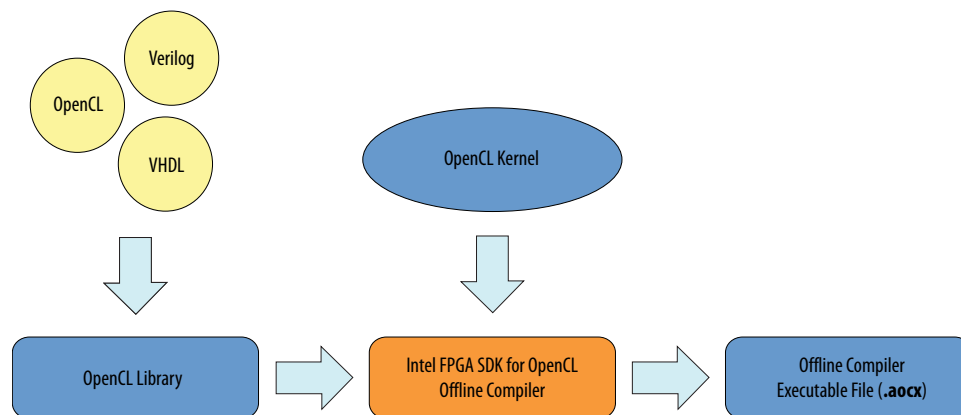
### 11.1. OpenCL Library

An OpenCL library is a single file that contains multiple functions. Each function is comprised of data processing logic that works at any clock frequency.

*Note:* To understand how to create a static library file using SYCL, refer to [Use of Static Library for FPGA](#) in the *Intel oneAPI Programming Guide*.

You can create an OpenCL library in OpenCL or register transfer level (RTL). You can then include this library file and use the functions inside your OpenCL kernels or in HLS components. For information about HLS libraries, refer to *Intel High Level Synthesis Compiler: Reference Manual*.

**Figure 27. Overview of Intel FPGA SDK for OpenCL's Library Support**



You may use a previously-created library or create your own library. To use an OpenCL library, you do not require in-depth knowledge in hardware design or in the implementation of library primitives. To create an OpenCL library, you need to create the following files:

**Table 10. Necessary Files and Components for Creating an OpenCL Library**

File or Component	Description
<b>RTL Components</b>	
RTL source files	Verilog, System Verilog, or VHDL files that define the RTL component. Additional files such as Intel Quartus Prime IP File (.qip), Synopsys Design Constraints File (.sdc), and Tcl Script File (.tcl) are not allowed.
eXtensible Markup Language File (.xml)	Describes the properties of the RTL component. The Intel FPGA SDK for OpenCL Offline Compiler uses these properties to integrate the RTL component into the OpenCL pipeline.
Header file (.hcl)	A header file that contains valid OpenCL kernel language and declares the signatures of function(s) that are implement by the RTL component.
OpenCL emulation model file (.cl)	Provides C model for the RTL component that is used only for emulation. Full hardware compilations use the RTL source files.
<b>OpenCL Functions</b>	
OpenCL source files (.cl)	Contains definitions of the OpenCL functions. These functions are used during emulation and full hardware compilations.
Header file (.hcl)	A header file describing the functions to be called from OpenCL in the OpenCL kernel language syntax.
<b>HLS Functions</b>	
HLS source files (.cpp)	Contains definitions of the OpenCL functions. These functions are used during emulation and full hardware compilations.
Header file (.hcl)	A header file describing the functions to be called from OpenCL in the OpenCL kernel language syntax.

**Note:** There is no difference in the header file used for RTL, OpenCL, and HLS library functions. A single header file can have all types of functions declared. A single library can contain any of the supported sources. You can create a library from mixed sources (OpenCL, HLS, or RTL) and target these Intel high-level design products:

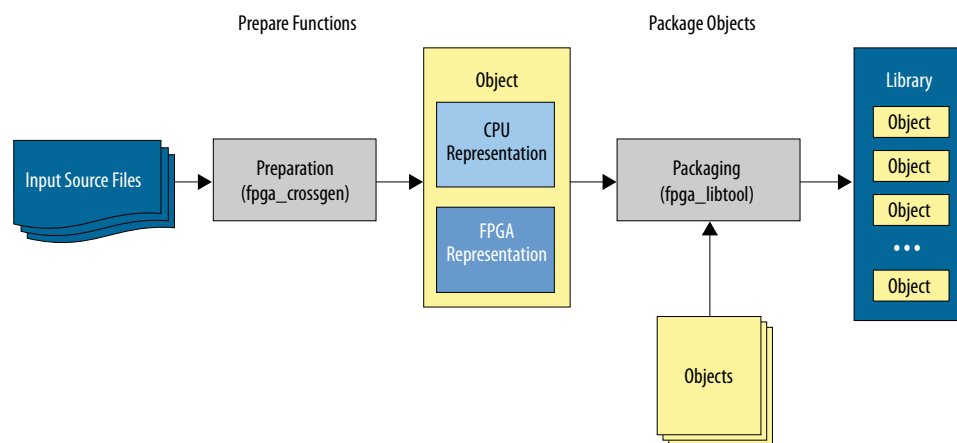
- Intel FPGA SDK for OpenCL Pro Edition
- Intel High Level Synthesis Compiler Pro Edition

Creating a OpenCL library is a two-step process:

1. Each object file is created from input source files using the `fpga_crossgen` command. The required input source files depend on the type of source code you are creating the object from.

- An object is effectively an intermediate representation of your source code with both a CPU representation and an FPGA representation of your code.
  - An object can be targeted for use with only one Intel high-level design product. If you want to target more than one high-level design product, you must generate a separate object for each target product.
2. Object files are combined into a library file using the `fpga_libtool` command.
- Objects created from different types of source code can be combined into a library, provided all objects target the same high-level design product.
- A library is assigned a version number, and can be used only with the targeted high-level design product with the same version number (for example, Intel FPGA SDK for OpenCL Pro Edition version 22.4).

**Figure 28. High-level View of the Library Creation Process**



[Creating Library Objects From OpenCL Code on page 172](#)

[Understanding RTL Modules and the OpenCL Pipeline on page 174](#)

[Packaging an OpenCL Helper Function File for an OpenCL Library on page 190](#)

[Packaging an RTL Component for an OpenCL Library on page 190](#)

[Verifying the RTL Modules on page 192](#)

[Specifying an OpenCL Library when Compiling an OpenCL Kernel on page 193](#)

[Debugging Your OpenCL Library Through Simulation \(Preview\) on page 194](#)

[Using an OpenCL Library that Works with Simple Functions \(Example 1\) on page 197](#)

[Using an OpenCL Library that Works with External Memory \(Example 2\) on page 197](#)

[OpenCL Library Command-Line Options on page 198](#)

### Related Information

[OpenCL Library Command-Line Options on page 198](#)

### 11.1.1.1. Creating Library Objects From OpenCL Code

You can create a library from object files from your OpenCL source code. An OpenCL-based object file includes code for CPU as well as hardware execution (CPU-capturing testbench and emulation use).

**Note:** Testbench is Intel HLS Compiler-specific. For information, refer to [Intel High Level Synthesis Compiler: Reference Manual](#).

A library can contain multiple object files. You can create object files for use in different Intel high-level design tools from the same OpenCL source code. Depending on the target high-level design tool, your source code might require adjustments to support tool-specific data types or constructs.

#### Intel FPGA SDK for OpenCL

No additional work is needed in your Intel FPGA SDK for OpenCL source code when you use the code to create objects for the offline compiler libraries.

**Restriction:** Creating object files from OpenCL code is supported only on Linux operating systems.

#### Intel HLS Compiler

The Intel FPGA SDK for OpenCL supports language constructs that are not natively supported by C++. Your component might need modifications to support those constructs (it is always preferred to allow OpenCL data types as library function call parameters).

The Intel HLS Compiler supports a limited set of OpenCL language constructs through the `ocl_types.h` header file. For details, review [Intel High Level Synthesis Compiler: Reference Manual](#).

**Restriction:** You cannot use systems of tasks in components intended for use in an OpenCL library object.

#### 11.1.1.1.1. Creating an Object File From OpenCL Code

Use the `fpga_crossgen` command to create library objects from your OpenCL code. An object created from OpenCL code contains information required both for emulating the functions in the object and synthesizing the hardware for the object functions.

**Note:** When the `aoc` command links non-library object files, it checks that all files are compiled with the same command line options to avoid conflicting directions to the compiler. The exception to this are the `-ffp-reassociate` and `-ffp-contract=fast` flags that must be passed at parse time (that is, when the `.cl` or `.cpp` file is passed into the `aoc` or the `fpga_crossgen` command). Different commands can have different settings, and the setting applied to a given source file are determined by the options provided when that source file was parsed.

The `fpga_crossgen` command creates one object file from one input source file. The object created can be used only in libraries that target the same Intel high-level design tool. Also, objects are versioned. That is, each object is assigned a compiler version number and be used only with Intel high-level design tools with the same version number.

Create a library object using the following command:

```
fpga_crosssgen <source_file> --target target_HLD_tool [-o <object_file>]
```

where, *target\_HLD\_tool* is the target Intel high-level design tool for this library. This parameter can have one of the following values:

- **aoc**  
Target this object to be included in libraries for kernels developed with the Intel FPGA SDK for OpenCL.  
Objects built for the Intel FPGA SDK for OpenCL are not operating system-specific. The objects are combined as Intel FPGA SDK for OpenCL object files (.aoco).
- **hls**  
Target this object to be included in libraries for components developed with the Intel HLS Compiler.  
Objects built for the Intel HLS Compiler are combined as operating system specific object files (.o on Linux). You cannot use objects created on one operating system with the Intel HLS Compiler running on a different operating system.

If you do not specify an object file name with the -o option, the object file name defaults to be the same name as the source code file name.

### 11.1.1.2. Packaging Object Files into a Library File

Gather the object files into a library file so that others can incorporate the library into their projects and call the functions that are contained in the objects in the library. To package object files into a library, use the `fpga_libtool` command.

Before you package object files into a library, ensure that you have the path information for all of the object files that you want to include in the library.

All objects that you want to package into a library must have the same version number (for example, Intel FPGA SDK for OpenCL Pro Edition version 19.3). The `fpga_libtool` command creates libraries encapsulated in operating system-specific archive files (.a on Linux and .lib on Windows). You cannot use libraries created on one operating system with an Intel high-level design product running on a different operating system.

Create the OpenCL library file using the following command:

```
fpga_libtool --target target_HLD_tool --create library_name object_file_1  
[object_file_2 ... object_file_n]
```

Where the command parameters are defined as follows:

**Table 11. Parameters of the `fpga_libtool` Command**

Parameter	Description
<i>target_HLD_tool</i>	The target Intel high-level design tool for this library. This parameter can have one of the following values:

*continued...*

Parameter	Description
	<ul style="list-style-type: none"> <li>• <code>aoc</code> Target this library for kernels developed with the Intel FPGA SDK for OpenCL. Libraries built for the Intel FPGA SDK for OpenCL are not operating system-specific. The objects are created as Intel FPGA SDK for OpenCL object files (<code>.aoclib</code>).</li> <li>• <code>hls</code> Target this library for components developed with the Intel HLS Compiler. Libraries built for the Intel HLS Compiler are encapsulated in operating system-specific archive files (<code>.a</code> on Linux and <code>.lib</code> on Windows).</li> </ul>
<code>library_name</code>	<p>The name of the library file.</p> <p>Specify the file extension of the library files as follows, depending on the target high-level design tool:</p> <ul style="list-style-type: none"> <li>• Intel FPGA SDK for OpenCL Specify <code>.aoclib</code> as the file extension for an OpenCL library. OpenCL libraries are not operating-system specific.</li> <li>• Intel HLS Compiler Specify the operating-system specific archive extension:               <ul style="list-style-type: none"> <li>– <code>.a</code> for Linux-platform libraries.</li> <li>– <code>.lib</code> for Windows-platform libraries.</li> </ul> </li> </ul>

You can specify one or more object files to include in the library.

For example, the following command packages three Linux-platform objects (`prim1.o`, `prim2.o`, and `prim3.o`) into an OpenCL library called `libdemo`:

```
fpga_libtool --create libdemo.a prim1.o prim2.o prim3.o --target aoc
```

### 11.1.2. Understanding RTL Modules and the OpenCL Pipeline

The OpenCL library feature allows you to use RTL modules, written in Verilog, SystemVerilog, or VHDL, inside OpenCL kernels. This section provides an overview of how the Intel FPGA SDK for OpenCL Offline Compiler integrates RTL modules into the Intel FPGA SDK for OpenCL pipeline architecture.

Use RTL modules under the following circumstances:

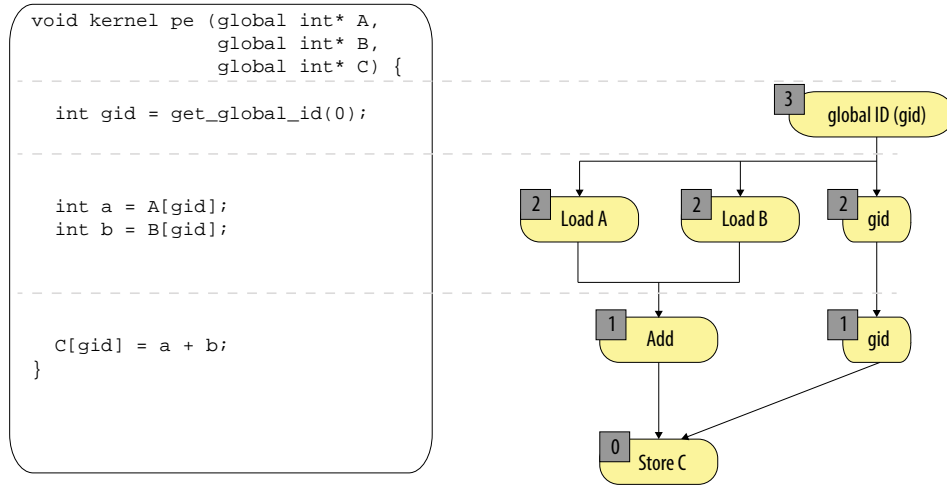
- You want to use optimized and verified RTL modules in OpenCL kernels without rewriting the modules as OpenCL functions.
- You want to implement OpenCL kernel functionality that you cannot express effectively in OpenCL.

#### 11.1.2.1. Overview: Intel FPGA SDK for OpenCL Pipeline Approach

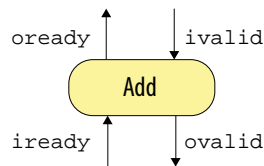
The following figure depicts the architecture of an Intel FPGA SDK for OpenCL pipeline:

**Figure 29. Parallel Execution Model of Intel FPGA SDK for OpenCL Pipeline Stages**

The operations on the right represent the SDK's pipeline implementation of the OpenCL kernel code on the left. Each yellow box is an operation or data value found in the pipeline. The number associated with each operation represents the number of threads in the pipeline.



Assume each level of operation is one stage in the pipeline. At each stage, the Intel FPGA SDK for OpenCL Offline Compiler executes all operations in parallel by the thread existing at that stage. For example, thread 2 executes Load A, Load B, and copies the current global ID (via `gid`) to the next pipeline stage. Similar to the pipelined execution on instructions in reduced instruction set computing (RISC) processors, the SDK's pipeline stages also execute in parallel. The threads advances to the next pipeline stage only after all the stages have completed execution.



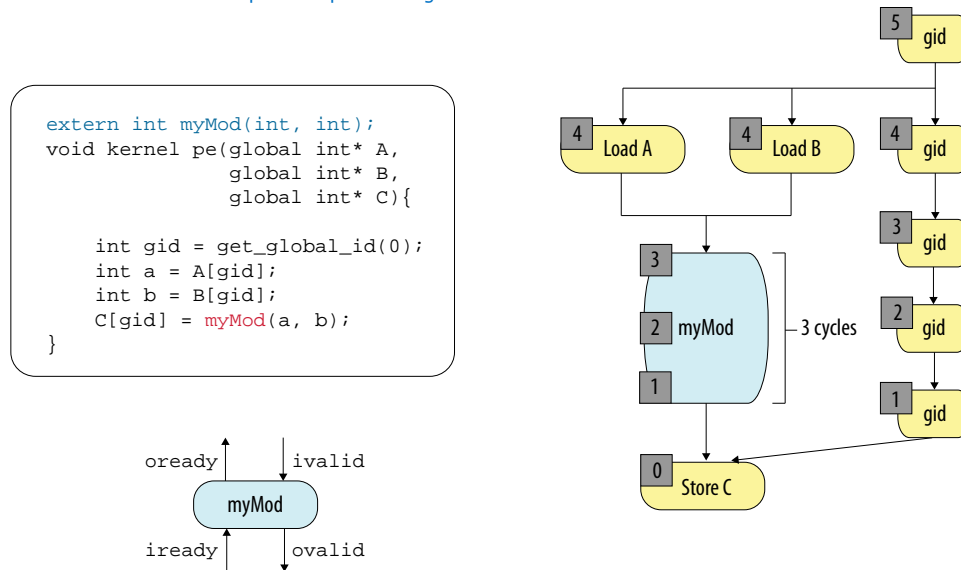
Some operations are capable of stalling the Intel FPGA SDK for OpenCL pipeline. Examples of such operations include variable latency operations like memory load and store operations. To support stalls, ready and valid signals need to propagate throughout the pipeline so that the offline compiler can schedule the pipeline stages. However, ready signals are not necessary if all operations have fixed latency. In these cases, the offline compiler optimizes the pipeline to statically schedule the operations, which significantly reduces the logic necessary for pipeline implementation.

### 11.1.2.2. Integration of an RTL Module into the Intel FPGA SDK for OpenCL Pipeline

When you specify an OpenCL library during kernel compilation, the offline compiler integrates the RTL module within the library into the overall pipeline.

**Figure 30. Integration of an RTL Module into an Intel FPGA SDK for OpenCL Pipeline**

This figure depicts the integration of the RTL module `myMod` into the pipeline depicted in [Parallel Execution Model of Intel FPGA SDK for OpenCL Pipeline Stages](#).



The depicted RTL module has a balanced latency where the threads of the RTL module match the number of pipeline stages. A balanced latency allows the threads of the RTL module to execute without stalling the SDK's pipeline.

Setting the latency of the RTL module in the RTL specification file allows the offline compiler to balance the pipeline latency. RTL supports Avalon® streaming interfaces; therefore, the latency of the RTL module can be variable (that is, not fixed). However, the variability in the latency should be small in order to maximize performance. In addition, specify the latency in the OpenCL library object manifest file so that the RTL module experiences a good approximation of the actual latency in steady state.

### Related Information

- [Avalon Interface Specifications](#)
- [Pipelined Read Transfer with Variable Latency](#)
- [Pipelined Read Transfers with Fixed Latency](#)
- [Avalon Streaming Interface](#) on page 178
- [Object Manifest File Syntax of an RTL Module](#) on page 183

### 11.1.2.3. Stall-Free RTL

The Intel FPGA SDK for OpenCL Offline Compiler can optimize hardware resource usage and performance by removing stall logic around an RTL module with fixed latency.

An RTL module that has a variable latency and uses Avalon streaming interface input and output signals can wait until input data is ready. Conversely, the Intel FPGA SDK for OpenCL pipeline can stall until it receives valid output data from the RTL module. For an RTL module with a fixed latency, you can remove an RTL stall by modifying the object manifest file (`.xml`), as described below.



1. To instruct the offline compiler to remove stall logic around the RTL module, if appropriate, set the `IS_STALL_FREE` attribute under the `FUNCTION` element to "yes".

This modification informs the offline compiler that the RTL module produces valid data every `EXPECTED_LATENCY` cycle(s).

**Note:** `EXPECTED_LATENCY` is an attribute you specify in the `.xml` file under the `FUNCTION` element.

2. Specify a value for `EXPECTED_LATENCY` such that the latency equals the number of pipeline stages in the module.

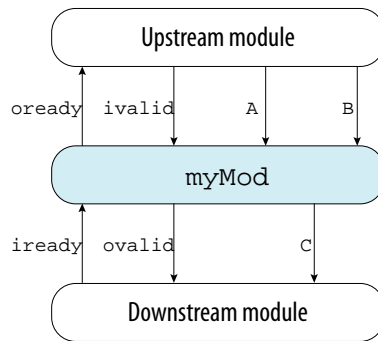
**Caution:** An inaccurate `EXPECTED_LATENCY` value causes the RTL module to be out of sync with the rest of the pipeline.

A stall-free RTL module might receive an invalid input signal (that is, `ivalid` is low). In this case, the module ignores the input and produces invalid data on the output. For a stall-free RTL module without an internal state, it might be easier to propagate the invalid input through the module. However, for an RTL module with an internal state, you must handle an `ivalid=0` input carefully.

#### 11.1.2.4. RTL Module Interfaces

For an RTL module to properly interact with other compiler-generated operations, you must support a simplified Avalon streaming interface at both input and output of an RTL module.

The following diagram shows the complete interface of the `myMod` RTL module shown in [Integration of an RTL Module into an Intel FPGA SDK for OpenCL Pipeline](#).



In this diagram, `myMod` interacts with the upstream module through data signals, `A` and `B`, and control signals, `ivalid` (input) and `oready` (output). The `ivalid` control signal equals 1 (`ivalid = 1`) if and only if data signal `A` and data signal `B` contain valid data. When the control signal `oready` equals 1 (`oready = 1`), it indicates that the `myMod` RTL module can process the data signals `A` and `B` if they are valid (that is, `ivalid = 1`). When `ivalid = 1` and `oready = 0`, the upstream module is expected to hold the values of `ivalid`, `A`, and `B` in the next clock cycle.

`myMod` interacts with the downstream module through data signal `C` and control signals, `ovalid` (output) and `iredy` (input). The `ovalid` control signal equals 1 (`ovalid = 1`) if and only if data signal `C` contains valid data. The `iredy` control signal equals 1 (`ivalid = 1`) indicates that the downstream module is able to

process data signal `C` if it is valid. When `ovalid = 1` and `iready = 0`, the `myMod` RTL module is expected to hold the valid of the `ovalid` and `C` signals in the next clock cycle.

`myMod` module asserts `oready` for a single clock cycle to indicate it is ready for an active cycle. Cycles during which `myMod` module is ready for data are called ready cycles. During ready cycles, the module above `myMod` module can assert `ivalid` to send data to `myMod`.

For a detailed explanation of data transfer under backpressure, refer to “[Data Transfer with Backpressure](#)” in *Avalon Interface Specifications*. Ignore the information pertaining to `readyLatency` option.

### Related Information

[Data Transfer with Backpressure](#)

#### 11.1.2.5. Avalon Streaming Interface

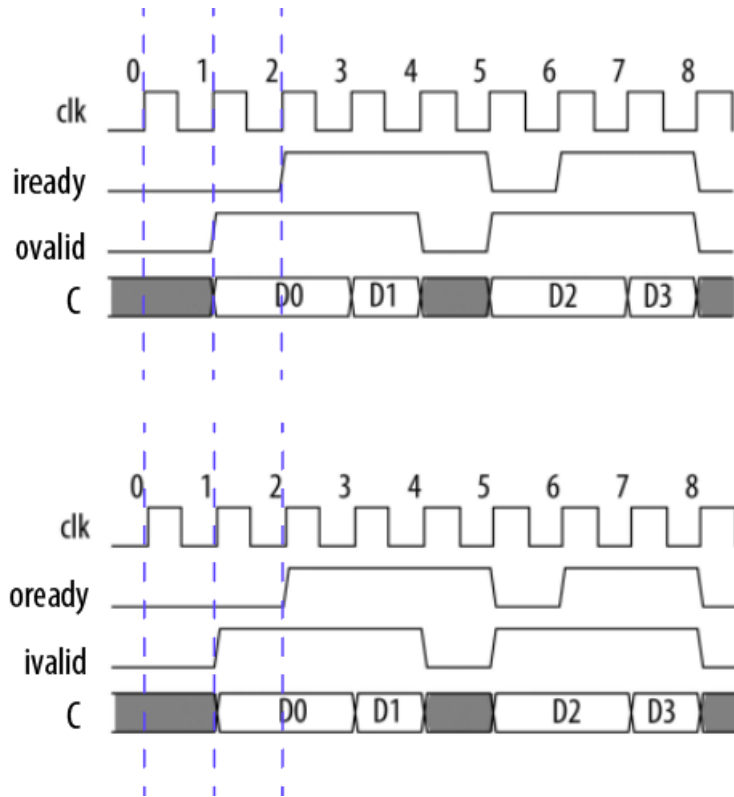
The offline compiler expects the RTL module to support Avalon streaming interface with `readyLatency = 0`, at both input and output.

As shown in [Integration of an RTL Module into an Intel FPGA SDK for OpenCL Pipeline](#), the RTL module must have four ports:

- `ivalid` and `iready` as the input Avalon streaming interface
- `ovalid` and `oready` as the output Avalon streaming interface

The following figure illustrates the timing diagram for input data transfer with back pressure. For more information about Avalon streaming interfaces, refer to the “[Avalon Streaming Interfaces](#)” section in *Avalon Interface Specifications*.

Figure 31. Timing Diagram for Input Data Transfer with Back Pressure



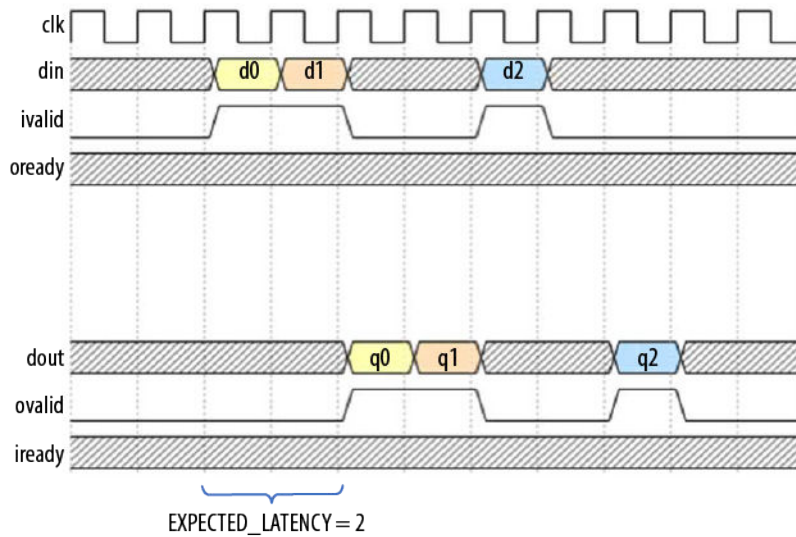
For an RTL module with a fixed latency, the output signals (*ovalid* and *oready*) can have constant high values, and the input ready signal (*iready*) can be ignored.

A stall-free RTL module might receive an invalid input signal (*ivalid* is low). In this case, the module ignores the input and produces invalid data on the output. For a stall-free RTL module without an internal state, it might be easier to propagate the invalid input through the module. However, for an RTL module with an internal state, you must handle an *ivalid* = 0 input carefully.

### Example Timing Diagram of a Stall-free RTL Component

Consider the following example timing diagram of a stall-free RTL component:

**Figure 32. Timing Diagram of a Stall-free RTL Component**



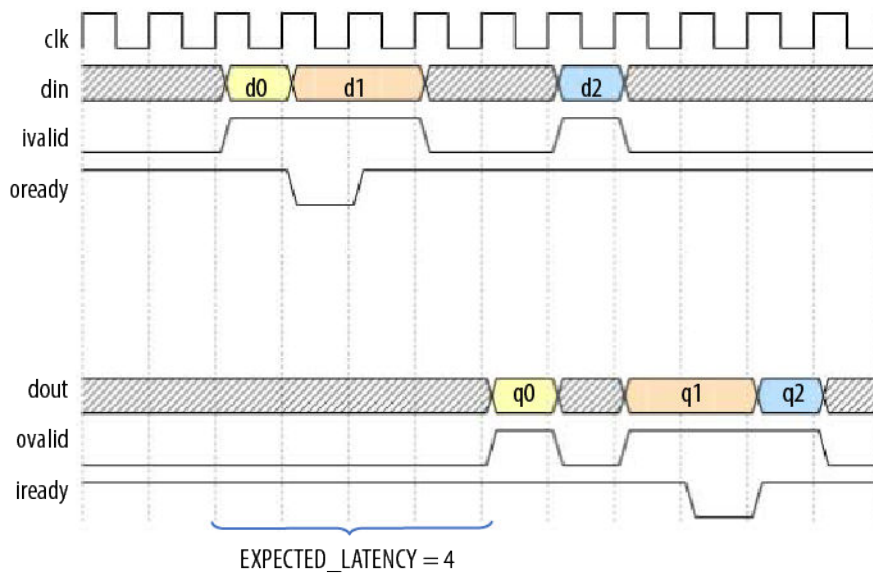
where,

- IS\_STALL\_FREE value = "yes"
- IS\_FIXED\_LATENCY value = "yes"
- EXPECTED\_LATENCY value = "2"

### Example Timing Diagram of a Non-stall-free RTL Component

Consider the following example timing diagram of a stallable RTL component:

**Figure 33. Timing Diagram of a Stallable RTL Component**



where,

- IS\_STALL\_FREE value = "no"
- IS\_FIXED\_LATENCY value = "no"
- EXPECTED\_LATENCY value = "4"

### Performing Advanced Compiler Optimizations

Both ALLOW\_MERGING and HAS\_SIDE\_EFFECTS parameters allow the offline compiler to perform advanced optimizations. Consider the following combinations to understand their impact completely:

*Note:* The combination you select depends on your design architecture. Allow the compiler to replicate the RTL block for multiple calls or vectorized code.

**Table 12.** ALLOW\_MERGING and HAS\_SIDE\_EFFECTS Parameters Combination and Their Impact

Combination	Description
ALLOW_MERGING value = "no" HAS_SIDE_EFFECTS value = "no"	Each call to an RTL library corresponds to one distinct instance in the hardware. Calls might be optimized away by the compiler if deemed redundant or unnecessary. Calls might be vectorized, with multiple instances in the hardware created for a single RTL library call.
ALLOW_MERGING value = "no" HAS_SIDE_EFFECTS value = "yes"	Each call to an RTL library corresponds to one distinct instance in hardware. Calls are not optimized away by the compiler. The compiler errors out if the attribute num_simd_work_items is greater than 1 for the kernel calling the RTL library.
ALLOW_MERGING value = "yes" HAS_SIDE_EFFECTS value = "no"	Multiple calls to an RTL library might be merged into one call, and hence correspond to one instance in the hardware. Calls might be optimized away by the compiler if deemed redundant or unnecessary. Calls might be vectorized, with multiple instances in the hardware created for a single RTL library call.
ALLOW_MERGING value = "yes" HAS_SIDE_EFFECTS value = "yes"	Multiple calls to an RTL library might be merged into one call, and hence correspond to one instance in hardware. Calls are not optimized away by the compiler. The compiler errors out if the attribute num_simd_work_items is greater than 1 for the kernel calling the RTL library.

#### 11.1.2.6. RTL Reset and Clock Signals

Resets and clocks of RTL modules are connected to the same clock and reset drivers as the rest of the OpenCL pipeline.

Because of the common clock and reset drivers, an RTL module runs in the same clock domain as the OpenCL kernel. The module is reset only when the OpenCL kernel is first loaded onto the FPGA, either via Intel FPGA SDK for OpenCL program utility or the clCreateProgramwithBinary host function. In particular, if the host restarts a kernel via successive clEnqueueNDRangeKernel or clEnqueueTask invocations, the associated RTL modules does not reset in between these restarts.

The following steps outline the process of setting the kernel clock frequency:

1. The Intel Quartus Prime software's Fitter applies an aggressive constraint on the kernel clock.
2. The Intel Quartus Prime software's Timing Analyzer performs static timing analysis to determine the frequency that the Fitter actually achieves.
3. The phase-locked loop (PLL) that drives the kernel clock sets the frequency determined in Step 2 to be the kernel clock frequency.

Optionally, an RTL module can access a system-wide clock that runs at twice the frequency of the OpenCL kernel clock. This system-wide clock can be connected to an input signal of the RTL module by including an `AVALON` element of type `clock2x`. The phase relationship between the `clock` and `clock2x` signals is such that the rising and falling edges of clock are coincident with rising edges of `clock2x`.

**Warning:** Timing failures may occur if one or more signals in your design are not able to satisfy all of the timing requirements of the device. All timing small timing failure can cause problems, so binaries that failed timing should not be used for development and production builds.

If your design failed timing, you have the following options:

- Timing failures can depend on how a design is placed on the FPGA, so running a sweep of different seeds (which results in different component placements) might lead to a passing binary.
- Decreasing the size of the design makes the component placement easier and timing failures less likely.
- Timing failures may be indicative of BSP problems, so if you are using a custom BSP, discuss with your BSP vendor. If you want to investigate it further, the Intel Quartus Prime Static Timing Analyzer outputs a `*.sta.rpt` file that contains more details about the timing analysis performed.

#### Related Information

- [Intel Quartus Prime Pro Edition User Guide: Timing Analyzer](#)
- [Design Evaluation for Timing Closure](#)

#### 11.1.2.6.1. Intel Stratix 10 Design-Specific Reset Requirements for Stall-Free and Stallable RTL Modules

When creating an RTL module for Intel Stratix 10 OpenCL designs, ensure that the module satisfies specific logic reset requirements.

##### Reset Requirements for Stall-Free RTL Modules

A stall-free RTL module is a fixed-latency module for which the Intel FPGA SDK for OpenCL Offline Compiler can optimize away stall logic.

- When creating a stall-free RTL module for a Intel Stratix 10 design, use synchronous clear signals only.
- After **deassertion** of the reset signal to the stall-free RTL module, the module must be operational within 15 clock cycles. If the reset signal is pipelined within the module, this requirement limits the reset pipelining to no more than 15 stages.

## Reset Requirements for Stallable RTL Modules

A stallable RTL module has a variable latency, and it relies on backpressured input and output interfaces to function correctly.

- When creating a stallable RTL module for a Intel Stratix 10 design, use synchronous clear signals only.
- After **assertion** of the reset signal to the stallable RTL module, the module must deassert its `oready` and `ovalid` interface signals within 40 clock cycles.
- After **deassertion** of the reset signal to the stallable RTL module, the module must be fully operational within 40 clock cycles. The module signals its readiness by asserting the `oready` interface signal.

## Related Information

[Stall-Free RTL](#) on page 176

### 11.1.2.7. Object Manifest File Syntax of an RTL Module

This section provides the syntax of a simple object manifest file for an RTL module that implements double-precision square root function. The RTL module is implemented in VHDL with a Verilog wrapper.

The following object manifest file is for an RTL module named `my_fp_sqrt_double` (line 2) that implements an OpenCL helper function named `my_sqrtfd` (line 2).

```
1: <RTL_SPEC>
2:   <FUNCTION name="my_sqrtfd" module="my_fp_sqrt_double">
3:     <ATTRIBUTES>
4:       <IS_STALL_FREE value="yes"/>
5:       <IS_FIXED_LATENCY value="yes"/>
6:       <EXPECTED_LATENCY value="31"/>
7:       <CAPACITY value="1"/>
8:       <HAS_SIDE_EFFECTS value="no"/>
9:       <ALLOW_MERGING value="yes"/>
10:    </ATTRIBUTES>
11:    <INTERFACE>
12:      <AVALON port="clock" type="clock"/>
13:      <AVALON port="resetn" type="resetn"/>
14:      <AVALON port="ivalid" type="ivalid"/>
15:      <AVALON port="iready" type="iready"/>
16:      <AVALON port="ovalid" type="ovalid"/>
17:      <AVALON port="oready" type="oready"/>
18:      <INPUT port="datain" width="64"/>
19:      <OUTPUT port="dataout" width="64"/>
20:    </INTERFACE>
21:    <C_MODEL>
22:      <FILE name="c_model.cl" />
23:    </C_MODEL>
24:    <REQUIREMENTS>
25:      <FILE name="my_fp_sqrt_double_s5.v" />
26:      <FILE name="fp_sqrt_double_s5.vhd" />
27:    </REQUIREMENTS>
28:    <RESOURCES>
29:      <ALUTS value="2057"/>
30:      <FFS value="3098"/>
31:      <RAMS value="15"/>
32:      <MLABS value="43"/>
33:      <DSPS value="1.5"/>
34:    </RESOURCES>
35:  </FUNCTION>
36: </RTL_SPEC>
```

**Table 13. Elements and Attributes in the Object Manifest File**

XML Element	Description
RTL_SPEC	Top-level element in the object manifest file. There can only be one such top-level element in the file. In this example, the name RTL_SPEC is historic and carries no file-specific meaning.
FUNCTION	Element that defines the OpenCL function that the RTL module implements. The name attribute within the FUNCTION element specifies the function's name. You may have multiple FUNCTION elements, each declaring a different function that you can call from the OpenCL kernel. The same RTL module can implement multiple functions by specifying different parameters.
ATTRIBUTES	Element containing other XML elements that describe various characteristics (for example, latency) of the RTL module. The example RTL module takes one PARAMETER setting named WIDTH, which has a value of 32. Refer to <a href="#">XML Elements Associated with the ATTRIBUTES Element in the Object Manifest File of an RTL Module</a> for more details other ATTRIBUTES-specific elements. <i>Note:</i> If you create multiple OpenCL helper functions for different modules, or use the same RTL module with different PARAMETER settings, you must create a separate FUNCTION element for each function.
INTERFACE	Element containing other XML elements that describe the RTL module's interface. The example object manifest file shows the Avalon streaming interface signals that every RTL module must provide (that is, clock, resetn, ivalid, iready, ovalid, and oready). The resetn signal is active low. Its synchronicity depends on the target device:  <i>Intel Arria 10</i> The resetn signal is asynchronous to the clock signal.  <i>Intel Stratix 10</i> The resetn signal is synchronous to the clock signal. For more information about reset signal timing, see <a href="#">Intel Stratix 10 Design-Specific Reset Requirements for Stall-Free and Stallable RTL Modules</a> on page 182.  The signal names must match the ones specified in the .xml file. An error occurs during library creation if a signal name is inconsistent.
C_MODEL	Element specifying one or more files that implement OpenCL C model for the function. The model is used only during emulation. However, the C_MODEL element and the associated file(s) must be present when you create the library file.
REQUIREMENTS	Element specifying one or more RTL resource files (that is, .v, .sv, .vhd, .hex, and .mif). The specified paths to these files are relative to the location of the object manifest file. Each RTL resource file becomes part of the associated Platform Designer component that corresponds to the entire OpenCL system. <i>Note:</i> The Intel FPGA SDK for OpenCL library feature does not support .qip files. An Intel FPGA SDK for OpenCL Offline Compiler error occurs if you compile an OpenCL kernel while using a library that includes an unsupported resource file type.
RESOURCES	Optional element specifying the FPGA resources that the RTL module uses. If you do not specify this element, the FPGA resources that the RTL module uses defaults to zero.

**11.1.2.7.1. XML Elements for ATTRIBUTES**

In the object manifest file of the RTL module within an Intel FPGA SDK for OpenCL library, there are XML elements under ATTRIBUTES that you can specify to set the module's characteristics.



**Table 14. XML Elements Associated with the ATTRIBUTES Element in the Object Manifest File of an RTL Module**

**Attention:** Except for IS\_STALL\_FREE and EXPECTED\_LATENCY, all elements have safe values. If you are unsure which value you should specify for an attribute, set it to the safe value. Compiling your kernel with a library that uses safe values results in functional hardware. However, the hardware might be larger than the actual size.

XML Element	Description
IS_STALL_FREE	<p>Instructs the Intel FPGA SDK for OpenCL Offline Compiler to remove all stall logic around the RTL module.</p> <p>Set IS_STALL_FREE to "yes" to indicate that the module neither generates stalls internally nor can it properly handle incoming stalls. The module simply ignores its stall input. If you set IS_STALL_FREE to "no", the module must properly handle all stall and valid signals.</p> <p><i>Note:</i> If you set IS_STALL_FREE to "yes", you must also set IS_FIXED_LATENCY to "yes". Also, if the RTL module has an internal state, it must properly handle invalid=0 inputs.</p> <p>An incorrect IS_STALL_FREE setting leads to incorrect results in hardware.</p>
IS_FIXED_LATENCY	<p>Indicates whether the RTL module has a fixed latency.</p> <p>Set IS_FIXED_LATENCY to "yes" if the RTL module always takes a known number of clock cycles to compute its output. The value you assign to the EXPECTED_LATENCY element specifies the number of clock cycles.</p> <p>The safe value for IS_FIXED_LATENCY is "no". When you set IS_FIXED_LATENCY="no", the EXPECTED_LATENCY value must be at least 1.</p> <p><i>Note:</i> For a given module, you may set IS_FIXED_LATENCY to "yes" and IS_STALL_FREE to "no". Such a module produces its output in a fixed number of clock cycles and handles stall signals properly.</p>
EXPECTED_LATENCY	<p>Specifies the expected latency of the RTL module.</p> <p>If you set IS_FIXED_LATENCY to "yes", the EXPECTED_LATENCY value indicates the number of pipeline stages inside the module. In this case, you must set this value to be the exact latency of the module. Otherwise, the offline compiler generates incorrect hardware.</p> <p>For a module with variable latency, the offline compiler balances the pipeline around this module to the EXPECTED_LATENCY value that you specify. For modules that can stall and require use of signals such as iready, the EXPECTED_LATENCY value must be set to at least 1. The specified value and the actual latency might differ, which might affect the number of stalls inside the pipeline. However, the resulting hardware is correct.</p>
CAPACITY	<p>Specifies the number of multiple inputs that this module can process simultaneously. You must specify a value for CAPACITY if you also set IS_STALL_FREE="no" and IS_FIXED_LATENCY="no". Otherwise, you do not need to specify a value for CAPACITY.</p> <p>If CAPACITY is strictly less than EXPECTED_LATENCY, the offline compiler automatically inserts capacity-balancing FIFO buffers after this module when necessary.</p> <p>The safe value for CAPACITY is 1.</p>
HAS_SIDE_EFFECTS	<p>Indicates whether the RTL module has side effects. Modules that have internal states or communicate with external memories are examples of modules with side effects.</p> <p>Set HAS_SIDE_EFFECTS to "yes" to indicate that the module has side effects. Specifying HAS_SIDE_EFFECTS to "yes" ensures that optimization efforts do not remove calls to modules with side effects.</p> <p>Stall-free modules with side effects (that is, IS_STALL_FREE="yes" and HAS_SIDE_EFFECTS="yes") must properly handle invalid=0 input cases because the module might receive invalid data occasionally.</p> <p>The safe value for HAS_SIDE_EFFECTS is "yes".</p>
ALLOW_MERGING	<p>Instructs the offline compiler to merge multiple instances of the RTL module.</p> <p>Set ALLOW_MERGING to "yes" to allow merging of multiple instances of the module. Intel recommends setting ALLOW_MERGING to "yes".</p> <p>The safe value for ALLOW_MERGING is "no".</p>

*continued...*

XML Element	Description
	<i>Note:</i> Marking the module with HAS_SIDE_EFFECTS="yes" does not prevent merging.
PARAMETER	<p>Specifies the value of an RTL module parameter.</p> <p>PARAMETER attributes:</p> <ul style="list-style-type: none"> <li>• name—Specifies the name of the RTL module parameter.</li> <li>• value—Specifies a decimal numeric value for the parameter.</li> <li>• type—Specifies a system parameter the value of which is used as the RTL module parameter value. You can use the following system parameter name with the type attribute: <ul style="list-style-type: none"> <li>– bspaddresswidth—Specifies the Avalon memory bus width required to address the memory range configured for OpenCL global memory in the board support package.</li> </ul> </li> </ul> <p><i>Note:</i> The value for an RTL module parameter can be specified using either a value or a type attribute.</p>

### 11.1.2.7.2. XML Elements for INTERFACE

In the object manifest file of the RTL module within an Intel FPGA SDK for OpenCL library, there are XML elements under INTERFACE that you can define to specify aspects of the RTL module's interface (for example, Avalon streaming interface).

**Table 15. Mandatory XML Elements Associated with the INTERFACE Element in the Object Manifest File of an RTL Module**

XML Element	Description
INPUT	<p>Specifies the input parameter of the RTL module.</p> <p>INPUT attributes:</p> <ul style="list-style-type: none"> <li>• port—Specifies the port name of the RTL module.</li> <li>• width—Specifies the width of the port in bits.</li> </ul> <p>The input parameters are concatenated to form the input stream. Aggregate data structures such as structs and arrays are not supported as input parameters.</p>
OUTPUT	<p>Specifies the output parameter of the RTL module.</p> <p>OUTPUT attributes:</p> <ul style="list-style-type: none"> <li>• port—Specifies the port name of the RTL module.</li> <li>• width—Specifies the width of the port in bits.</li> </ul> <p>The return value from the input stream is sent out via the output parameter on the output stream. Aggregate data structures such as structs and arrays are not supported as input parameters.</p>

**Table 16. Additional XML Elements to Support External Memory Access**

XML Element	Description
MEM_INPUT	<p>Describes a pointer input to the RTL module.</p> <p>MEM_INPUT attributes:</p> <ul style="list-style-type: none"> <li>• port—Specifies the name of the pointer input.</li> <li>• access—Specifies to the Intel FPGA SDK for OpenCL Offline Compiler how the RTL module uses this pointer. Valid access values are <code>readonly</code> and <code>readwrite</code>. If the RTL module only writes with this pointer, assign <code>readwrite</code> to access.</li> </ul> <p>Because all pointers to external memory must be 64 bits, there is no width attribute associated with MEM_INPUT.</p>
AVALON_MEM	Declares the Avalon memory-mapped interface for your RTL module.
<i>continued...</i>	

XML Element	Description
	<p>AVALON_MEM attributes:</p> <ul style="list-style-type: none"> <li>• <code>port</code>—Specifies the root of the corresponding port names in the RTL module. For example, if <code>port</code> has a value of <code>avm_port0_</code>, the names of all Avalon memory-mapped interface ports for the RTL module starts with <code>avm_port0_</code>.</li> <li>• <code>width</code>—Specifies the data width, which must match the corresponding width value in the accelerator board's <code>board_spec.xml</code> file. Within the <code>board_spec.xml</code> file, the width value is specified in the interface element under <code>global_mem</code>. For more information, refer to the <code>global_mem</code> section under <i>XML Elements, Attributes, and Parameters in the board_spec.xml File</i> in the <i>Intel FPGA SDK for OpenCL Custom Platform Toolkit User Guide</i>.</li> <li>• <code>burstwidth</code>—Specifies the number of bits required to represent burst size. Use <math>\text{burstwidth} = \log(\text{maxburst}) + 1</math> to calculate the burst size, where <code>maxburst</code> is the corresponding maximum burst size specified in the <code>board_spec.xml</code> file. For example, if <code>maxburst=16</code>, <code>burstwidth=5</code>.</li> <li>• <code>optype</code>—Specifies either the Avalon memory-mapped port is reading (<code>read</code>) or writing (<code>write</code>) from external memory. You can only assign either <code>read</code> or <code>write</code> to <code>optype</code>.</li> </ul>

For the AVALON\_MEM element defined in the code example above, the corresponding RTL module ports are as follows:

```
output      avm_port0_enable,
input [511:0] avm_port0_readdata,
input      avm_port0_readdatavalid,
input      avm_port0_waitrequest,
output [31:0] avm_port0_address,
output      avm_port0_read,
output      avm_port0_write,
input      avm_port0_writeack,
output [511:0] avm_port0_writedata,
output [63:0] avm_port0_byteenable,
output [4:0] avm_port0_burstcount,
```

There is no assumed correspondence between pointers that you specify with MEM\_INPUT and the Avalon memory-mapped interfaces that you specify with AVALON\_MEM. An RTL module can use a single pointer to address zero to multiple Avalon memory-mapped interfaces.

**Related Information**

[XML Elements, Attributes, and Parameters in the board\\_spec.xml File: global\\_mem](#)

**11.1.2.7.3. XML Elements for RESOURCES**

In the object manifest file of the RTL module within an Intel FPGA SDK for OpenCL library, there are optional elements under RESOURCES that you can define to specify the FPGA resource utilization of the module. If you do not specify a particular element, it has a default value of zero.

**Table 17. XML Elements Associated with the RESOURCES Element in the Object Manifest File of an RTL Module**

XML Element	Description
ALUTS	Specifies the number of combinational adaptive look-up tables (ALUTs) that the module uses.
FFS	Specifies the number of dedicated logic registers that the module uses.
RAMS	Specifies the number of block RAMs that the module uses.
DSPS	Specifies the number of digital signal processing (DSP) blocks that the module uses.
MLABS	Specifies the number of memory logic arrays (MLABs) that the module uses. This value is equal to the number of adaptive logic modules (ALMs) that is used for memory divided by 10 because each MLAB consumes 10 ALMs.

### 11.1.2.8. Interaction between RTL Module and External Memory

Allow your RTL module to interact with external memory only if the interaction is necessary and unavoidable.

*Important:* The preferred method for having your RTL module interact with external memory is to have the OpenCL kernel access global memory and then feed that memory content to the RTL module. For operations like reading from and writing to external memory on every kernel invocation, instruct the OpenCL kernel to perform the operation. To do so, you can create an OpenCL helper function for the OpenCL kernel in the same Intel FPGA SDK for OpenCL library as the RTL module.

The following examples demonstrate how to structure code in an RTL module for easy integration into an OpenCL library:

**Table 18. Example Code in an RTL Module that Interacts with External Memory**

Complex RTL Module	Simplified RTL Module
<pre>// my_rtl_fn does: // out_ptr[idx] = fn(in_ptr[idx]) my_rtl_fn (in_ptr, out_ptr,idx);</pre>	<pre>int in_value = in_ptr[idx]; // my_rtl_fn now does: out = fn(in) int out_value = my_rtl_fn (in_value); out_ptr[idx] = out_value;</pre>

The complex RTL module on the left reads a value from external memory, performs a scalar function on the value, and then writes the value back to global memory. Such an RTL module is difficult to describe when you integrate it into an OpenCL library. In addition, this RTL module is harder to verify and causes very conservative pointer analysis in the Intel FPGA SDK for OpenCL Offline Compiler.

The simplified RTL module on the right provides the same overall functionality as the complex RTL module. However, the simplified RTL module only performs a scalar-to-scalar calculation without connecting to global memory. Integrating this simplified RTL module into the OpenCL library makes it much easier for the offline compiler to analyze the resulting OpenCL kernel.

There are times when an RTL module requires an Avalon memory-mapped interface port to communicate with external memory. This Avalon memory-mapped interface port connects to the same arbitration network to which all other global load and store units in the OpenCL kernels connect.

If an RTL module receives a memory pointer as an argument, the offline compiler enforces the following memory model:

- If an RTL module writes to a pointer, nothing else in the OpenCL kernel can read from or write to this pointer.
- If an RTL module reads from a pointer, the rest of the OpenCL kernel and other RTL modules may also read from this pointer.
- You may set the `access` field of the `MEM_INPUT` attribute to specify how the RTL module uses the memory pointer. Ensure that you set the value for `access` correctly because there is no way to verify the value.

#### 11.1.2.9. Order of Threads Entering an RTL Module

Do not assume that threads entering an RTL module follow a defined order. As a result, thread 0 does not necessarily enter the module before thread 1.

#### 11.1.2.10. OpenCL C Model of an RTL Module

To run in emulation, each RTL module within an OpenCL library must have an OpenCL C model. The OpenCL C model verifies the overall OpenCL system during emulation.

Example OpenCL C model file for a square root function:

```
double my_sqrtfd (double a)
{
    return sqrt(a);
}
```

Intel recommends that you emulate your OpenCL system. If you decide not to emulate your OpenCL system, no C model is required.

#### Related Information

[Object Manifest File Syntax of an RTL Module](#) on page 183

#### 11.1.2.11. Potential Incompatibility between RTL Modules and Partial Reconfiguration

When creating an OpenCL library using RTL modules, you might encounter Partial Reconfiguration-related issues.

Consider a situation where you create and verify your library on a device that does not support Partial Reconfiguration (PR). If a library user then uses the library's RTL module inside a PR region, the module might not function correctly after PR.

To ensure that the RTL modules function correctly on a device that uses PR:

- The RTL modules do not use memory logic array blocks (MLABs) with initialized content.
- The RTL modules do not make any assumptions regarding the power-up values of any logic.

For complete PR coding guidelines, refer to [Creating a Partial Reconfiguration Design](#) in the *Partial Reconfiguration User Guide*.

### 11.1.3. Packaging an OpenCL Helper Function File for an OpenCL Library

Before creating an OpenCL library file, package each OpenCL source file with helper functions into a `.aoco` file. Unlike RTL modules, you do not need to create an object manifest file.

In general, you do not need to create a library to share helper functions written in OpenCL. You can distribute a helper function in source form (for example, `<shared_file>.cl`) and then insert the line `#include "<shared_file>.cl"` in the OpenCL kernel source code.

Consider creating a library under the following circumstances:

- The helper functions are in multiple files and you want to simplify distribution.
- You do not want to expose the helper functions' source code.

The helper functions are stored as LLVM IR, an assembly-like language, without comments inside the associated library.

Hardware generation is not necessary for the creation of a `.aoco` file. Compile the OpenCL source file using the `-c` offline compiler command option.

*Note:*

A library can only include OpenCL helper functions. The Intel FPGA SDK for OpenCL Offline Compiler issues an error message if the library contains OpenCL kernels.

- To package an OpenCL source file into a `.aoco` file, invoke the following command: `aoc -c -shared <OpenCL_source_file_name>.cl -o <OpenCL_object_file_name>.aoco`

where the `-shared` offline compiler command option instructs the compiler to create a `.aoco` file that is suitable for inclusion into an OpenCL library.

#### Related Information

- [Packaging Object Files into a Library File](#) on page 173
- [Specifying an OpenCL Library when Compiling an OpenCL Kernel](#) on page 193

### 11.1.4. Packaging an RTL Component for an OpenCL Library

Before creating an OpenCL library file, package each RTL component into a `.aoco` file.

Hardware generation is not necessary for the creation of a `.aoco` file. Compile the OpenCL source file using the `-c` Intel FPGA SDK for OpenCL Offline Compiler command option.

- To package an RTL component into a `.aoco` file, invoke the following command: `aoc -c <RTL component description file name>.xml -o <RTL object file name>.aoco`

#### Related Information

- [Packaging Object Files into a Library File](#) on page 173
- [Verifying the RTL Modules](#) on page 192
- [Specifying an OpenCL Library when Compiling an OpenCL Kernel](#) on page 193

### 11.1.4.1. Restrictions and Limitations in RTL Support for the Intel FPGA SDK for OpenCL Library Feature

The Intel FPGA SDK for OpenCL supports the use of RTL modules in an OpenCL library with some restrictions and limitations.

When creating your RTL module, ensure that it operates within the following restrictions:

- An RTL module must use a single input Avalon streaming interface. That is, a single pair of `ready` and `valid` logic must control all the inputs.

You have the option to provide the necessary Avalon streaming interface ports but declare the RTL module as stall-free. In this case, you do not have to implement proper stall behavior because the Intel FPGA SDK for OpenCL Offline Compiler creates a wrapper for your module. Refer to *XML Syntax of an RTL Module* and *Using an OpenCL Library that Works with Simple Functions (Example 1)* for more syntax and usage information, respectively.

*Note:* You must handle `ivalid` signals properly if your RTL module has an internal state. Refer to *Stall-Free RTL* for more information.

- The RTL module must work correctly regardless of the kernel clock frequency.
- RTL modules cannot connect to external I/O signals. All input and output signals must come from an OpenCL kernel.
- An RTL module must have a `clock` port, a `resetn` port, and Avalon-ST input and output ports (that is, `ivalid`, `ovalid`, `iready`, `oready`). Name the ports as specified here.
- RTL modules that communicate with external memory must have Avalon memory-mapped interface port parameters that match the corresponding Custom Platform parameters. The offline compiler does not perform any width or burst adaptation.
- RTL modules that communicate with external memory must behave as follows:
  - They cannot burst across the burst boundary.
  - They cannot make requests every clock cycle and stall the hardware by monopolizing the arbitration logic. An RTL module must pause its requests regularly to allow other load or store units to execute their operations.
- RTL modules cannot act as stand-alone OpenCL kernels. RTL modules can only be helper functions and be integrated into an OpenCL kernel during kernel compilation.
- Every function call that corresponds to RTL module instantiation is completely independent of other instantiations. There is no hardware sharing.
- Do not incorporate kernel code (that is, functions marked as `kernel`) into a `.aoelib` library file. Incorporating kernel code into the library file causes the offline compiler to issue an error message. You may incorporate helper functions into the library file.
- An RTL component must receive all its inputs at the same time. A single `ivalid` input signifies that all inputs contain valid data.
- You can only set RTL module parameters in the `<RTL module description file name>.xml` specification file, not the OpenCL kernel source file. To use the same RTL module with multiple parameters, create a separate `FUNCTION` tag for each parameter combination.

- You can only pass data inputs to an RTL module by value via the OpenCL kernel code. Do not pass data inputs to an RTL module via pass-by reference, structs, or channels. In the case of channel data, extract the data from the channel first and then pass the extracted the scalar data to the RTL module.

*Note:* Passing data inputs to an RTL module via pass-by reference or structs causes a fatal error to occur in the offline compiler.

- The debugger (for example, GDB for Linux) cannot step into a library function during emulation if the library is built without the debug information. However, irrespective of whether the library is built with or without the debug data, optimization and area reports are not mapped to the individual code line numbers inside a library.
- Names of RTL module source files cannot conflict with the file names of Intel FPGA SDK for OpenCL Offline Compiler IP. Both the RTL module source files and the offline compiler IP files are stored in the `<kernel file name>/system/synthesis/submodules` directory. Naming conflicts causes existing offline compiler IP files in the directory to be overwritten by the RTL module source files.
- The SDK does not support `.qip` files. You must manually parse nested `.qip` files to create a flat list of RTL files.

*Tip:* It is very difficult to debug an RTL module that works correctly on its own but works incorrectly as part of an OpenCL kernel. Double check all parameters under the `ATTRIBUTES` element in the `<RTL module description file name>.xml` file.

- All offline compiler area estimation tools assume that RTL module area is 0. The SDK does not currently support the capability of specifying an area model for RTL modules.

#### Related Information

- [Object Manifest File Syntax of an RTL Module](#) on page 183
- [Using an OpenCL Library that Works with Simple Functions \(Example 1\)](#) on page 197
- [Stall-Free RTL](#) on page 176

### 11.1.5. Verifying the RTL Modules

The creator of an OpenCL library is responsible for verifying the RTL modules within the library, both as stand-alone entities and as part of an OpenCL system.

1. Verify each RTL module using standard hardware verification methods.
2. Modify one of Intel FPGA SDK for OpenCL library design examples to test your RTL modules inside the overall OpenCL system.

This testing step is critical to prevent library users from encountering hardware problems.

It is crucial that you set the values for the `ATTRIBUTES` elements in the object manifest file correctly. Because you cannot simulate the entire OpenCL system, you likely cannot discover problems caused by interface-level errors until hardware runs.

3. Invoke the `aocl library [<command option>]` command.



*Note:* The Intel FPGA SDK for OpenCL library utility performs consistency checks on the object manifest file and source files, with some limitations.

- For a list of supported *<command options>*, invoke the `aocl library` command.
- The `library` utility does not detect errors in values assigned to elements within the `ATTRIBUTES`, `MEM_INPUT`, and `AVALON_MEM` elements in the object manifest file.
- The `library` utility does not detect RTL syntax errors. You must check the *<your\_kernel\_filename>/quartus\_sh\_compile.log* file for RTL syntax errors. However, parsing the errors might be time consuming.

### 11.1.6. Specifying an OpenCL Library when Compiling an OpenCL Kernel

To use an OpenCL library in an OpenCL kernel, specify the library file name and directory when you compile the kernel.

*Important:* Using a library does not reduce kernel compilation time.

- To specify an OpenCL library to the Intel FPGA SDK for OpenCL Offline Compiler, invoke the following command: `aoc -l <library_file_name>.aoclib [-L <library_directory>] <kernel_file_name>.cl`

where the `-l <library_file_name>.aoclib` command option specifies the library file name, and the `-L <library_directory>` command option specifies the directory containing the library files.

You may include multiple instances of `-l <library_file_name>` and `-L <library_directory>` in the offline compiler command.

For example, if you create a library that includes the functions `my_divfd()`, `my_sqrtfd()`, and `my_rsqrtfd()`, the OpenCL kernel code might resemble the following:

```
#include "lib_header.hcl"

kernel void test_lib (
    global double * restrict in,
    global double * restrict out,
    int N) {
    int i = get_global_id(0);
    for (int k = 0; k < N; k++) {
        double x = in[i*N + k];
        out[i*N + k] = my_divfd
            (my_rsqrtfd(x),
             my_sqrtfd(my_rsqrtfd (x)));
    }
}
```

*Note:* Library-related lines are highlighted in bold.

The corresponding `lib_header.h` file might resemble the following:

```
double my_sqrtfd (double x);
double my_rsqrtfd(double x);
double my_divfd(double a, double b);
```

### 11.1.7. Debugging Your OpenCL Library Through Simulation (Preview)

The Intel FPGA SDK for OpenCL simulator assesses the functionality of your OpenCL library.

The Intel FPGA SDK for OpenCL simulator generates a `.aocx` file that runs on an x86-64 Windows or a Linux host. This feature allows you to simulate the functionality of your kernel and iterate on your design without needing to compile your library to hardware and running on the FPGA each time.

**Note:** The performance of the simulator is very slow when compared to that of hardware. So, Intel recommends using a smaller data set for testing.

Use the simulator when you want insight into the dynamic performance of your OpenCL library and more information about the functional correctness of your OpenCL library than emulation or the OpenCL reporting tools provide.

The simulator is cycle accurate, has a netlist identical to generate hardware, and can provide full waveforms for debugging. View the waveforms with Siemens\* EDA QuestaSim\* software.

#### 11.1.7.1. Compiling a Library for Simulation (`-march=simulator`)

To compile an OpenCL library for simulation, include the `-march=simulator` option in your `aoc` command. To enable collecting the waveform during the simulation, include the `-ghdl` option in your `aoc` command. Use the optional `<depth>` attribute with the `-ghdl` option as `-ghdl[=<depth>]` to specify how many levels of hierarchy are logged. If you do not specify a value for the `<depth>` attribute, all waveform signals are logged.

- Before you perform library simulation, perform the following tasks:
  - Install a Custom Platform from your board vendor for your FPGA accelerator boards.
  - Verify that the environment variable `QUARTUS_ROOTDIR_OVERRIDE` points to Intel Quartus Prime Pro Edition software installation folder.
- To simulate library on Windows systems, you need the Microsoft linker and additional compilation time libraries. Verify that the `PATH` environment variable setting includes all the paths described in the *Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables* section of the *Intel FPGA SDK for OpenCL Pro Edition Getting Started Guide*.

The `PATH` environment variable setting must include the path to the `LINK.EXE` file in Microsoft Visual Studio.

- Ensure that your `LIB` environment variable setting includes the path to the Microsoft compilation time libraries.

The compilation time libraries are available with Microsoft Visual Studio.
- Verify that the `LD_LIBRARY_PATH` environment variable setting includes all the paths described in the *Setting the Intel FPGA SDK for OpenCL Pro Edition User Environment Variables* section in the *Intel FPGA SDK for OpenCL Pro Edition Getting Started Guide*.

- To compile a simulation that targets a specific board, invoke the `aoc -march=simulator -ghdl -board=<board_name> <your_kernel_filename>.cl` command.
- For Linux systems, the Intel FPGA SDK for OpenCL Offline Compiler offers symbolic debug support for the debugger.  
The offline compiler debug support allows you to pinpoint the origins of functional errors in your kernel source code.

### 11.1.7.2. Simulating Your OpenCL Library

If you want to view the waveforms generated during simulation, you must install and configure Siemens EDA QuestaSim software.

**Important:** You can also run the emulator and simulator from separate terminal or command prompt sessions.

To run your OpenCL library through the simulator:

1. Run the utility command `aocl linkflags` to find out which libraries are necessary for building a host application. The software lists the libraries for both emulation and regular kernel compilation flows.
2. Build a host application and link it to the libraries from Step 1.

*Tip:* To emulate multiple devices alongside other OpenCL SDKs, link your host application to the Khronos ICD Loader Library *before* you link it to the host runtime libraries. Link the host application to the ICD Loader Library by modifying the `Makefile` for the host application. For more information, see *Linking Your Host Application to the Khronos ICD Loader Library*.

3. If necessary, move the `.aocx` file to a location where the host can find easily, preferably the current working directory.
4. Set the `CL_CONTEXT_MPSIM_DEVICE_INTELFPGA` environment variable to enable the simulation device:

— Windows:

```
set CL_CONTEXT_MPSIM_DEVICE_INTELFPGA=1
```

— Linux:

```
env CL_CONTEXT_MPSIM_DEVICE_INTELFPGA=1
```

**Remember:** When the environment variable `CL_CONTEXT_MPSIM_DEVICE_INTELFPGA` is set, only the simulation devices are available. That is, access to physical boards and the emulation device is disabled.

You might need to set `CL_CONTEXT_COMPILER_MODE_INTELFPGA=3` if the host program cannot find the simulator device.

5. Run your host program.  
To debug your host code and device, you can run your host code in `gdb` or Eclipse.

Running the host program gives you a waveform file, `vsim.wif`, that you can view in QuestaSim software as your host code executes. The `vsim.wif` file is written to the same directory that you run your host program from.

6. If you change your host or kernel program and you want to test it, only recompile the modified host or kernel program and then rerun simulation.

### 11.1.7.3. Troubleshooting Simulator Issues

Review this section to troubleshoot simulator problems you might have when attempting to run a simulation.

#### Windows Compilation Fails - Host Program Reports Corrupt `.aocx` file

During the compilation of the `device.cl` file, your directory path is likely too long. Use the `-o` compiler option to output your compilation results to a shorter path.

#### A `socket=-11` Error Is Logged to `transcript.log`

If you receive the following error message, you are mixing resources from Questa\*-Intel FPGA Edition and ModelSim\* SE:

```
Message: "src/hls_cosim_ipc_socket.cpp:202: void IPCSocketMaster::connect():  
Assertion `sockfd != -1 && "IPCSocketMaster::connect() call to accept() failed" '  
failed."
```

An example of mixing QuestaSim resources is compiling a device with ModelSim SE and then running the host program in Questa\*-Intel FPGA Edition.

#### Running the Host Program Generates a Segmentation Fault

If you receive a segmentation fault when you run your host program, you might be running the emulator and the simulator from the same terminal or command prompt sessions. Remember to unset emulator environment variables before trying to run the simulator.

Try to avoid compiling your device and your host program in the same terminal or command prompt sessions. By using separate sessions, you can avoid possible environment variable conflicts.

#### Simulator Backward Compatibility

In software releases prior to Intel Quartus Prime Pro Edition software version 19.3, the simulator does not work with the Platform Designer.

#### Compatibility with Questa\*-Intel FPGA Starter Edition Software

Questa\*-Intel FPGA Starter Edition software has limitations on design size that prevent it from simulating OpenCL designs. When trying to launch a simulation using Questa\*-Intel FPGA Starter Edition software, you may encounter the following error message:

```
Error: The simulator's process ended unexpectedly.
```

Simulate the designs with Questa\*-Intel FPGA Edition or ModelSim SE software.

### 11.1.8. Using an OpenCL Library that Works with Simple Functions (Example 1)

Intel provides an OpenCL library design example of a simple kernel that uses a library containing RTL implementations of three double-precision functions: `sqrt`, `rsqrt`, and `divide`.

The `library_example1` includes a library, a kernel, and a host system. The `example1.cl` kernel source file includes two kernels. The kernel `test_lib` uses library functions; the kernel `test_builtin` uses built-in functions. The host runs both kernels and then compares their outputs and runtimes. Intel recommends that you use the same strategy to verify your own library functions.

To compile this design example, perform the following tasks:

1. Obtain the `library_example1` from the OpenCL design examples in the `$INTELFPGAOCSDKROOT/examples_aoc` directory.
2. Copy it to a local directory.
3. Follow the instructions in the `README.html` file, which is located in the top-level of the example directory.  
When you run the compiled host program, it should produce the following output:

```
Loading example1.aocx ...
Create buffers
Generate random data for conversion...
Enqueuing both library and builtin in kernels 4 times with global size 65536
Kernel computation using library function took 5.35333 seconds
Kernel computation using built-in function took 5.39949 seconds
Reading results to buffers...
Checking results...
Library function throughput is within 5% of builtin throughput.
PASSED
```

### 11.1.9. Using an OpenCL Library that Works with External Memory (Example 2)

Intel provides an OpenCL library design example of a simple kernel that uses a library containing two RTL modules that communicate with global memory.

The `library_example2` includes a library, a kernel, and a host system. In this example, the RTL code that communicates with global memory is Custom Platform- or Reference Platform-dependent. Ensure that the compilation targets the board that corresponds to the Stratix V Network Reference Platform.

Intel generated the RTL modules `copyElement()` and `sumOfElements()` using the Intel FPGA SDK for OpenCL Offline Compiler, which explains the extra inputs in the code.

The `example2.cl` kernel source file includes two kernels. The kernel `test6` is an NDRange kernel that calls the `copyElement()` RTL function, which copies data from `B[]` to `A[]` and then stores `global_id+100` in `C[]`. The kernel `test11` is a single work-item kernel that uses an RTL function. The `sumOfElements()` RTL function determines the sum of the elements of `A[]` in range `[i, N]` and then adds the rest to `C[i]`.

**Note:** First invocations of `sumOfElements(i=0)` takes more time to execute than later invocations.

To compile this design example, perform the following tasks:

1. Obtain the `library_example2` from the OpenCL design examples in the `$INTELFPGAOCSDKROOT/examples_aoc` directory.
  2. Copy it into a local directory.
  3. Follow the instructions in the `README.html` file, which is located in the top-level of the example directory.
- When you run the compiled host program, it should produce the following output:

```

Loading example2.aocx ...
Running test6
Launching the kernel test6 with globalSize=128 localSize=16
Loading example2.aocx ...
Running test11
Launching the kernel test11 with globalSize=1 localSize=1
PASSED
    
```

**Related Information**

- [Compiling a Kernel for a Specific FPGA Board and Custom Platform \(-board=<board\\_name>\) and \(-board-package=<board\\_package\\_path>\) on page 122](#)
- [Intel FPGA SDK for OpenCL Stratix V Network Reference Platform Porting Guide](#)

**11.1.10. OpenCL Library Command-Line Options**

Both the Intel FPGA SDK for OpenCL Offline Compiler's set of commands and the SDK utility include options you can invoke to perform OpenCL library-related tasks.

**Table 19. Library-Related Intel FPGA SDK for OpenCL Offline Compiler Command Options**

Command Option	Description
<code>-shared</code>	In conjunction with the <code>-rtl</code> command option, compiles an OpenCL source file into an object file ( <code>.aoco</code> ) that you can then include into a library.  <code>aoc -rtl -shared &lt;OpenCL source file name&gt;.cl -o &lt;OpenCL object file name&gt;.aoco</code>
<code>-I=&lt;library_directory&gt;</code>	Adds <code>&lt;library_directory&gt;</code> to the header file search path.  <code>aocl -I &lt;library_header_file_directory&gt; -l &lt;library_file_name&gt;.aoclib &lt;kernel_file_name&gt;.cl</code>
<code>-L=&lt;library_directory&gt;</code>	Adds <code>&lt;library_directory&gt;</code> to the OpenCL library search path. Space after <code>"-L"</code> is optional.  <code>aoc -l=&lt;library_file_name&gt;.aoclib [-L=&lt;library_directory&gt;] &lt;kernel_file_name&gt;.cl</code>
<code>-l=&lt;library_file_name&gt;.aoclib</code>	Specifies the OpenCL library file ( <code>&lt;library_file_name&gt;.aoclib</code> ). Space after <code>-l</code> is optional.
<i>continued...</i>	

Command Option	Description
	<code>aoc -l=&lt;library_file_name&gt;.aoclib [-L=&lt;library directory&gt;] &lt;kernel_file_name&gt;.cl</code>
<code>-library-debug</code>	Generates debug output that relates to libraries. Part of the additional output appears in stdout, the other part appears in the <code>&lt;kernel_file_name&gt;/&lt;kernel_file_name&gt;.log</code> file. <code>aoc -l=&lt;library_file_name&gt;.aoclib -library-debug &lt;kernel_file_name&gt;.cl</code>

**Table 20. Intel FPGA SDK for OpenCL Library Utility (aocl library) Command Options**

Command Option	Description
<code>hdl-comp-pkg &lt;XML_specification_file&gt;.xml</code>	Packages a single HDL component into a <code>.aoco</code> file that you then include into a library. Invoking this command option is similar to invoking <code>aoc -rtl &lt;XML_specification_file&gt;.xml</code> . However, the processing time is faster because the <code>aocl</code> utility does not perform any environment checks. <code>aocl library hdl-comp-pkg &lt;XML_specification_file&gt;.xml -o &lt;output_file&gt;.aoco</code>
<code>-rtl &lt;XML_specification_file&gt;.xml</code>	Same function as <code>hdl-comp-pkg &lt;XML_specification_file&gt;.xml</code> . <code>aocl library -rtl &lt;XML_specification_file&gt;.xml</code>
<code>create</code>	Creates a library file from the <code>.aoco</code> files that you created by invoking the <code>hdl-comp-pkg</code> utility option or the <code>aoc -shared</code> command, and any other <code>.aoclib</code> libraries. <code>aocl library create [-name &lt;library_name&gt;] [-vendor &lt;library_vendor&gt;] [-version &lt;library_version&gt;] [-o &lt;output_file&gt;.aoclib] [.aoco...] [.aoclib...]</code> where <code>-name</code> , <code>-vendor</code> , and <code>-version</code> are optional information strings you can specify and add to the library.
<code>list &lt;library_name&gt;</code>	Lists all the RTL components in the library. Currently, this option is not available for use to list OpenCL functions. <code>aocl library list &lt;library_name&gt;</code>
<code>help</code>	Prints the list of Intel FPGA SDK for OpenCL library utility options and their descriptions on screen. <code>aocl library help</code>

## 11.2. Memory Attributes for Configuring Kernel Memory Systems

The Intel FPGA SDK for OpenCL provides kernel memory attributes that you can apply to constant, local, and private variables or arrays in your design to customize the on-chip memory architecture of the local and private memory systems.

**Attention:** Only apply these local memory kernel attributes to constant, local or private variables.

**Table 21. OpenCL Attributes for Configuring Local Memory**

Attribute	Description
register	Specifies that the variable or array must be carried through the pipeline in registers. Registers can be implemented either exclusively in FFs or in a combination of FFs and RAM-based FIFOs.
memory("impl_type")	Specifies that the variable or array must be implemented in a memory system. Including the memory kernel attribute is equivalent to declaring the variable or array with the <code>__local</code> qualifier. You can pass an optional string argument to specify the memory implementation type. Specify <i>impl_type</i> as either <code>BLOCK_RAM</code> or <code>MLAB</code> to implement the memory using memory blocks (such as M20K) or memory logic array blocks (MLABs), respectively.
numbanks( <i>N</i> )	Specifies that the memory system implementing the variable or array must have <i>N</i> banks, where <i>N</i> is a power-of-2 integer value greater than zero.
bankwidth( <i>N</i> )	Specifies that the memory system implementing the variable or array must have banks that are <i>N</i> bytes wide, where <i>N</i> is a power-of-2 integer value greater than zero.
singlepump	Specifies that the memory system implementing the variable or array must be clocked at the same rate as the component accessing it.
doublepump	Specifies that the memory system implementing the variable or array must be clocked at twice the rate as the component accessing it.
merge("label", "direction")	Forces two or more variables or arrays to be implemented in the same memory system. <i>label</i> is an arbitrary string. Assign the same label to all variables that you want to merge. Specify <i>direction</i> as either <code>width</code> or <code>depth</code> to identify whether the memories should be merged width-wise or depth-wise, respectively.
bank_bits( <i>b</i> <sub>0</sub> , <i>b</i> <sub>1</sub> , ..., <i>b</i> <sub><i>n</i></sub> )	Forces the memory system to split into 2 <sup><i>n</i></sup> banks, with { <i>b</i> <sub>0</sub> , <i>b</i> <sub>1</sub> , ..., <i>b</i> <sub><i>n</i></sub> } forming the bank-select bits. <i>Important:</i> <i>b</i> <sub>0</sub> , <i>b</i> <sub>1</sub> , ..., <i>b</i> <sub><i>n</i></sub> must be consecutive, positive integers. <i>Note:</i> If you specify the <code>numbanks(<i>n</i>)</code> attribute without the <code>bank_bits</code> attribute, the compiler automatically infers the bank-select bits based on the memory access pattern.
private_copies( <i>N</i> )	Specifies that the variable or array declared or accessed inside a pipelined loop has a maximum of <i>N</i> private copies to allow <i>N</i> simultaneous iterations of the loop at any given time, where <i>N</i> is an unsigned integer value. Apply this attribute when the scope of a variable (through its declaration or access pattern) is limited to a loop. If the loop also has a <code>#pragma max_concurrency <i>M</i></code> , the number of private copies created is <code>min(<i>M</i>,<i>N</i>)</code> .
max_replicates( <i>N</i> )	Specifies that the memory implementing the variable or array has no more than <i>N</i> replicates, where <i>N</i> is an integer value greater than 0, to enable simultaneous reads from the datapath.
simple_dual_port	Specifies that the memory implementing the variable or array should have no port that services both reads and writes.
force_pow2_depth( <i>N</i> )	Specifies that the memory implementing the variable or array has a power-of-2 depth. This option is enabled if <i>N</i> is 1 and disabled if <i>N</i> is 0. The default value is 1.

**Table 22. Code Examples for Memory Attributes**

Example Use Case	Syntax
Implements a variable in a register	<code>int __attribute__((register)) a[12];</code>
<i>continued...</i>	



Example Use Case	Syntax
Implements a memory system with eight banks, each with a width of 8 bytes	<pre data-bbox="699 321 1089 380">int __attribute__((memory,     numbanks(8),     bankwidth(8))) b[16];</pre>
Implements a double-pumped memory system with one 128-byte wide bank, and a maximum of two replicates.	<pre data-bbox="699 415 1159 506">int __attribute__((memory,     numbanks(1),     bankwidth(128),     doublepump,     max_replicates(2))) c[32];</pre>

You can also apply memory attributes to data members of a struct. Specify attributes for struct data members in the struct declaration. If you apply attributes to an object instantiation of a struct, then those attributes override the attributes specified in the declaration for struct data members. For example, consider the following code:

```
struct State {
    int array[100] __attribute__((__memory__));
    int reg[4] __attribute__((__register__));
};
__kernel void sum(...) {
    struct State S1;
    struct State S2 __attribute__((__memory__));
    // some uses
}
```

The offline compiler splits *S1* into two variables as *S1.array[100]* (implemented in memory) and *S1.reg[4]* (implemented in registers). However, the compiler ignores attributes applied at struct declaration for object *S2* and does not split it as the *S2* has the attribute *memory* applied to it.

### Related Information

- [Improve Kernel Performance by Banking the Local Memory](#)
- [Optimize Accesses to Local Memory by Controlling the Memory Replication Factor](#)

## 11.2.1. Restrictions on the Use of Variable-specific Attributes

The Intel FPGA SDK for OpenCL Offline Compiler errors out or issues warnings if it detects unsupported usages of the variable-specific attributes or incorrect memory configurations.

Unsupported uses of variable-specific attributes that cause compilation errors:

- You use the kernel attributes in declarations other than constant, local, or private variable declarations (for example, declarations for function parameters, global variable declarations, or function declarations).
- You use the `register` attribute in conjunction with any of the other variable-specific attributes.
- You include both the `singlepump` and `doublepump` attributes in the same variable declaration.

Incorrect memory configurations that cause the offline compiler to issue warnings during compilation:

- The memory configuration that is defined by the variable-specific attributes exceeds the available storage size (for example, specifying eight banks of local memory for an integer variable).

Incorrect memory configurations that cause compilation errors:

- The bank width is smaller than the data access size (for example, bank width is 2 bytes for an array of 4-byte integers).

## 11.3. Kernel Attributes for Reducing the Overhead on Hardware Usage

The Intel FPGA SDK for OpenCL includes kernel attributes that you can include in a single work-item kernel to reduce logic utilization and improve kernel performance. These kernel attributes enable the Intel FPGA SDK for OpenCL Offline Compiler to omit the generation of unnecessary hardware to increase efficiency.

### 11.3.1. Hardware for Kernel Interface

The Intel FPGA SDK for OpenCL Offline Compiler generates hardware around the kernel pipeline. For some OpenCL applications, these interface hardware components are not necessary.

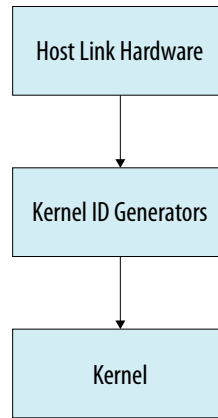
Hardware around the kernel pipeline is necessary for functions such as the following:

- Dispatching IDs for work-items and work-groups
- Communicating with the host regarding kernel arguments and work-group sizes

[Intel FPGA SDK for OpenCL Offline Compiler-Generated Interface Hardware around a Kernel Pipeline](#) illustrates the hardware that the offline compiler generates when it compiles the following kernel:

```
__kernel void my_kernel(global int* arg)
{
    ...
    int sum = 0;
    for(unsigned i = 0; i < n; i++)
    {
        if(sum < m) sum += val;
    }
    *arg = sum;
    ...
}
```

**Figure 34. Intel FPGA SDK for OpenCL Offline Compiler-Generated Interface Hardware around a Kernel Pipeline**



### 11.3.1.1. Omit Hardware that Generates and Dispatches Kernel IDs

The `max_global_work_dim(0)` kernel attribute instructs the Intel FPGA SDK for OpenCL Offline Compiler to omit logic that generates and dispatches global, local, and group IDs into the compiled kernel.

*Tip:* For oneAPI SYCL-specific kernel attribute, refer to [Omit Hardware that Generates and Dispatches Kernel IDs](#) topic in the *Intel oneAPI DPC++ FPGA Optimization Guide*.

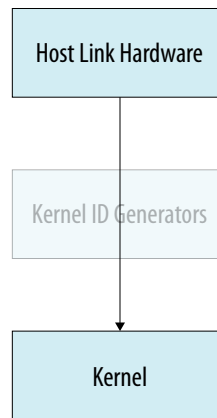
Semantically, the `max_global_work_dim(0)` kernel attribute specifies that the global work dimension of the kernel is zero. Setting this kernel attribute means that the kernel does not use any global, local, or group IDs. The presence of this attribute in the kernel code serves as a guarantee to the offline compiler that the kernel is a single work-item kernel.

When compiling the following kernel, the offline compiler generates interface hardware as illustrated in [Intel FPGA SDK for OpenCL Offline Compiler-Generated Interface Hardware for a Kernel with the `max\_global\_work\_dim\(0\)` Attribute](#).

```
channel int chan_in;
channel int chan_out;

__attribute__((max_global_work_dim(0)))
__kernel void plusK (int N, int k) {
    for (int i = 0; i < N; ++i) {
        int data_in = read_channel_intel(chan_in);
        write_channel_intel(chan_out, data_in + k);
    }
}
```

**Figure 35. Intel FPGA SDK for OpenCL Offline Compiler-Generated Interface Hardware for a Kernel with the `max_global_work_dim(0)` Attribute**



If your current kernel implementation has multiple work-items but does not use global, local, or group IDs, you can use the `max_global_work_dim(0)` kernel attribute if you modify the kernel code accordingly:

1. Wrap the kernel body in a `for` loop that iterates as many times as the number of work-items.
2. Launch the modified kernel with only one work-item.

### 11.3.1.2. Omit Communication Hardware between the Host and the Kernel

The `autorun` kernel attribute instructs the Intel FPGA SDK for OpenCL Offline Compiler to omit logic that is used for communication between the host and the kernel. A kernel that uses the `autorun` attribute starts executing automatically before any kernel that the host launches explicitly. In addition, this kernel restarts automatically as soon as it finishes its execution.

The `autorun` kernel attribute notifies the offline compiler that the kernel runs on its own and does not get enqueued by any host.

To leverage the `autorun` attribute, a kernel must meet all of the following criteria:

- Does not use I/O channels  
*Note:* Kernel-to-kernel channels are supported.
- Does not have any arguments
- Has either the `max_global_work_dim(0)` attribute or the `reqd_work_group_size(X,Y,Z)` attribute. If the `reqd_work_group_size(X,Y,Z)` attribute is used, then  $X$ ,  $Y$ , and  $Z$  must be divisors of  $2^{32}$ .

As mentioned above, kernels with the `autorun` attribute cannot have any arguments and start executing without the host launching them explicitly. As a result, the offline compiler does not need to generate the logic for communication between the host and the kernel. Omitting this logic reduces logic utilization and allows the offline compiler to apply additional performance optimizations.

A typical use case for the `autorun` attribute is a kernel that reads data from one or more kernel-to-kernel channels, processes the data, and then writes the results to one or more channels. When compiling the kernel, the offline compiler generates hardware as illustrated in [Single Work-item Kernel with No Interface Hardware](#).

```
channel int chan_in;
channel int chan_out;

__attribute__((max_global_work_dim(0)))
__attribute__((autorun))
__kernel void plusOne () {
    while(1) {
        int data_in = read_channel_intel(chan_in);
        write_channel_intel(chan_out, data_in + 1);
    }
}
```

**Figure 36. Single Work-item Kernel with No Interface Hardware**



### 11.3.1.3. Omit Hardware to Support the `global_work_offset` Argument in the `clEnqueueNDRangeKernel` API

The `uses_global_work_offset(0)` kernel attribute instructs the Intel FPGA SDK for OpenCL Offline Compiler to omit hardware required to support a non-zero or non-NULL `global_work_offset` argument in NDRange kernels.

This kernel attribute is recommended for all kernels that are always enqueued with a zero or NULL `global_work_offset` argument. When this kernel attribute is set, the Intel FPGA host runtime returns with `CL_INVALID_GLOBAL_OFFSET` error code if a non-zero or non-null `global_work_offset` argument is used to enqueue the kernel.

## 11.4. Kernel Replication Using the `num_compute_units(X,Y,Z)` Attribute

You can replicate your single work-item OpenCL kernel by including the `num_compute_units(X,Y,Z)` kernel attribute.

As mentioned in *Specifying Number of Compute Units*, including the `num_compute_units(N)` kernel attribute in your kernel instructs the Intel FPGA SDK for OpenCL Offline Compiler to generate multiple compute units to process data. The `num_compute_unit(N)` attribute instructs the offline compiler to generate  $N$  identical copies of the kernel in hardware.

**Remember:** To identify the specific compute unit controlling the data-dependent kernel processing, call the `get_compute_id()` intrinsic function.

#### Related Information

- [Customization of Replicated Kernels Using the `get\_compute\_id\(\)` Function](#) on page 206
- [Specifying Number of Compute Units](#) on page 40

### 11.4.1. Customization of Replicated Kernels Using the `get_compute_id()` Function

To create compute units that are slightly different from one another but share a lot of common code, call the `get_compute_id()` intrinsic function in a kernel that also uses the `num_compute_units (X,Y,Z)` attribute.

**Attention:** You can only use the `get_compute_id()` function in a kernel that also uses the `autorun` and `max_global_work_dim(0)` kernel attributes.

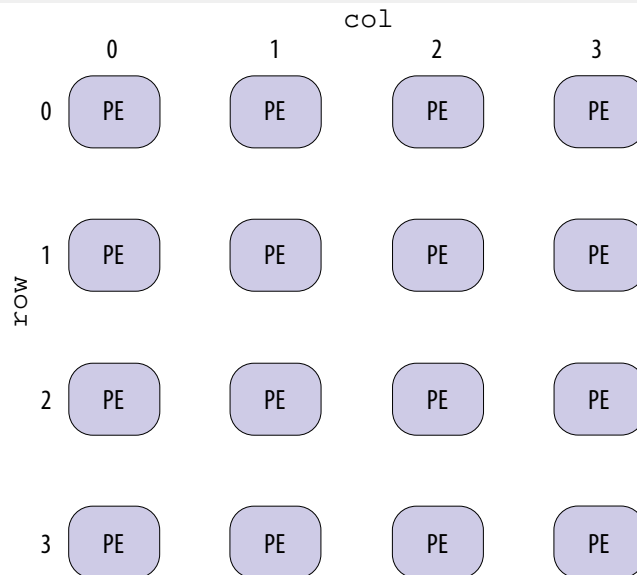
Retrieving compute IDs is a convenient alternative to replicating your kernel in source code and then adding specialized code to each kernel copy. When a kernel uses the `num_compute_units(X,Y,Z)` attribute and calls the `get_compute_id()` function, the Intel FPGA SDK for OpenCL Offline Compiler assigns a unique compute ID to each compute unit. The `get_compute_id()` function then retrieves these unique compute IDs. You can use the compute ID to specify how the associated compute unit should behave differently from the other compute units that are derived from the same kernel source code. For example, you can use the return value of `get_compute_id()` to index into an array of channels to specify which channel each compute unit should read from or write to.

The `num_compute_units` attribute accepts up to three arguments (that is, `num_compute_units(X,Y,Z)`). In conjunction with the `get_compute_id()` function, this attribute allows you to create one-dimensional, two-dimensional, and three-dimensional logical arrays of compute units. An example use case of a 1D array of compute units is a linear pipeline of kernels (also called a daisy chain of kernels). An example use case of a 2D array of compute units is a systolic array of kernels.

### Figure 37. Schematic Diagram of a 4x4 Array of Compute Units

The following example code specifies `num_compute_units(4,4)` in a single work-item kernel results in a 4x4 array that consists of  $4 \times 4 = 16$  compute units.

```
__attribute__((max_global_work_dim(0)))  
__attribute__((autorun))  
__attribute__((num_compute_units(4,4)))  
__kernel void PE() {  
  
    row = get_compute_id(0);  
    col = get_compute_id(1);  
  
    ...  
}
```



For a 3D array of compute units, you can retrieve the X, Y, and Z coordinates of a compute unit in the logical compute unit array using `get_compute_id(0)`, `get_compute_id(1)`, and `get_compute_id(2)`, respectively. In this case, the API is very similar to the API of the work-item's intrinsic functions (that is, `get_global_id()`, `get_local_id()`, and `get_group_id()`).

Global IDs, local IDs, and group IDs can vary at runtime based on how the host invokes the kernel. However, compute IDs are known at compilation time, allowing the offline compiler to generate optimized hardware for each compute unit.

#### 11.4.2. Using Channels with Kernel Copies

To implement channels within compute units (that is, replicated kernel copies), create an array of channels and then index into that array using the return value of `get_compute_id()`.

The example code below implements channels within multiple compute units.

```
#define N 4  
channel int chain_channels[N+1];  
  
__attribute__((max_global_work_dim(0)))  
__kernel void reader(global int *data_in, int size) {  
    for (int i = 0; i < size; ++i) {
```

```

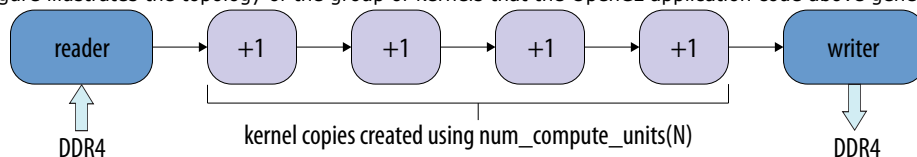
        write_channel_intel(chain_channels[0], data_in[i]);
    }
}

__attribute__((max_global_work_dim(0)))
__attribute__((autorun))
__attribute__((num_compute_units(N)))
__kernel void plusOne() {
    int compute_id = get_compute_id(0);
    int input = read_channel_intel(chain_channels[compute_id]);
    write_channel_intel(chain_channels[compute_id+1], input + 1);
}

__attribute__((max_global_work_dim(0)))
__kernel void writer(global int *data_out, int size) {
    for (int i = 0; i < size; ++i) {
        data_out[i] = read_channel_intel(chain_channels[N]);
    }
}
    
```

**Figure 38. Example Topology of Kernel Copies that Implement Channels**

This figure illustrates the topology of the group of kernels that the OpenCL application code above generates.



**Note:** The implementation of kernel copies is functionally equivalent to defining four separate kernels in your source code and then hard-coding unique indexes for the accesses to `chain_channels[N]`.

## 11.5. Intra-Kernel Registered Assignment Built-In Function

The Intel FPGA SDK for OpenCL Pro Edition provides the built-in function `__fpga_reg()` that you can include in your OpenCL kernel code. The `__fpga_reg()` function directs the Intel FPGA SDK for OpenCL Offline Compiler to insert at least one register between the operand and the return value of the function call.

In general, it is not necessary to include the `__fpga_reg()` function in your kernel code to achieve desired performance.

**Attention:** Intel strongly recommends that you use the `__fpga_reg()` function only if you are an experienced user of the Intel Quartus Prime Pro Edition software performing advanced optimization for a specific target device. You must have sufficient knowledge about the placement of portions of the data path on the FPGA.

Prototype of the `__fpga_reg()` built-in function:

```
T __fpga_reg(T op)
```

where *T* may be any sized type, such as standard OpenCL device data types, or a user-defined struct containing OpenCL types.



Use the `__fpga_reg()` function for the following purposes:

- Break the critical paths between spatially distant portions of a data path, such as between processing elements of a large systolic array.
- Reduce the pressure on placement and routing efforts caused by spatially distinct portions of the kernel implementation.

The `__fpga_reg()` function directs the Intel FPGA SDK for OpenCL Offline Compiler to insert at least one hardware pipelining register on the signal path that assigns the operand to the return value. This built-in function operates as an assignment in the OpenCL programming language, where the operand is assigned to the return value. The assignment has no implicit semantic or functional meaning beyond a standard C assignment. Functionally, you can think of the `__fpga_reg()` function being always optimized away by the offline compiler.

*Note:* The offline compiler does not provide feedback on where you should insert the `__fpga_reg()` function calls in your code. Use the Intel Quartus Prime Pro Edition software to determine where you should insert the calls to address specific aspects of performance.

You may introduce nested `__fpga_reg()` function calls in your kernel code to increase the minimum number of registers that the offline compiler inserts on the assignment path. Because each function call guarantees the insertion of at least one register stage, the number of calls provides a lower limit on the number of registers.

Consider the following example:

```
int out=__fpga_reg(__fpga_reg(in));
```

This line of code directs the offline compiler to insert at least two registers on the assignment path. The offline compiler may insert more than two registers on the path.

## A. Support Statuses of OpenCL Features

---

The Intel FPGA SDK for OpenCL host runtime conforms with the OpenCL platform layer and application programming interface (API), with clarifications and exceptions.

[Support Statuses of OpenCL 1.0 Features](#) on page 210

[Support Statuses of OpenCL 1.2 Features](#) on page 215

[Support Statuses of OpenCL 2.0 Features](#) on page 217

[Intel FPGA SDK for OpenCL Allocation Limits](#) on page 219

### A.1. Support Statuses of OpenCL 1.0 Features

The following sections outline the support statuses of the OpenCL features described in the *OpenCL Specification version 1.0*.

#### A.1.1. OpenCL 1.0 C Programming Language Implementation

OpenCL is based on C99 with some limitations. Section 6 of the *OpenCL Specification version 1.0* describes the OpenCL C programming language. The Intel FPGA SDK for OpenCL conforms with the OpenCL C programming language with clarifications and exceptions. The table below summarizes the support statuses of the features in the OpenCL programming language implementation. OpenCL programming language implementations that are supported with no additional clarifications are not shown.

Support Status column legend:

Symbol	Description
•	The feature is supported, and there might be a clarification for the supported feature in the Notes column
□	The feature is supported with exceptions identified in the Notes column.
X	The feature is not supported.

Section	Feature	Support Status	Notes
6.1.1	<i>Built-in Scalar Data Types</i>		
	double precision float	□	Preliminary support for all double precision float built-in scalar data type. This feature might not conform with the OpenCL Specification version 1.0.
<i>continued...</i>			

Section	Feature	Support Status	Notes
			Currently, the following double precision floating-point functions are expected to conform with the OpenCL Specification version 1.0: add / subtract / multiply / divide / ceil / floor / rint / trunc / fabs / fmax / fmin / sqrt / rsqrt / exp / exp2 / exp10 / log / log2 / log10 / sin / cos / asin / acos / sinh / cosh / tanh / asinh / acosh / atanh / pow / pown / powr / tanh / atan / atan2 / ldexp / log1p / sincos
	half precision float	☐	Support for scalar addition, subtraction and multiplication. Support for conversions to and from single-precision floating point. This feature might not conform with the OpenCL Specification version 1.0. This feature is supported in the Emulator.
6.1.2	<i>Built-in Vector Data Types</i>	☐	Preliminary support for vectors with three elements. Three-element vector support is a supplement to the OpenCL Specification version 1.0.
6.1.3	<i>Other Built-in Data Types</i>	☐	The SDK does not support image or sampler types because the SDK does not support images.
6.2.1	<i>Implicit Conversions</i>	•	Refer to Section 6.2.6: <i>Usual Arithmetic Conversions</i> in the <i>OpenCL Specification version 1.2</i> for an important clarification of implicit conversions between scalar and vector types.
6.2.2	<i>Explicit Casts</i>	•	The SDK allows scalar data casts to a vector with a different element type.
6.5	<i>Address Space Qualifiers</i>	☐	Function scope <code>__constant</code> variables are not supported.
6.6	<i>Image Access Qualifiers</i>	X	The SDK does not support images.
6.7	<i>Function Qualifiers</i>		
6.7.2	<i>Optional Attribute Qualifiers</i>	•	Refer to the <i>Intel FPGA SDK for OpenCL Best Practices Guide</i> for tips on using <code>reqd_work_group_size</code> to improve kernel performance. The SDK parses but ignores the <code>vec_type_hint</code> and <code>work_group_size_hint</code> attribute qualifiers.
6.9	<i>Preprocessor Directives and Macros</i>		
	<code>#pragma directive: #pragma unroll</code>	•	The Intel FPGA SDK for OpenCL Offline Compiler supports only <code>#pragma unroll</code> . You may assign an integer argument to the unroll directive to control the extent of loop unrolling. For example, <code>#pragma unroll 4</code> unrolls four iterations of a loop. By default, an unroll directive with no unroll factor causes the offline compiler to attempt to unroll the loop fully. Refer to the <i>Intel FPGA SDK for OpenCL Best Practices Guide</i> for tips on using <code>#pragma unroll</code> to improve kernel performance.
	<code>__ENDIAN_LITTLE__</code> defined to be value 1	•	The target FPGA is little-endian.
	<code>__IMAGE_SUPPORT__</code>	X	<code>__IMAGE_SUPPORT__</code> is undefined; the SDK does not support images.
6.10	<i>Attribute Qualifiers</i> —The offline compiler parses attribute qualifiers as follows:		
6.10.3	<i>Specifying Attributes of Variables</i> <code>—endian</code>	X	
<b>continued...</b>			

Section	Feature	Support Status	Notes
6.10.4	<i>Specifying Attributes of Blocks and Control-Flow-Statements</i>	X	
6.10.5	<i>Extending Attribute Qualifiers</i>	•	The offline compiler can parse attributes on various syntactic structures. It reserves some attribute names for its own internal use. Refer to the <i>Intel FPGA SDK for OpenCL Best Practices Guide</i> for tips on how to optimize kernel performance using these kernel attributes.
6.11.2	<i>Math Functions</i>		
	built-in math functions	□	Preliminary support for double precision floating-point built-in math functions that might not conform with the OpenCL Specification version 1.0.
	built-in half_ and native_ math functions	□	Preliminary support for double precision floating-point built-in half_ and native_ math functions that might not conform with the OpenCL Specification version 1.0.
6.11.5	<i>Geometric Functions</i>	□	Preliminary support for double precision floating-point built-in geometric functions. These functions might not conform with the OpenCL Specification version 1.0. Refer to <i>Argument Types for Built-in Geometric Functions</i> for a list of built-in geometric functions supported by the SDK.
6.11.8	<i>Image Read and Write Functions</i>	X	The SDK does not support images.
6.11.9	<i>Synchronization Functions—the barrier synchronization function</i>	•	Clarifications and exceptions: <ul style="list-style-type: none"> <li>• If a kernel specifies the reqd_work_group_size or max_work_group_size attribute, barrier supports the corresponding number of work-items.</li> <li>• If neither attribute is specified, a barrier is instantiated with a default limit of 128 work-items.</li> </ul> The work-item limit is the maximum supported work-group size for the kernel; this limit is enforced by the runtime.
6.11.1 1	<i>Async Copies from Global to Local Memory, Local to Global Memory, and Prefetch</i>	•	The implementation is naive: Work-item (0,0,0) performs the copy and the wait_group_events is implemented as a barrier. <ul style="list-style-type: none"> <li>• If a kernel specifies the reqd_work_group_size or max_work_group_size attribute, wait_group_events supports the corresponding number of work-items.</li> <li>• If neither attribute is specified, wait_group_events is instantiated with a default limit of 256 work-items.</li> </ul>

### Related Information

- [Intel FPGA SDK for OpenCL Pro Edition Best Practices Guide](#)
- [Argument Types for Built-in Geometric Functions](#) on page 213

## A.1.2. OpenCL C Programming Language Restrictions

The Intel FPGA SDK for OpenCL conforms with the OpenCL Specification restrictions on specific programming language features, as described in section 6.8 of the *OpenCL Specification version 1.0*.

**Important:** The Intel FPGA SDK for OpenCL Offline Compiler does not enforce restrictions on certain disallowed programming language features. Ensure that your kernel code does not contain features that the OpenCL Specification version 1.0 does not support.

The following table only describes additional clarifications for the OpenCL C Programming language restrictions enforced by the Intel FPGA SDK for OpenCL. All other restrictions exactly match section 6.8 of the *OpenCL Specification version 1.0*.

Feature	Support Status	Notes
structure-type kernel arguments	X	Convert structure arguments to a pointer to a structure in global memory.
irreducible control flow	X	The offline compiler returns a diagnostic error for this restriction, but this feature is not supported.
writes to memory of built-in types less than 32 bits in size	□	Store operations less than 32 bits in size might result in lower memory performance but are supported.
elements of a <code>struct</code> or a <code>union</code> belonging to different address spaces	X	The offline compiler does not enforce this restriction. <b>Warning:</b> Assigning elements of a <code>struct</code> or a <code>union</code> to different address spaces might cause a fatal error.

Support Status column legend:

Symbol	Description
•	The feature is supported, and there might be a clarification for the supported feature in the Notes column
□	The feature is supported with exceptions identified in the Notes column.
X	The feature is not supported.

### A.1.3. Argument Types for Built-in Geometric Functions

The Intel FPGA SDK for OpenCL supports scalar and vector argument built-in geometric functions with certain limitations.

Function	Argument Type	
	float	double
<code>cross</code>	•	•
<code>dot</code>		•
<code>distance</code>		•
<code>length</code>		•
<code>normalize</code>		•
<code>fast_distance</code>		—
<code>fast_length</code>		—
<code>fast_normalize</code>		—

Support Status legend:

Symbol	Description
•	The feature is supported.
–	The feature is not defined in the OpenCL Specification.

### A.1.4. Numerical Compliance Implementation

Section 7 of the *OpenCL Specification version 1.0* describes features of the C99 and IEEE 754 standards that OpenCL-compliant devices must support. The Intel FPGA SDK for OpenCL operates on 32-bit and 64-bit floating-point values in IEEE Standard 754-2008 format, but not all floating-point operators have been implemented.

The table below summarizes the implementation statuses of the floating-point operators:

Section	Feature	Support Status	Notes
7.1	<i>Rounding Modes</i>	<input type="checkbox"/>	Conversion between integer and single and half precision floating-point types support all rounding modes. Conversions between integer and double precision floating-point types support all rounding modes on a preliminary basis that might not conform with the OpenCL Specification version 1.0.
7.2	<i>INF, NaN and Denormalized Numbers</i>	<input type="checkbox"/>	Infinity (INF) and Not a Number (NaN) results for single precision operations are generated in a manner that conforms with the OpenCL Specification version 1.0. Most operations that handle denormalized numbers are flushed prior to and after a floating-point operation. Preliminary support for double precision floating-point operation that might not conform with the OpenCL Specification version 1.0.
7.3	<i>Floating-Point Exceptions</i>	X	
7.4	<i>Relative Error as ULPs</i>	<input type="checkbox"/>	Single precision floating-point operations conform with the numerical accuracy requirements for an embedded profile of the OpenCL Specification version 1.0. Preliminary support for double precision floating-point operation that might not conform with the OpenCL Specification version 1.0.
7.5	<i>Edge Case Behavior</i>	•	

### A.1.5. Image Addressing and Filtering Implementation

The Intel FPGA SDK for OpenCL does not support image addressing and filtering. The SDK does not support images.

### A.1.6. Atomic Functions

Section 9 of the *OpenCL Specification version 1.0* describes a list of optional features that some OpenCL implementations might support. The Intel FPGA SDK for OpenCL supports atomic functions conditionally. The implementation of the supported functions might not conform with the *OpenCL Specification Version 1.0*.

- Section 9.5: *Atomic Functions for 32-bit Integers*—The SDK supports all 32-bit global and local memory atomic functions. The SDK also supports 32-bit atomic functions described in Section 6.11.11 of the *OpenCL Specification version 1.1* and Section 6.12.11 of the *OpenCL Specification version 1.2*.
  - The SDK does not support 64-bit atomic functions described in Section 9.7 of the *OpenCL Specification version 1.0*.

**Attention:** The use of atomic functions might lower the performance of your design. The operating frequency of the hardware might decrease further if you implement more than one type of atomic functions (for example, `atomic_add` and `atomic_sub`) in the kernel.

### A.1.7. Embedded Profile Implementation

Section 10 of the *OpenCL Specification version 1.0* describes the OpenCL embedded profile. The Intel FPGA SDK for OpenCL conforms with the OpenCL embedded profile with clarifications and exceptions.

The table below summarizes the clarifications and exceptions to the OpenCL embedded profile:

Clause	Feature	Support Status	Notes
1	64-bit integers	•	64-bit integers are supported as a supplement to the <i>OpenCL Specification version 1.0</i> for the Embedded Profile.
2	3D images	X	The SDK does not support images.
3	Create 2D and 3D images with <code>image_channel_data_type</code> values	X	The SDK does not support images.
4	Samplers	X	The SDK does not support samplers.
5	Rounding modes	•	The default rounding mode for <code>CL_DEVICE_SINGLE_FP_CONFIG</code> is <code>CL_FP_ROUND_TO_NEAREST</code> .
6	Restrictions listed for single precision basic floating-point operations	X	
7	half type	X	This clause of the <i>OpenCL Specification version 1.0</i> does not apply to the SDK.
8	Error bounds listed for conversions from <code>CL_UNORM_INT8</code> , <code>CL_SNORM_INT8</code> , <code>CL_UNORM_INT16</code> and <code>CL_SNORM_INT16</code> to float	•	

## A.2. Support Statuses of OpenCL 1.2 Features

The following sections outline the support statuses of the OpenCL features described in the *OpenCL Specification version 1.2*.

### A.2.1. OpenCL 1.2 Runtime Implementation

The Intel FPGA SDK for OpenCL supports the implementation of sub-buffer objects. For more information about sub-buffer objects, refer to section 5.2 of the *OpenCL Specification version 1.2*.

The SDK also supports the implementation of the following APIs:

- `clSetMemObjectDestructorCallback`
- `clGetKernelArgInfo`
- `clSetEventCallback`

For more information about these APIs, refer to sections 5.4.1, 5.7.3, and 5.9 of the *OpenCL Specification 1.2*, respectively.

#### Related Information

[OpenCL Specification version 1.2](#)

### A.2.2. OpenCL 1.2 C Programming Language Implementation

The Intel FPGA SDK for OpenCL supports a number of OpenCL C programming language features that are specified section 6 of the *OpenCL Specification version 1.2*. The SDK conforms with the OpenCL C programming language with clarifications and exceptions.

**Attention:** The support status "●" means that the feature is supported, and there might be a clarification for the supported feature in the Notes column. The support status "□" means that the feature is supported with exceptions identified in the Notes column.

**Table 23. Support Statuses of OpenCL 1.2 C Programming Language Features**

Section	Feature	Support Status	Notes
6.1.3	Other Built-in Data Types	●	This feature might not conform with the OpenCL Specification version 1.2.
6.12.12	<i>Miscellaneous Vector Functions</i>	●	The SDK supports implementations of the following additional built-in vector functions: <ul style="list-style-type: none"> <li>• <code>vec_step</code></li> <li>• <code>shuffle</code></li> <li>• <code>shuffle2</code></li> </ul>
6.12.13	<i>printf</i>	●	This feature might not conform with the OpenCL Specification version 1.2. See below for details.
<p>The <code>printf</code> function in OpenCL has syntax and features similar to the <code>printf</code> function in C99, with a few exceptions. For details, refer to the <i>OpenCL Specification version 1.2</i>.</p> <p>To use a <code>printf</code> function, there are no requirements for special compilation steps, buffers, or flags. You can compile kernels that include <code>printf</code> instructions with the usual <code>aoc</code> command.</p> <p>During kernel execution, <code>printf</code> data is stored in a global <code>printf</code> buffer that the Intel FPGA SDK for OpenCL Offline Compiler allocates automatically. The size of this buffer is 64 kB; the total size of data arguments to a <code>printf</code> call should not exceed this size. When kernel execution completes, the contents of the <code>printf</code> buffer are printed to standard output. The format string for a <code>printf</code> statement cannot exceed 256 characters.</p> <p>Buffer overflows are handled seamlessly; <code>printf</code> instructions can be executed an unlimited number of times. However, if the <code>printf</code> buffer overflows, kernel pipeline execution stalls until the host reads the buffer and prints the buffer contents.</p>			
<b>continued...</b>			



Section	Feature	Support Status	Notes
			<p>Because <code>printf</code> functions store their data into a global memory buffer, the performance of your kernel drops if it includes such functions.</p> <p>There are no use limitations on <code>printf</code> functions. You can use <code>printf</code> instructions inside <code>if-then-else</code> statements, loops, and so on. A kernel can contain multiple <code>printf</code> instructions executed by multiple work-items.</p> <p>Format string arguments and literal string arguments of <code>printf</code> calls are transferred to the host system from the FPGA using a special memory region. This memory region can overflow if the total size of the <code>printf</code> string arguments is large (3000 characters or less is usually safe in a typical OpenCL application). If there is an overflow, the error message cannot parse auto-discovery string at byte offset 4096 is printed during host program execution.</p> <p>Output from <code>printf</code> is never intermixed, even though work-items may execute <code>printf</code> functions concurrently. However, the order of concurrent <code>printf</code> execution is not guaranteed. In other words, <code>printf</code> outputs might not appear in program order if the <code>printf</code> instructions are in concurrent datapaths.</p>

#### Related Information

[OpenCL Specification version 1.2](#)

### A.3. Support Statuses of OpenCL 2.0 Features

The following sections outline the support statuses of the OpenCL features described in the *OpenCL Specification version 2.0*.

#### A.3.1. OpenCL 2.0 Headers

The Intel FPGA SDK for OpenCL provides unified headers that support all OpenCL versions.

To use OpenCL 2.0 features, include the following two lines in your code:

```
#define CL_TARGET_OPENCL_VERSION 200
#include <CL/opencl.h>
```

**Attention:** The SDK currently does not support all OpenCL 2.0 APIs. If you make a call to an unsupported API, the call returns with an error code to indicate that the API is not fully supported.

#### A.3.2. OpenCL 2.0 Runtime Implementation

The Intel FPGA SDK for OpenCL offers preliminary support for shared virtual memory implementation, as described in section 5.6 of the *OpenCL Specification version 2.0*. For more information on shared virtual memory, refer to section 5.6 of the *OpenCL Specification version 2.0*.

**Important:** Refer to your board's specifications to verify that your board supports shared virtual memory.

#### Related Information

[OpenCL Specification version 2.0 \(API\)](#)

#### A.3.3. OpenCL 2.0 C Programming Language Restrictions for Pipes

The Intel FPGA SDK for OpenCL offers preliminary support of OpenCL pipes. The following table lists the support statuses of pipe-specific OpenCL C programming language implementations, as described in the *OpenCL Specification version 2.0*

**Attention:** The support status "•" means that the feature is supported. There might be a clarification for the supported feature in the Notes column. A feature that is not supported by the SDK is identified with an "X".

**Table 24. Support Statuses of Built-in Pipe Read and Write Functions**

Details of the built-in pipe read and write functions are available in section 6.13.16.2 of the *OpenCL Specification version 2.0*.

Function	Support Status
<code>int read_pipe (pipe gentype p, gentype *ptr)</code>	•
<code>int write_pipe (pipe gentype p, const gentype *ptr)</code>	•
<code>int read_pipe (pipe gentype p, reserve_id_t reserve_id, uint index, gentype *ptr)</code>	X
<code>int write_pipe (pipe gentype p, reserve_id_t reserve_id, uint index, const gentype *ptr)</code>	X
<code>reserve_id_t reserve_read_pipe (pipe gentype p, uint num_packets)</code> <code>reserve_id_t reserve_write_pipe (pipe gentype p, uint num_packets)</code>	X
<code>void commit_read_pipe (pipe gentype p, reserve_id_t reserve_id)</code> <code>void commit_write_pipe (pipe gentype p, reserve_id_t reserve_id)</code>	X
<code>bool is_valid_reserve_id (reserve_id_t reserve_id)</code>	X

**Table 25. Support Statuses of Built-in Work-Group Pipe Read and Write Functions**

Details of the built-in pipe read and write functions are available in section 6.13.16.3 of the *OpenCL Specification version 2.0*.

Function	Support Status
<code>reserve_id_t work_group_reserve_read_pipe (pipe gentype p, uint num_packets)</code> <code>reserve_id_t work_group_reserve_write_pipe (pipe gentype p, uint num_packets)</code>	X
<code>void work_group_commit_read_pipe (pipe gentype p, reserve_id_t reserve_id)</code> <code>void work_group_commit_write_pipe (pipe gentype p, reserve_id_t reserve_id)</code>	X

**Table 26. Support Statuses of Built-in Pipe Query Functions**

Details of the built-in pipe query functions are available in section 6.13.16.4 of the *OpenCL Specification version 2.0*.

Function	Support Status
<code>uint get_pipe_num_packets (pipe gentype p)</code>	X
<code>uint get_pipe_max_packets (pipe gentype p)</code>	X

**Related Information**

[OpenCL Specification version 2.0 \(C Language\)](#)

## A.4. Intel FPGA SDK for OpenCL Allocation Limits

Item	Limit
Maximum number of contexts	Limited only by host memory size
Maximum number of devices	128
Minimum global memory allocation by runtime	The runtime allocates 64 kB of device memory when the context is created. This memory is reserved for program variables in global address space and for static variables inside functions. If the OpenCL kernel uses the <code>printf</code> function, the runtime allocates an additional 64 kB of device memory.
Maximum number of queues	No static limit
Maximum number of program objects per context	No static limit
Maximum number of event objects per context	Limited only by host memory size
Maximum number of dependencies between events within a context	No static limit
Maximum number of event dependencies per command	No static limit
Maximum number of concurrently running kernels	The total number of queues
Maximum number of kernels per FPGA device	Hardware: no static limit Emulator: 256
Maximum number of arguments per kernel	128
Maximum total size of kernel arguments	256 bytes per kernel
Maximum number of declared variables in the local memory per kernel	No static limit



## B. Intel FPGA SDK for OpenCL Pro Edition Programming Guide Archives

---

For the latest and previous versions of this release notes, refer to [Intel FPGA SDK for OpenCL Pro Edition Programming Guide](#). If a software version is not listed, the guide for the previous software version applies.

## C. Document Revision History of the Intel FPGA SDK for OpenCL Pro Edition Programming Guide

Document Version	Intel Quartus Prime Version	Changes
2022.12.19	22.4	<ul style="list-style-type: none"> <li>Maintenance release.</li> <li>Added a warning about the <code>-fast-compile</code> flag leading to functional failures in <i>Speeding Up Your OpenCL Compilation (-fast-compile)</i>.</li> </ul>
2022.09.26	22.3	<ul style="list-style-type: none"> <li>Maintenance release.</li> <li>Changed the emulator environment variable <code>OCL_TBB_NUM_WORKERS</code> to <code>CL_CONFIG_CPU_TBB_NUM_WORKERS</code> in <i>Emulator Environment Variables</i>.</li> </ul>
2022.06.21	22.2	<ul style="list-style-type: none"> <li>Maintenance release.</li> </ul>
2022.03.28	22.1	<ul style="list-style-type: none"> <li>Maintenance release.</li> </ul>
2021.12.13	21.4	<ul style="list-style-type: none"> <li>Maintenance release.</li> </ul>
2021.10.04	21.3	<ul style="list-style-type: none"> <li>Changed all occurrences of "ModelSim - Intel FPGA Edition" to "Questa*- Intel FPGA Edition".</li> <li>Changed all occurrences of Mentor Graphics to Siemens EDA*.</li> <li>Changed all occurrences of "ModelSim - Intel FPGA Starter Edition" to "Questa*-Intel FPGA Starter Edition".</li> </ul>
2021.06.23	21.2	Changed <code>simple_dual_port_memory</code> to <code>simple_dual_port</code> .
2021.06.21	21.2	<ul style="list-style-type: none"> <li>Added back the offline compiler command option <code>-ffp-reassoc</code> that was removed in 21.1 release.</li> <li>Renamed the offline compiler command option <code>-ffp-reassoc</code> to <code>-ffp-reassociate</code>.</li> </ul>
2021.03.29	21.1	<ul style="list-style-type: none"> <li>Made minor updates to the instructions in <i>Debugging Your OpenCL Kernel on Linux</i>.</li> <li>Added a new topic <i>Pipelining Loops in Non-task Kernels (-auto-pipeline)</i> to describe the <code>-auto-pipeline aoc</code> command option.</li> <li>Renamed the cluster-control attribute <code>__attribute__((stall_enable))</code> as <code>use_stall_enable_clusters</code>. <code>__attribute__((stall_enable))</code> is deprecated.</li> <li>Made minor update to an important note in <i>Specifying Number of SIMD Work-Items</i>.</li> <li>Made minor update to a restriction in <i>Omit Communication Hardware between the Host and the Kernel</i>.</li> <li>Removed the flag <code>-ffp-reassoc</code>.</li> <li>Updated the descriptions of the floating-point pragmas in <i>Floating Point Optimizations (fp contract and fp reassociate Pragma)</i>.</li> </ul>
2020.12.14	20.4	<ul style="list-style-type: none"> <li>Made minor update in the description in <i>Specifying the stall_enable Cluster-control Attribute</i>.</li> <li>Added a minor note about including the <code>CL/cl_ext_intelfpga.h</code> header file, in <i>Partitioning Buffers Across Different Memory Types (Heterogeneous Memory)</i>.</li> </ul>
<b>continued...</b>		

Intel Corporation. All rights reserved. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Intel warrants performance of its FPGA and semiconductor products to current specifications in accordance with Intel's standard warranty, but reserves the right to make changes to any products and services at any time without notice. Intel assumes no responsibility or liability arising out of the application or use of any information, product, or service described herein except as expressly agreed to in writing by Intel. Intel customers are advised to obtain the latest version of device specifications before relying on any published information and before placing orders for products or services.

\*Other names and brands may be claimed as the property of others.

Document Version	Intel Quartus Prime Version	Changes
2020.09.28	20.3	<ul style="list-style-type: none"> <li>Removed <code>-fpc</code> and <code>-fp-relaxed</code> compilation flags because they are deprecated.</li> <li>Updated the description of irreducible control flow feature in <i>OpenCL C Programming Language Restrictions</i>.</li> <li>Added a note about <code>-ffp-reassoc</code> and <code>-ffp-contract=fast</code> flags exception in <i>Creating an Object File From OpenCL Code</i>.</li> <li>Added information about <code>-ghdl [&lt;depth&gt;]</code> attribute in <i>Compiling a Library for Simulation (-march=simulator)</i>.</li> <li>Added a note about the placement of <code>fp contract</code> and <code>fp reassoc</code> pragma statements in <i>Floating Point Optimizations (fp contract and fp reassoc Pragma)</i>.</li> <li>Added a note about unrolling of nested loops with large bounds and an example about partial unroll in <i>Unrolling a Loop (unroll Pragma)</i>.</li> <li>Added a note about emulation performance expectations in the <i>Emulating and Debugging Your OpenCL Kernel</i> topic.</li> <li>Removed the output in <i>Displaying the Compilation Environment of an OpenCL Binary (env)</i> topic since the output was no longer valid.</li> <li>Renamed the topic title <i>Collecting Profile Data During Kernel Execution</i> to <i>Triggering Collection of Profiling Data During Kernel Execution</i> and made minor update to the topic description.</li> <li>Renamed the topic title <i>Profiling Enqueued and Autorun Kernels</i> to <i>Profiling Autorun Kernels</i> and made minor update to the topic description.</li> <li>Made minor update in <i>Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow</i>.</li> <li>Removed topics related to legacy emulator.</li> <li>Minor description update in the <i>Setting up the Emulator</i>.</li> <li>Minor update in the instructions and removed a note related to legacy emulator in the <i>Debugging Your OpenCL Kernel on Linux</i>.</li> <li>Updated a note about the lack of support for Windows BSPs in <i>Listing the Available FPGA Boards and Custom Platforms (-list-boards and -list-board-packages)</i>, <i>Managing an FPGA Board</i>, and <i>Installing an FPGA Board</i>.</li> </ul>
2020.06.22	20.2	<ul style="list-style-type: none"> <li>Minor update to <i>Programming Strategies for Optimizing Pointer-to-Local Memory Size</i> topic.</li> <li>Added a new topic <i>Specifying the stall_enable Cluster-control Attribute</i>.</li> <li>Updated the output sample of <code>-list-board-packages</code> and added a note in <i>Listing the Available FPGA Boards and Custom Platforms (-list-boards and -list-board-packages)</i> and <i>Compiling a Kernel for a Specific FPGA Board and Custom Platform (-board=&lt;board_name&gt; and (-board-package=&lt;board_package_path&gt;)</i> topics.</li> <li>Made minor updates with respect to BSP download instructions in <i>Installing an FPGA Board</i>.</li> <li>Removed the mention of default FCD path in <i>Uninstalling an FPGA Board</i>.</li> <li>Changed the BSP path in the following topics: <ul style="list-style-type: none"> <li><i>Querying the Device Name of Your FPGA Board (diagnose)</i></li> <li><i>Emulating and Debugging Your OpenCL Kernel</i>.</li> <li><i>Using the Legacy Emulator</i></li> </ul> </li> <li>Updated the instruction to accessing the design example in <i>Using an OpenCL Library that Works with Simple Functions (Example 1)</i>.</li> <li>Modified texts to remove references to <code>ap_int&lt;d&gt;</code> and <code>ap_uint&lt;d&gt;</code> in <i>Implementing Arbitrary Precision Integers</i>.</li> <li>Minor update to the description of pipe execution in a loop with multiple work items in <i>Work-item Serial Execution of Pipes</i></li> </ul>
<b>continued...</b>		

Document Version	Intel Quartus Prime Version	Changes
2020.04.13	20.1	<ul style="list-style-type: none"> <li>• Updated the GCC versions in <i>Setting up the Emulator</i>, <i>Emulating your OpenCL Kernel</i>, and <i>Limitations of the Intel FPGA SDK for OpenCL Emulator</i>.</li> <li>• Removed <code>numreadports</code> and <code>numwriteports</code> attributes from <i>Memory Attributes for Configuring Kernel Memory Systems</i>.</li> <li>• Changed all occurrences of "XML specification file" as "object manifest file".</li> <li>• Added a new topic <i>Sharing Multiple Devices Across Multiple Host Programs</i>.</li> <li>• Added more information about multiple command queues in <i>Requirement for Multiple Command Queues to Execute Kernels Concurrently</i>.</li> <li>• Added return values to the APIs and recommendations in <i>API Functions for Interacting with cl_mem Pipe Objects Bound to Host-Accessible Pipe Kernel Arguments</i>.</li> <li>• Added a note about synchronization in <i>Direct Communication with Kernels via Host Pipes</i>.</li> <li>• Changed <code>a10gx_hostch</code> to <code>a10gx_hostpipe</code> in <i>Direct Communication with Kernels via Host Pipes</i>.</li> <li>• Added examples and additional information about performing advanced compiler optimizations in <i>Avalon Streaming Interface</i>.</li> <li>• Added <code>force_pow2_depth</code> memory attribute to <i>Memory Attributes for Configuring Kernel Memory Systems</i>.</li> <li>• Changed the default work group size for kernels with barriers from 256 to 128 in <i>OpenCL 1.0 C Programming Language Implementation and Specifying Work-Group Sizes</i>.</li> <li>• Added a new topic <i>Emulating Applications with a Channel That Reads or Writes to an I/O Channel</i>.</li> <li>• Added a new topic <i>Loop Fusion</i>.</li> <li>• Added a new topic <i>Fusing Adjacent Loops (loop_fuse Pragma)</i>.</li> <li>• Added a new topic <i>Marking Loops to Prevent Automatic Fusion (nofusion Pragma)</i>.</li> <li>• Changed all instances of the memory attribute <code>max_concurrency</code> to <code>private_copies</code>.</li> <li>• Removed lines in the description of memory attribute <code>private_copies</code> that tells N is rounded up to the next power of two.</li> <li>• Added a warning about timing failure and possible solutions in <i>RTL Reset and Clock Signals</i>.</li> <li>• Updated a sentence about clock phase relationship in <i>RTL Reset and Clock Signals</i>.</li> <li>• Removed an entry about in-phase 2x clock in <i>Restrictions and Limitations in RTL Support for the Intel FPGA SDK for OpenCL Library Feature</i>.</li> <li>• Changed <code>-fmax=&lt;fmax target in MHz&gt;</code> to <code>-clock=&lt;clock target in Hz/KHz/MHz/GHz or s/ms/us/ns/ps&gt;</code>.</li> <li>• Removed description about predefined preprocessor macro <code>AOCL_BOARD_&lt;board_name&gt;</code> from <i>Using Predefined Preprocessor Macros in Conditional Compilation</i> topic.</li> <li>• Removed a statement about <code>f<sub>max</sub></code> benefit from <i>Forcing Ring Interconnect for Global Memory (-global-ring)</i>.</li> <li>• Removed <i>Duplicating the Store Ring to Improve the Write Throughput to Global Memory (-duplicate-ring)</i> topic.</li> <li>• Added a new topic <i>Forcing a Single Store Ring to Reduce Area at the Expense of Write Throughput to Global Memory (-force-single-store-ring)</i>.</li> <li>• Added a new topic <i>Forcing Fewer Read Data Reorder Units to Reduce Area at the Expense of Read Throughput to Global Memory (-num-reorder)</i>.</li> </ul>

continued...

Document Version	Intel Quartus Prime Version	Changes
		<ul style="list-style-type: none"> <li>Reorganized and updated <i>Ensuring Compatibility with Other OpenCL SDKs</i> topic for clarity.</li> <li>Removed the details about OpenCL data types in <i>XML Elements for INTERFACE</i>.</li> <li>Added more information to <i>Loop Interleaving Control (max_interleaving Pragma)</i> topic.</li> </ul>
2019.12.16	19.4	<ul style="list-style-type: none"> <li>Updated <i>Guideline for Naming the Kernel</i> about using alphanumeric character, dash, underscore, or dot in file names.</li> <li>Added a subsection about compatibility with ModelSim - Intel FPGA Starter Edition software in the <i>Limitations of the Simulator</i> topic.</li> <li>Updated a limitation about debugger in <i>Restrictions and Limitations in RTL Support for the Intel FPGA SDK for OpenCL Library Feature</i>.</li> <li>Added a note about BSP support in <i>Direct Communication with Kernels via Host Pipes</i>.</li> <li>Removed <code>buffer_location</code> example above table <i>Additional XML Elements to Support External Memory Access</i> as well as description in the named table.</li> <li>Added a note about how to disable debugging and enable optimizations in <i>Emulating Your OpenCL Kernel</i>.</li> <li>Added a note about some source lines going out of order in <i>Debugging Your OpenCL Kernel on Linux</i>.</li> <li>Updated information about the <code>resethn</code> signal in <i>XML Syntax of an RTL Module</i>.</li> <li>Updated the description of the <code>.aoco</code> file with a mention of <code>-save-temps</code> option in <i>Kernel Compilation Flows</i>.</li> <li>Removed <code>.aoco</code> file in <i>One-step Compilation for Simple Kernels</i>.</li> <li>Minor update in <i>Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow</i>.</li> <li>Added a statement about <code>-save-temps</code> in <i>Compiling and Linking Your Kernels or Object Files without Building Hardware (-rtl)</i>.</li> <li>Included the option <code>-save-temps</code> in <i>Specifying the Name of an Intel FPGA SDK for OpenCL Offline Compiler Output File (-o &lt;filename&gt;)</i> and <i>Compiling a Kernel for a Specific FPGA Board and Custom Platform (-board=&lt;board_name&gt;)</i> and <i>(-board-package=&lt;board_package_path&gt;)</i>.</li> <li>Changed <code>cl_ext.h</code> to <code>cl_ext_intelfpga.h</code> in <i>Accessing Custom Platform-Specific Functions</i>.</li> </ul>
2019.10.30	19.3	Removed the pragma directive <code>#pragma OPENCL EXTENSION cl_intel_arbitrary_precision_integers : enable</code> and details about declaring the arbitrary precision integers in the <i>Implementing Arbitrary Precision Integers</i> topic.
2019.09.30	19.3	<ul style="list-style-type: none"> <li>Removed the topic <i>Debugging Your OpenCL System That is Gradually Slowing Down</i> since runtime performance no longer degrades due to event mismanagement.</li> <li>Added a new topic <i>Enabling All Global Memory</i>.</li> <li>Made a minor update in the description by removing some texts in the <i>Forcing Ring Interconnect for Global Memory (-global-ring)</i> topic.</li> <li>Updated <i>Out-of-Order Command Queues</i> topic to describe the support for out-of-order command queues in this release.</li> <li>Updated <i>Requirement for Multiple Command Queues to Execute Kernels Concurrently</i> with more information.</li> </ul>
<b>continued...</b>		



Document Version	Intel Quartus Prime Version	Changes
		<ul style="list-style-type: none"> <li>• In <i>Memory Attributes for Configuring Kernel Memory Systems</i> topic, made the following updates: <ul style="list-style-type: none"> <li>– Deprecated <code>numreadports</code> and <code>numwriteports</code> attributes.</li> <li>– Added <code>max_replicates(N)</code> and <code>simple_dual_port_memory</code> attributes.</li> <li>– Updated descriptions of the remaining attributes.</li> </ul> </li> <li>• Removed some of the restrictions since they no longer were a restriction in the <i>Restrictions on the Use of Variable-specific Attributes</i> topic.</li> <li>• Added another limitation about OpenCL kernels calling HLS tasks in the <i>Limitations of the Incremental Compilation Feature</i> topic.</li> <li>• Added <i>Loop Interleaving Control (max_interleaving Pragma)</i>.</li> <li>• Added a subsection about simulator backward compatibility in the <i>Limitations of the Simulator</i>.</li> <li>• Added a note about GCC 6.3.0 and made some minor updates in the <i>Emulating your OpenCL Kernel</i> topic.</li> <li>• Added information about the library creation tools in the following topics: <ul style="list-style-type: none"> <li>– <i>OpenCL Library</i></li> <li>– <i>Creating Library Objects From OpenCL Code</i></li> <li>– <i>Creating an Object File from OpenCL Code</i></li> <li>– <i>Packaging Multiple Object Files into a Library File</i></li> </ul> </li> <li>• Added <i>Relaxing the Order of Floating-Point Operations (-ffp-reassoc)</i> to describe the new aoc command option.</li> <li>• Added <i>Reducing Floating-Point Rounding Operations (-ffp-contract)</i> to describe the new aoc command option.</li> <li>• Added <i>Floating Point Optimizations (fp contract and fp reassoc Pragma)</i> to describe the new pragmas.</li> <li>• Added a note in <i>Relaxing the Order of Floating-Point Operations (-fp-relaxed)</i> and <i>Reducing Floating-Point Rounding Operations (-fpc)</i> about their deprecation in a future release.</li> <li>• Updated the limits for some of the items in the <i>Allocation Limits</i> topic.</li> <li>• Changed <code>-fp-relaxed</code> to <code>-ffp-reassoc</code> in the <i>Programming Strategies for Inferring the Accumulator</i> topic.</li> <li>• Changed <code>-fp-relaxed</code> to <code>-ffp-reassoc</code> in the <i>Discrepancies in Hardware and Emulator Results</i> topic.</li> <li>• Updated a code snippet in <i>Using OpenCL ICD Extension APIs</i>.</li> <li>• Changed two instances of "emulator" as "legacy emulator" in <i>Simulating Your OpenCL Library</i>.</li> <li>• Changed <code>march=emulator</code> to <code>-march=emulator -legacy-emulator</code> in <i>Compiling a Kernel for Emulation (Legacy Emulator) (-march=emulator -legacy-emulator)</i>.</li> <li>• Made minor updates in <i>Using the Legacy Emulator</i>.</li> <li>• Changed the topic title from <i>Fast Emulator Environment Variables</i> to <i>Emulator Environment Variables</i>. Also added <code>CL_CONFIG_CPU_EMULATE_DEVICES</code> variable description in the topic.</li> <li>• Changed the topic title from <i>Extensions Supported by the Fast Emulator</i> to <i>Extensions Supported by the Emulator</i>. Also changed references to fast emulator as emulator within the topic.</li> <li>• Changed the topic title from <i>Fast Emulator Known Issues</i> to <i>Emulator Known Issues</i>. Also changed references to fast emulator as emulator within the topic.</li> <li>• Removed unrecognized attribute warning for <code>uses_global_work_ofset</code> kernel attribute in the <i>Emulator Known Issues</i> topic.</li> </ul>

**continued...**

Document Version	Intel Quartus Prime Version	Changes
		<ul style="list-style-type: none"> <li>Added another limitation about GCC version in the <i>Limitations of the Emulator</i> topic.</li> <li>Made the following changes in the <i>Compiling a Kernel for Emulation (-march=emulator)</i> topic:               <ul style="list-style-type: none"> <li>Changed <code>-march=emulator -fast-emulator</code> to <code>march=emulator</code>.</li> <li>Removed a bullet point about <code>QUARTUS_ROOTDIR_OVERRIDE</code> variable.</li> <li>Updated the code snippet.</li> </ul> </li> <li>Made the following updates in the <i>Setting up the Emulator</i> topic:               <ul style="list-style-type: none"> <li>Changed the title from <i>Setting up the Fast Emulator</i> to <i>Setting up the Emulator</i>.</li> <li>Removed manual installation of the Intel FPGA SDK for OpenCL.</li> <li>Changed fast emulator references as emulator within the topic.</li> <li>Added a note about emulator requiring GCC6.3.0 or later.</li> </ul> </li> </ul>
2019.07.01	19.2	<ul style="list-style-type: none"> <li>Topics from the following chapters were moved to <i>Intel FPGA SDK for OpenCL Pro Edition Best Practices Guide</i>:               <ul style="list-style-type: none"> <li><i>Reviewing Your Kernel's report.html File</i></li> <li>Entire <i>Profiling Your OpenCL Kernel</i> section</li> </ul> </li> <li>In <a href="#">Intel FPGA SDK for OpenCL Allocation Limits</a> on page 219, limit of maximum number of declared variables in the local memory per kernel updated to 512.</li> <li>In <a href="#">Using Predefined Preprocessor Macros in Conditional Compilation</a> on page 83, updated the description of <code>INTELFPGA_CL</code> macro.</li> <li>Added <a href="#">Disabling Pipelining of a Loop (disable_loop_pipelining Pragma)</a> on page 28 topic to describe the <code>disable_loop_pipelining</code> pragma.</li> </ul>
2019.04.22	19.1	Fixed broken links.
2019.04.01	19.1	<ul style="list-style-type: none"> <li>Added the following new topics:               <ul style="list-style-type: none"> <li><a href="#">Loop Speculation (speculated_iterations Pragma)</a> on page 35</li> <li><a href="#">Specifying the private_copies Memory Attribute</a> on page 41</li> </ul> </li> <li>Updated the <a href="#">Loop Concurrency (max_concurrency Pragma)</a> on page 34 to include a note about unlimited concurrency.</li> <li>Updated <a href="#">Linking Your Host Application to the Khronos ICD Loader Library</a> on page 107 topic for the process to be followed when using ICD for non-standard calls that include <code>IntelFPGA</code> in their function name.</li> <li>In <i>Instrumenting the Kernel Pipeline with Performance Counters (-profile)</i>, added a note about instrumenting the OpenCL kernel pipeline with performance counters.</li> <li>Added <a href="#">Disabling Hardware Kernel Invocation Queue (-no-hardware-kernel-invocation-queue)</a> on page 138 to describe the <code>-no-hardware-kernel-invocation-queue</code> option.</li> <li>Added <a href="#">Forcing Ring Interconnect for Global Memory (-global-ring)</a> on page 130 to describe the <code>-global-ring</code> option.</li> <li>Added <a href="#">Duplicating the Store Ring to Improve the Write Throughput to Global Memory (-duplicate-ring)</a> to describe the <code>-duplicate-ring</code> option.</li> <li>Added <a href="#">Modifying the Handshaking Protocol (-hyper-optimized-handshaking)</a> on page 138 to describe the <code>-hyper-optimized-handshaking</code> option.</li> <li>Added <a href="#">Integer Promotion Rules</a> on page 91 to describe integer promotion rules.</li> <li>In <a href="#">Direct Communication with Kernels via Host Pipes</a> on page 75 and <a href="#">Triggering Collection Profiling Data During Kernel Execution</a> on page 99, added a note about including the <code>CL/cl_ext_intelfpga.h</code> header file.</li> </ul>

continued...

Document Version	Intel Quartus Prime Version	Changes
		<ul style="list-style-type: none"> <li>• Added a new topic <a href="#">Omit Hardware to Support the global_work_offset Argument in the clEnqueueNDRangeKernel API</a> on page 205 about <code>uses_global_work_offset(0)</code> kernel attribute.</li> <li>• In <a href="#">OpenCL 1.0 C Programming Language Implementation</a> on page 210: <ul style="list-style-type: none"> <li>— Changed the support status for 6.1.3 and added a description.</li> <li>— Removed 6.1.4, 6.1.5 and 6.10.2.</li> <li>— Added description for 6.6 and 6.11.8.</li> <li>— For 6.11.2 and 6.11.5, updated the descriptions.</li> <li>— Changed the support status for 6.11.9 and 6.11.11.</li> </ul> </li> <li>• In <a href="#">OpenCL C Programming Language Restrictions</a> on page 212: <ul style="list-style-type: none"> <li>— removed notes for the following features: <ul style="list-style-type: none"> <li>• pointers to functions</li> <li>• bit fields</li> <li>• <code>extern</code>, <code>static</code>, <code>auto</code>, and <code>register</code> storage-class specifiers</li> <li>• recursion</li> <li>• declaration of arguments to <code>__kernel</code> functions of type <code>event_t</code></li> </ul> </li> <li>— removed all features that had no special notes.</li> <li>— removed the following features: <ul style="list-style-type: none"> <li>• pointer assignments between address spaces</li> <li>• images</li> <li>• predefined identifiers</li> </ul> </li> <li>— updated the description of "writes to memory of built-in types less than 32 bits in size".</li> </ul> </li> <li>• In <a href="#">Argument Types for Built-in Geometric Functions</a> on page 213, added the legend about the support status.</li> <li>• In <a href="#">Numerical Compliance Implementation</a> on page 214, updated the description of 7.1</li> <li>• In <a href="#">Embedded Profile Implementation</a> on page 215, added description to clause 4.</li> <li>• In <a href="#">OpenCL 1.2 Runtime Implementation</a> on page 216, removed information about support for image objects.</li> <li>• In <a href="#">OpenCL 1.2 C Programming Language Implementation</a> on page 216, updated the support status and notes of <code>printf</code>, and notes of other built-in data types features.</li> <li>• Removed outdated information and added new information to <a href="#">OpenCL 2.0 Headers</a> on page 217.</li> <li>• In <a href="#">Listing the Available FPGA Boards and Custom Platforms (-list-boards and -list-board-packages)</a> on page 17, added details about <code>-list-board-packages</code> option.</li> <li>• In <a href="#">Installing an FPGA Board (install)</a> on page 20, added more information to step 6 about FCD and added an example in step 5.</li> <li>• In <a href="#">Uninstalling an FPGA Board (uninstall)</a> on page 22, added more information about FCD.</li> <li>• In <a href="#">Compiling a Kernel for a Specific FPGA Board and Custom Platform (-board=&lt;board_name&gt; and (-board-package=&lt;board_package_path&gt;)</a> on page 122, added more information about <code>-board-package=&lt;board_package_path&gt;</code></li> <li>• In <a href="#">Emulator Known Issues</a> on page 152, AutoRun changed to Autorun.</li> <li>• Added <a href="#">Specifying Schedule Fmax Target for Kernels (-clock=&lt;clock_target&gt;)</a> on page 125 to describe the <code>-fmax=&lt;fmax target in MHz&gt;</code></li> </ul>

*continued...*

Document Version	Intel Quartus Prime Version	Changes
		<ul style="list-style-type: none"> <li>• Removed topics <a href="#">Linking to the ICD Loader Library on Windows</a> and <a href="#">Linking to the ICD Loader Library on Linux</a> and added information relevant to ICD and FCD drivers on Windows and Linux in the following topics:               <ul style="list-style-type: none"> <li>– <a href="#">Linking Your Host Application to the Khronos ICD Loader Library</a> on page 107</li> <li>– <a href="#">Displaying Paths to OpenCL Host Runtime and MMD Libraries (ldflags)</a> on page 109</li> <li>– <a href="#">Listing OpenCL Host Runtime and MMD Libraries (ldlibs)</a> on page 109</li> <li>– <a href="#">Displaying Information on OpenCL Host Runtime and MMD Libraries (link-config or linkflags)</a> on page 109</li> </ul> </li> <li>• Added <a href="#">Using OpenCL ICD Extension APIs</a> on page 110 topic describe how to use OpenCL ICD extension APIs.</li> <li>• In <a href="#">Memory Attributes for Configuring Kernel Memory Systems</a> on page 199, made the following changes:               <ul style="list-style-type: none"> <li>– Title modified</li> <li>– Introductory description changed.</li> <li>– Description of <code>register</code> updated with information about FF and FIFO.</li> <li>– <code>memory</code> attribute name and description updated with more information about memory implementation type.</li> <li>– <code>max_concurrency(N)</code> attribute and its description added.</li> <li>– Added information about specifying memory attributes in struct data members.</li> </ul> </li> <li>• In <a href="#">Restrictions on the Use of Variable-specific Attributes</a> on page 201, mentioned constant variable and removed a bullet about <code>numbanks</code> attribute.</li> <li>• Added the topic <a href="#">Specifying Schedule Fmax Target for Kernels (-clock=&lt;clock_target&gt;)</a> on page 125 to describe the <math>f_{max}</math> target attribute and <code>aoc</code> command option.</li> <li>• Revised and updated <a href="#">Emulating and Debugging Your OpenCL Kernel</a> on page 140 to reflect that the fast emulator is now the default and preferred emulator.</li> </ul>
2018.12.24	18.1.1	<ul style="list-style-type: none"> <li>• Updated <a href="#">RTL Reset and Clock Signals</a> on page 181 topic to discuss RTL module access of the system-wide clock that runs at twice the frequency of the OpenCL kernel clock.</li> <li>• Updated the description of the <code>EXPECTED_LATENCY</code> element in <a href="#">XML Elements for ATTRIBUTES</a> on page 184 to include the <code>EXPECTED_LATENCY</code> value requirement for modules that can stall and require signals such as <code>iready</code>.</li> <li>• Updated the description of the <code>IS_FIXED_LATENCY</code> XML element in <a href="#">XML Elements for ATTRIBUTES</a> on page 184 to include the <code>EXPECTED_LATENCY</code> value requirement when <code>IS_FIXED_LATENCY="no"</code> is set.</li> <li>• Added <code>PARAMETER</code> elements to <a href="#">XML Elements for ATTRIBUTES</a> on page 184.</li> <li>• In the second bullet of <a href="#">#unique_230/unique_230_Connect_42_ul_rlt_5pj_w5</a> on page 191, removed "exactly one clock".</li> <li>• Revised <a href="#">Specifying the Name of an Intel FPGA SDK for OpenCL Offline Compiler Output File (-o &lt;filename&gt;)</a> on page 122 to remove the <code>=</code> from the syntax of the <code>-o</code> option. The correct way to specify the name of the output file in the <code>aoc</code> command is <code>-o &lt;filename&gt;</code>. This option was documented incorrectly as <code>-o=&lt;filename&gt;</code>. Also, the syntax for this command option was corrected in examples throughout this publication.</li> </ul>

*continued...*

Document Version	Intel Quartus Prime Version	Changes
2018.09.27	18.1	<ul style="list-style-type: none"> <li>Removed duplicate content in <a href="#">RTL Module Interfaces</a> on page 177.</li> </ul>
2018.09.24	18.1	<ul style="list-style-type: none"> <li>In Intel FPGA SDK for OpenCL Pro Edition, the Intel FPGA SDK for OpenCL Offline Compiler has a new front end. For a summary of changes introduced by this new front end, see <i>Improved Intel FPGA SDK for OpenCL Compiler Front End</i> in the <a href="#">Intel FPGA SDK for OpenCL Pro Edition Release Notes</a>.</li> <li>Added information about new OpenCL simulator preview in <a href="#">Debugging Your OpenCL Library Through Simulation (Preview)</a> on page 194 and the following subtopics: <ul style="list-style-type: none"> <li><a href="#">Compiling a Library for Simulation (-march=simulator)</a> on page 194</li> <li><a href="#">Simulating Your OpenCL Library</a> on page 195</li> <li><a href="#">Troubleshooting Simulator Issues</a> on page 196</li> </ul> </li> <li>Added information about the new OpenCL fast emulator preview in <a href="#">Setting up the Emulator</a> on page 140 and the following subtopics: <ul style="list-style-type: none"> <li><a href="#">Emulator Environment Variables</a> on page 150</li> <li><a href="#">Extensions Supported by the Emulator</a> on page 151</li> <li><a href="#">Emulator Known Issues</a> on page 152</li> </ul> </li> <li>Revised <a href="#">Including Structure Data Types as Arguments in OpenCL Kernels</a> on page 84 to indicate that <code>struct</code> kernel parameters must be passed either by value or as a pointer to <code>struct</code>.</li> <li>Removed references to passing channels by reference because passing channels by reference is no longer supported.</li> </ul>
2018.08.03	18.0	<ul style="list-style-type: none"> <li>Corrected a typo in <a href="#">Partitioning Buffers Across Different Memory Types (Heterogeneous Memory)</a> on page 96. The correct flag <code>CL_MEM_HETEROGENEOUS_INTELFPGA</code>. Previously, the flag was incorrectly called <code>CL_MEM_HETEROGENEOUS_INTEL</code>.</li> </ul>
2018.06.14	18.0	<ul style="list-style-type: none"> <li>Corrected the description of the behavior of the <code>aoc -c</code> command in <a href="#">Compiling Your Kernel without Building Hardware (-c)</a> on page 120. With the <code>-c</code> command option, no folder or subdirectory is created.</li> </ul>
2018.05.23	18.0	<ul style="list-style-type: none"> <li>Corrected code example in <a href="#">Using Channels with Kernel Copies</a> on page 207 so that example compiles successfully.</li> </ul>
2018.05.04	18.0	<ul style="list-style-type: none"> <li>Removed Intel Quartus Prime Standard Edition software-related information.</li> <li>Updated <code>aoc</code> commands for intermediate compilation, and all relevant information, from <code>-c</code> to <code>-rtl</code>. Intel FPGA SDK for OpenCL Pro Edition provides the <code>-rtl</code> flag for intermediate compilation.</li> <li>Increased the maximum number of devices to 128, as documented in the following topics: <ul style="list-style-type: none"> <li><a href="#">Installing an FPGA Board (install)</a> on page 20</li> <li><a href="#">Querying the Device Name of Your FPGA Board (diagnose)</a> on page 22</li> <li><a href="#">Running a Board Diagnostic Test (diagnose &lt;device_name&gt;)</a> on page 23</li> <li><a href="#">Programming the FPGA Offline or without a Host (program &lt;device_name&gt;)</a> on page 23</li> <li><a href="#">Programming the Flash Memory (flash &lt;device_name&gt;)</a> on page 23</li> <li><a href="#">Programming Multiple FPGA Devices</a> on page 112</li> </ul> </li> <li>In <a href="#">One-Step Compilation for Simple Kernels</a> on page 11, updated the figure <i>One-Step OpenCL Kernel Compilation Flow</i> and content of the topic.</li> <li>In <a href="#">Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow</a> on page 12, updated the figure <i>Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow</i> and inserted information on <code>-rtl</code>, fast compilation and incremental compilation.</li> <li>Under <a href="#">Programming Strategies for Optimizing Data Processing Efficiency</a> on page 26, added the topic <a href="#">Loop Concurrency (max_concurrency Pragma)</a> on page 34.</li> </ul>

Document Version	Intel Quartus Prime Version	Changes
		<ul style="list-style-type: none"> <li>• Modified topic title <i>Programming Strategies for Optimizing Local Memory Efficiency</i> to <a href="#">Programming Strategies for Optimizing Pointer-to-Local Memory Size</a> on page 42.</li> <li>• In <a href="#">Emulating I/O Channels</a> on page 53, removed the section <i>Emulating Communication Between a Kernel and a Host or Other Process (Linux only)</i>.</li> <li>• Updated <a href="#">Emulating a Kernel that Passes Pipes or Channels by Value</a> on page 142 to add information on passing pipes or channels by value, and that support for passing channels by reference is deprecated.</li> <li>• Modified the topic title Requirement for Multiple Command Queues in Channels or Pipes Implementation to <a href="#">Requirement for Multiple Command Queues to Execute Kernels Concurrently</a> on page 94, and updated its content.</li> <li>• Added a <a href="#">Compiling Your Kernel Incrementally (-incremental)</a> on page 132 topic with the following subtopics:             <ul style="list-style-type: none"> <li>— <a href="#">The Incremental Compile Report</a> on page 133</li> <li>— <a href="#">Additional Command Options for Incremental Compilation</a> on page 135</li> <li>— <a href="#">Limitations of the Incremental Compilation Feature</a> on page 137</li> </ul> </li> <li>• Added a <a href="#">Compiling Your Kernel with Memory Error Correction Coding (-ecc)</a> on page 137 topic for the ECC early feature.</li> <li>• In <a href="#">Direct Communication with Kernels via Host Pipes</a> on page 75:             <ul style="list-style-type: none"> <li>— Retitled <i>New Optional Kernel Argument Attribute</i> to <a href="#">Optional intel_host_accessible Kernel Argument Attribute</a> on page 76</li> <li>— Retitled <i>New API Functions</i> to <a href="#">API Functions for Interacting with cl_mem Pipe Objects Bound to Host-Accessible Pipe Kernel Arguments</a> on page 76</li> </ul> </li> <li>• In <a href="#">Intel FPGA SDK for OpenCL Allocation Limits</a> on page 219:             <ul style="list-style-type: none"> <li>— Removed reference to the environment variable <code>CL_CONTEXT_PROGRAM_VARIABLES_TOTAL_SIZE_INTELFPGA</code> because the runtime no longer supports it.</li> <li>— Updated the maximum number of queues</li> </ul> </li> <li>• In <a href="#">Object Manifest File Syntax of an RTL Module</a> on page 183, added information on the <code>RESOURCES</code> element.             <ul style="list-style-type: none"> <li>— Added the topic <a href="#">XML Elements for RESOURCES</a> on page 187.</li> </ul> </li> <li>• In <a href="#">Intel FPGA SDK for OpenCL Advanced Features</a> on page 169, added the topic <a href="#">Intel Stratix 10 Design-Specific Reset Requirements for Stall-Free and Stallable RTL Modules</a> on page 182.</li> <li>• In <a href="#">Intel FPGA SDK for OpenCL Advanced Features</a> on page 169, added the topic <a href="#">Intra-Kernel Registered Assignment Built-In Function</a> on page 208</li> </ul>

Date	Version	Changes
December 2017	2017.12.08	<ul style="list-style-type: none"> <li>• Added the following new topics:                             <ul style="list-style-type: none"> <li>– <a href="#">Profiling Autorun Kernels</a></li> <li>– <a href="#">Profiling Autorun Kernels</a> on page 101</li> <li>– <a href="#">Profile Data Acquisition</a> on page 103</li> <li>– <a href="#">Multiple Autorun Profiling Calls</a> on page 102</li> <li>– <a href="#">Developing OpenCL Applications Using Intel Code Builder for OpenCL</a> on page 165</li> <li>– <a href="#">Configuring the Intel Code Builder for OpenCL Offline Compiler Plug-in for Microsoft Visual Studio</a> on page 165</li> <li>– <a href="#">Configuring the Intel Code Builder for OpenCL Offline Compiler Plug-in for Eclipse</a> on page 165</li> <li>– <a href="#">Creating a Session in the Intel Code Builder for OpenCL</a> on page 166</li> <li>– <a href="#">Configuring a Session</a> on page 167</li> </ul> </li> <li>• In <a href="#">Object Manifest File Syntax of an RTL Module</a> on page 183, removed <code>&lt;PARAMETER name="WIDTH" value="32" /&gt;</code> from the XML specification file.</li> </ul>
November 2017	2017.11.06	<ul style="list-style-type: none"> <li>• Moved topics into separate chapters.</li> <li>• Rebranded references to the following:                             <ul style="list-style-type: none"> <li>– The macro ALTERA_CL to INTELFGPA_CL.</li> <li>– The environment variable ALTERAOCLSDKROOT to INTELFGPAOCLSDKROOT.</li> <li>– The environment variable <code>CL_CONTEXT_PROGRAM_VARIABLES_TOTAL_SIZE_ALTERA</code> to <code>CL_CONTEXT_PROGRAM_VARIABLES_TOTAL_SIZE_INTELFGPA</code></li> <li>– <code>clGetExtensionFunctionAddress</code> to <code>clGetExtensionFunctionAddressIntelFGPA</code></li> <li>– The environment variable <code>CL_CONTEXT_EMULATOR_DEVICE_ALTERA</code> to <code>CL_CONTEXT_EMULATOR_DEVICE_INTELFGPA</code></li> <li>– <code>write_channel_altera</code> to <code>write_channel_intel</code></li> <li>– <code>write_channel_nb_altera</code> to <code>write_channel_nb_intel</code></li> <li>– <code>CL_MEM_BANK</code> to <code>CL_CHANNEL</code></li> <li>– <code>CL_MEM_BANK_1_INTEL</code> to <code>CL_CHANNEL_1_INTELFGPA</code></li> <li>– <code>CL_MEM_BANK_2_INTEL</code> to <code>CL_CHANNEL_2_INTELFGPA</code></li> <li>– Arria 10 to Intel Arria 10</li> <li>– Quartus Prime to Intel Quartus Prime</li> <li>– Intel FPGA SDK for OpenCL Profiler to Intel FPGA Dynamic Profiler for OpenCL</li> <li>– TimeQuest Timing Analyzer to Timing Analyzer</li> <li>– Qsys Pro to Platform Designer</li> </ul> </li> <li>• In <a href="#">Intel FPGA SDK for OpenCL FPGA Programming Flow</a> on page 8, added FPGA data flow architecture diagram and related text.</li> <li>• In <a href="#">Intel FPGA SDK for OpenCL Advanced Features</a> section, added <a href="#">RTL Module Interfaces</a> on page 177 to provide example of how RTL module interfaces operate.</li> <li>• Updated the timing diagram in <a href="#">Avalon Streaming Interface</a> on page 178.</li> <li>• In <a href="#">Implementing Blocking Channel Writes</a> on page 49 and <a href="#">Implementing Blocking Channel Reads</a> on page 50, removed "which cannot be a constant" in the definition of <code>&lt;type&gt;</code>.</li> <li>• Added the topic <a href="#">Debugging Your OpenCL System That is Gradually Slowing Down</a>.</li> </ul>

*continued...*

Date	Version	Changes
		<ul style="list-style-type: none"> <li>• Added a link to PLDA website in <a href="#">Compiling Your OpenCL Kernel</a> on page 119.</li> <li>• Updated the last bullet point of <a href="#">Guidelines for Naming the Kernel</a> on page 25 to include the keywords VHDL and Verilog.</li> <li>• In <a href="#">Intel FPGA SDK for OpenCL Advanced Features</a> on page 169, listed the aspects of the design that can be controlled.</li> <li>• In <a href="#">OpenCL Library</a> on page 169, added the expansion of RTL.</li> <li>• Split the sections of <a href="#">Understanding RTL Modules and the OpenCL Pipeline</a> on page 174 into individual topics <a href="#">Overview: Intel FPGA SDK for OpenCL Pipeline Approach</a> on page 174 and <a href="#">Integration of an RTL Module into the Intel FPGA SDK for OpenCL Pipeline</a> on page 175.</li> <li>• In <a href="#">Overview: Intel FPGA SDK for OpenCL Pipeline Approach</a> on page 174, aligned the left-hand example code with image on the right-hand. Moved the bottom portion of the image above the paragraph, which explains the image.</li> <li>• In <a href="#">Integration of an RTL Module into the Intel FPGA SDK for OpenCL Pipeline</a> on page 175, added related links about Avalon-ST.</li> <li>• In <a href="#">Stall-Free RTL</a> on page 176, split the paragraph into steps and added related links.</li> <li>• In <a href="#">Requirements for Deterministic Multiple Work-Item Ordering</a> on page 45, added a third requirement for work-item ordering.</li> <li>• Updated <a href="#">Implementing Nonblocking Channel Reads</a> on page 51.</li> <li>• Added new topic <a href="#">Speeding Up Your OpenCL Compilation (-fast-compile)</a> on page 132 and implemented the convention <code>-option=&lt;value&gt;</code>.</li> <li>• In <a href="#">One-Step Compilation for Simple Kernels</a> on page 11 and <a href="#">Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow</a> on page 12, replaced references to .log file with HTML report and double dash command option with single dash.</li> <li>• In <a href="#">Compiling a Kernel for Emulation (-march=emulator)</a> on page 144, added support for Stratix 10</li> <li>• In <a href="#">Obtaining General Information on Software, Compiler, and Custom Platform</a> on page 16, added a Notice to highlight that double dash and <code>-option &lt;value&gt;</code> conventions of aoc command are deprecated.</li> </ul>

*continued...*



Date	Version	Changes
		<ul style="list-style-type: none"> <li>• Implemented the conventions single dash and <code>-option=&lt;value&gt;</code> in the following topics: <ul style="list-style-type: none"> <li>— <a href="#">Displaying the Compiler Version (-version)</a> on page 16</li> <li>— <a href="#">Listing the Intel FPGA SDK for OpenCL Offline Compiler Command Options (no argument, -help, or -h)</a> on page 17</li> <li>— <a href="#">Listing the Available FPGA Boards and Custom Platforms (-list-boards and -list-board-packages)</a> on page 17</li> <li>— <a href="#">Partitioning Buffers Across Multiple Interfaces of the Same Memory Type</a> on page 95</li> <li>— <a href="#">Specifying the Location of Header Files (-I=&lt;directory&gt;)</a> on page 121</li> <li>— <a href="#">Specifying the Name of an Intel FPGA SDK for OpenCL Offline Compiler Output File (-o &lt;filename&gt;)</a> on page 122</li> <li>— <a href="#">Compiling a Kernel for a Specific FPGA Board and Custom Platform (-board=&lt;board_name&gt;) and (-board-package=&lt;board_package_path&gt;)</a> on page 122</li> <li>— <a href="#">Resolving Hardware Generation Fitting Errors during Kernel Compilation (-high-effort)</a> on page 124</li> <li>— <a href="#">Defining Preprocessor Macros to Specify Kernel Parameters (-D&lt;macro_name&gt;)</a> on page 125</li> <li>— <a href="#">Displaying the Estimated Resource Usage Summary On-Screen (-report)</a> on page 128</li> <li>— <a href="#">Disabling Burst-Interleaving of Global Memory (-no-interleaving=&lt;global_memory_type&gt;)</a> on page 129</li> <li>— <a href="#">Configuring Constant Memory Cache Size (-const-cache-bytes=&lt;N&gt;)</a> on page 130</li> <li>— <a href="#">Relaxing the Order of Floating-Point Operations (-fp-relaxed)</a></li> <li>— <a href="#">Reducing Floating-Point Rounding Operations (-fpc)</a></li> <li>— <a href="#">Emulating Channel Depth</a> on page 143</li> <li>— <a href="#">Compiling a Kernel for Emulation (-march=emulator)</a></li> <li>— <a href="#">Packaging an OpenCL Helper Function File for an OpenCL Library</a> on page 190</li> <li>— <a href="#">OpenCL Library Command-Line Options</a> on page 198</li> <li>— <a href="#">Instrumenting the Kernel Pipeline with Performance Counters (-profile)</a></li> </ul> </li> <li>• In <a href="#">Accessing Custom Platform-Specific Functions</a> on page 103, added related links to ICD loader.</li> <li>• In <a href="#">Specifying Number of Compute Units</a> on page 40, <a href="#">Kernel Replication Using the num_compute_units(X,Y,Z) Attribute</a> on page 205, and <a href="#">Emulating Your OpenCL Kernel</a> on page 144, added a note on compute unit.</li> <li>• In <a href="#">Compiling a Kernel for Emulation (-march=emulator)</a>, added support for Intel Stratix 10.</li> <li>• Added a new topic <a href="#">Discrepancies in Hardware and Emulator Results</a> on page 148.</li> <li>• Added support status column legend <a href="#">OpenCL C Programming Language Restrictions</a> on page 212</li> <li>• Simplified the flowcharts in <a href="#">One-Step Compilation for Simple Kernels</a> on page 11 and <a href="#">Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow</a> on page 12 and updated the texts relevantly.</li> <li>• Removed references to <code>AOCL_BOARD_PACKAGE_ROOT</code> throughout the guide since it is deprecated.</li> <li>• Updated instances of <code>aocl install</code> to <code>aocl install &lt;path_to_customplatform&gt;</code>.</li> <li>• Updated instances of <code>aocl uninstall</code> to <code>aocl uninstall &lt;path_to_customplatform&gt;</code></li> </ul>

*continued...*

Date	Version	Changes
		<ul style="list-style-type: none"> <li>• Added the following new topics for host pipes:               <ul style="list-style-type: none"> <li>– <a href="#">Direct Communication with Kernels via Host Pipes</a> on page 75</li> <li>– <a href="#">Optional intel_host_accessible Kernel Argument Attribute</a> on page 76</li> <li>– <a href="#">API Functions for Interacting with cl_mem Pipe Objects Bound to Host-Accessible Pipe Kernel Arguments</a> on page 76</li> <li>– <a href="#">Creating a Host Accessible Pipe</a> on page 80</li> <li>– <a href="#">Example Use of the cl_intel_fpga_host_pipe Extension</a> on page 81</li> </ul> </li> <li>• In <a href="#">Enabling the Intel FPGA SDK for OpenCL Channels for OpenCL Kernel</a> on page 48, added the pragma to enable the channel extension.</li> <li>• Updated the design example compilation procedures in <a href="#">Using an OpenCL Library that Works with Simple Functions (Example 1)</a> on page 197 and <a href="#">Using an OpenCL Library that Works with External Memory (Example 2)</a> on page 197.</li> <li>• In <a href="#">Restrictions in the Implementation of Intel FPGA SDK for OpenCL Channels Extension</a> on page 46, replaced the Single Site call subsection with Multiple Channel Call Site.</li> </ul>
May 2017	2017.05.08	<ul style="list-style-type: none"> <li>• Rebranded some functions in code examples as follows:               <ul style="list-style-type: none"> <li>– Rebranded <code>read_channel_altera</code> to <code>read_channel_intel</code>.</li> <li>– Rebranded <code>write_channel_altera</code> to <code>write_channel_intel</code>.</li> <li>– Rebranded <code>read_channel_nb_altera</code> to <code>read_channel_nb_intel</code>.</li> <li>– Rebranded <code>write_channel_nb_altera</code> to <code>write_channel_nb_intel</code>.</li> <li>– Rebranded <code>clGetBoardExtensionFunctionAddressAltera</code> to <code>clGetBoardExtensionFunctionAddressIntelFPGA</code>.</li> </ul> </li> <li>• Added <a href="#">Emulating I/O Channels</a> on page 53.</li> <li>• Added <a href="#">Implementing Arbitrary Precision Integers</a> on page 82.</li> <li>• Added <a href="#">Coalescing Nested Loops</a> on page 29.</li> <li>• Added <a href="#">Specifying a Loop Initiation interval (II)</a> on page 33.</li> <li>• Added <a href="#">Emulating Channel Depth</a> on page 143.</li> <li>• Added <a href="#">Avalon Streaming Interface</a> on page 178.</li> <li>• Removed all references to <code>#pragma OPENCL EXTENSION cl_altera_channels : enable</code> because this pragma is not required to implement channels.</li> <li>• Reorganized information related to heterogeneous memory as follows:               <ul style="list-style-type: none"> <li>– Merged <a href="#">Specifying Pointer Size in Memory</a> content into <a href="#">Programming Strategies for Optimizing Pointer-to-Local Memory Size</a> on page 42.</li> <li>– Restructured into three topics:                   <ul style="list-style-type: none"> <li>• <a href="#">Allocating OpenCL Buffers for Manual Partitioning of Global Memory</a> on page 94</li> <li>• <a href="#">Partitioning Buffers Across Multiple Interfaces of the Same Memory Type</a> on page 95</li> <li>• <a href="#">Partitioning Buffers Across Different Memory Types (Heterogeneous Memory)</a> on page 96</li> </ul> </li> <li>– Moved <a href="#">Specifying Buffer Location in Global Memory</a> (previously under <a href="#">Programming Strategies for Optimizing Access Efficiency</a>) content into <a href="#">Partitioning Buffers Across Different Memory Types (Heterogeneous Memory)</a> on page 96.</li> </ul> </li> <li>• Updated <a href="#">Triggering Collection Profiling Data During Kernel Execution</a> on page 99 with a warning about the affect of collecting profile data on kernel launch times.</li> <li>• Updated <a href="#">Compiling Your OpenCL Kernel</a> on page 119 with restrictions on compiling an encrypted <code>.cl</code> file.</li> <li>• Updated <a href="#">Restrictions and Limitations in RTL Support for the Intel FPGA SDK for OpenCL Library Feature</a> on page 191 to indicate that an RTL module must use a single-input Avalon-ST interface to control inputs.</li> </ul>

**continued...**



Date	Version	Changes
		<ul style="list-style-type: none"> <li>• Updated topics affected by changes to the OpenCL profiler as follows:                             <ul style="list-style-type: none"> <li>— Updated <i>Launching the GUI (report)</i> with new command options.</li> <li>— Updated <i>Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow</i> diagram in <i>Multistep Intel FPGA SDK for OpenCL Pro Edition Design Flow</i> topic to reflect the new command options.</li> </ul> </li> <li>• Corrected code example in <a href="#">Implementing Nonblocking Channel Reads</a> on page 51.</li> <li>• Corrected code example in <i>Channel Execution in Loop with Multiple Work-Items</i> section of <a href="#">Work-Item Serial Execution of Channels</a> on page 46.</li> <li>• In <i>Intel FPGA SDK for OpenCL Advanced Features</i> section, made the following updates:                             <ul style="list-style-type: none"> <li>— Updated <a href="#">Interaction between RTL Module and External Memory</a> on page 188 to indicate preferred method for RTL module and external memory interactions.</li> <li>— Updated <a href="#">Potential Incompatibility between RTL Modules and Partial Reconfiguration</a> on page 189 to include link to the partial reconfiguration guidelines in the <i>Quartus Prime Pro Edition Handbook</i>.</li> <li>— Added information about bankbits and mergein <a href="#">Memory Attributes for Configuring Kernel Memory Systems</a> on page 199.</li> <li>— Rebranded some functions in code examples as follows:                                     <ul style="list-style-type: none"> <li>• Rebranded <code>read_channel_altera</code> to <code>read_channel_intel</code>.</li> <li>• Rebranded <code>write_channel_altera</code> to <code>write_channel_intel</code>.</li> </ul> </li> </ul> </li> </ul>
October 2016	2016.10.31	<ul style="list-style-type: none"> <li>• Rebranded the Altera SDK for OpenCL to Intel FPGA SDK for OpenCL.</li> <li>• Rebranded the Altera Offline Compiler to Intel FPGA SDK for OpenCL Offline Compiler.</li> <li>• Deprecated and removed support for big-endian system, resulting in the following documentation changes:                             <ul style="list-style-type: none"> <li>— Removed the topic <i>Compiling a Kernel for a Big-Endian System (--big-endian)</i>.</li> <li>— Removed big-endian (64-bit) from the list of architectures that the host application can target.</li> </ul> </li> <li>• Added the topic <i>Displaying the Compilation Environment of an OpenCL Binary</i> to introduce the <code>aoc env</code> command.</li> <li>• Removed <i>Adding Source References to Optimization Reports (-g)</i> because the offline compiler automatically includes source information in the compiler reports and enables symbolic debug during emulation on an x86 Linux machine.</li> <li>• Added the topic <i>Removing Debug Data from Compiler Reports and Source Code from the .aocx File (-g0)</i> to introduce the <code>-g0 aoc</code> command option.</li> <li>• In <i>Limitations of the Intel FPGA SDK for OpenCL Emulator</i>, removed the limitation "The Emulator does not support half data type".</li> <li>• In <i>Linking Your Host Application to the Khronos ICD Loader Library</i>, provided an update that the Intel-supplied ICD Loader Library supports OpenCL Specification version 1.0 as well as implemented APIs from the OpenCL Specification versions 1.1, 1.2, and 2.0.</li> <li>• In <i>Managing an FPGA Board</i>, provided the following updates:                             <ul style="list-style-type: none"> <li>— Noted that the SDK supports installation of multiple Custom Platforms. To use the SDK utilities on each board in a multi-board installation, the <code>AOCL_BOARD_PACKAGE_ROOT</code> environment variable setting must correspond to the Custom Platform subdirectory of the associated board.</li> <li>— Noted that in a system with multiple Custom Platforms, the host program should use ACD to discover the boards instead of directly linking to the MMD libraries.</li> </ul> </li> </ul>

continued...

Date	Version	Changes
		<ul style="list-style-type: none"> <li>• Added the topic <i>Reviewing Your Kernel's report.html File</i> and included deprecation notice for the <code>analyze-area</code> utility option. As a result of introducing the HTML report, removed the following topics:               <ul style="list-style-type: none"> <li>– <i>Reviewing Your Kernel's Resource Usage Information in the Area Report</i></li> <li>– <i>Accessing the Area Report</i></li> <li>– <i>Layout of the Area Report</i></li> </ul> </li> <li>• In <i>Multistep Design Flow</i>, updated the design steps and the figure <i>The Multistep Intel FPGA SDK for OpenCL Design Flow</i> to replace area report with the HTML report, and remove information on enabling <code>-g</code>.</li> <li>• In <i>Inferring a Register</i>, corrected the text following the code snippet that explained how the offline compiler decide on the implementation of the array in hardware.</li> <li>• In <i>Linking to the ICD Loader Library</i> on Windows, updated the text to improve clarity.</li> <li>• In <i>Support Statuses of OpenCL Features</i> section, made the following updates:               <ul style="list-style-type: none"> <li>– Rebranded the Altera SDK for OpenCL to Intel FPGA SDK for OpenCL.</li> <li>– Rebranded the Altera Offline Compiler to Intel FPGA SDK for OpenCL Offline Compiler.</li> <li>– Modified information in the <i>Intel FPGA SDK for OpenCL Allocation Limits</i> section:                   <ul style="list-style-type: none"> <li>• Updated information regarding minimum global memory allocation by runtime.</li> <li>• Updated the maximum number of queues from 70 to 256.</li> <li>• Updated the maximum number of kernels per FPGA device from 64 to no static limit when compiling to hardware and 256 when compiling to emulator.</li> </ul> </li> <li>– In <i>OpenCL 1.0 C Programming Language Implementation</i>, under the Description column for half-precision float, added a note that this feature is supported in the Emulator. In addition, updated the support status of half-precision float from X to □.</li> <li>– Under <i>Support Statuses of OpenCL 2.0 Features</i>, added the topic <i>OpenCL 2.0 Headers</i> to explain that using the OpenCL 2.0 headers to call unsupported APIs results in an error.</li> </ul> </li> </ul>
May 2016	2016.05.02	<ul style="list-style-type: none"> <li>• Added a schematic diagram of the AOCL programming model in the <i>Altera SDK for OpenCL FPGA Programming Flow</i> section.</li> <li>• Moved the figure <i>The AOCL FPGA Programming Flow</i> to the <i>Altera Offline Compiler Kernel Compilation Flows</i> section.</li> <li>• Updated the figure <i>The Multistep AOCL Design Flow</i> and associated text to include the Review Area Report step.</li> <li>• Added information on the single-cycle floating-point accumulator feature for single work-item kernels. Refer to the <i>Single-Cycle Floating-Point Accumulator for Single Work-Item Kernels</i> section for more information.</li> <li>• Added information in the <i>Emulating Your OpenCL Kernel</i> section on multi-device support for emulation alongside other OpenCL SDKs using ICD.</li> </ul>

**continued...**

Date	Version	Changes
		<ul style="list-style-type: none"> <li>• Included information on the enhanced area report feature:                             <ul style="list-style-type: none"> <li>— Added the option to invoke the <code>analyze-area</code> AOCL utility command to generate an HTML area report.</li> <li>— Included a topic that describes the layout of the HTML area report.</li> </ul> </li> <li>• In <i>Linking to the ICD Loader Library on Windows</i>, removed <code>\$ (AOCL_LDLIBS)</code> from the code example for the modified Makefile.</li> <li>• In the <i>Multiple Work-Item Ordering</i> sections for channels and pipes, modified the characteristics that the AOCL uses to check whether the channel or pipe call is work-item invariant.</li> <li>• Added <a href="#">Intel FPGA SDK for OpenCL Advanced Feature</a> section.</li> <li>• In <i>OpenCL 1.2 Runtime Implementation</i> under <i>Support Statuses of OpenCL Features</i> sections, noted that AOCL supports the <code>clSetEventCallback</code>, <code>clGetKernelArgInfo</code>, and <code>clSetMemObjectDestructorCallback</code> APIs.</li> </ul>
November 2015	2015.11.02	<ul style="list-style-type: none"> <li>• Added the option to invoke the <code>aoc</code> command with no argument to access the Altera Offline Compiler help menu.</li> <li>• Updated the <i>Multiple Host Threads</i> section to specify that the OpenCL host runtime is thread-safe.</li> <li>• Updated the following figure and sections to reflect multiple kernel source file support:                             <ul style="list-style-type: none"> <li>— The figure <i>The AOCL FPGA Programming Flow</i> in the <i>AOCL FPGA Programming Flow</i> section</li> <li>— The <i>Compiling Your Kernel to Create Hardware Configuration File</i> section</li> <li>— The <i>Compiling Your Kernel without Building Hardware (-c)</i> section</li> </ul> </li> <li>• In <i>Multiple Work-Item Ordering for Channels</i>, removed misleading text.</li> <li>• Updated the <i>Overview of Channels Implementation</i> figure.</li> <li>• Updated the following sections on OpenCL pipes:                             <ul style="list-style-type: none"> <li>— <i>Overview of a Pipe Network Implementation</i> figure in <i>Overview of the OpenCL Pipe Functions</i></li> <li>— Emulation support in <i>Restrictions in OpenCL Pipes Implementation</i> section</li> <li>— Replaced erroneous code with the correct syntax</li> <li>— Added link to <i>Implementing I/O Pipes Using the io Attribute in Declaring the Pipe Handle</i></li> </ul> </li> <li>• Added a reminder in <i>Programming an FPGA via the Host</i> that you should release an event object after use to prevent excessive memory usage.</li> <li>• In <i>Support Statuses of OpenCL Features</i> section, made the following updates:                             <ul style="list-style-type: none"> <li>— Categorized feature support statuses and limitations based on OpenCL Specification versions.</li> <li>— Added the following functions to the list of OpenCL-conformant double precision floating-point functions: <code>sinh / cosh / tanh / asinh / acosh / atanh / pow / pown / powr / tanh / atan / atan2 / ldexp / log1p / sincos</code></li> <li>— In <i>OpenCL 1.2 Runtime Implementation</i>, added sub-buffer object support.</li> <li>— In <i>OpenCL 2.0 Runtime Implementation</i>, added preliminary shared virtual memory support.</li> <li>— In <i>Altera SDK for OpenCL Allocation Limits</i>, added a minimum global memory allocation limit by the runtime.</li> </ul> </li> </ul>

**continued...**

Date	Version	Changes
May 2015	15.0.0	<ul style="list-style-type: none"> <li>• In <i>Guidelines for Naming the Kernel</i>, added entry that advised against naming an OpenCL kernel <code>kernel.cl</code>.</li> <li>• In <i>Instrumenting the Kernel Pipeline with Performance Counters (--profile)</i>, specified that you should run the host application from a local disk to avoid potential delays caused by slow network disk accesses.</li> <li>• In <i>Emulating and Debugging Your OpenCL Kernel</i>, modified Caution note to indicate that you must emulate a design targeting an SoC on a non-SoC board.</li> <li>• In <i>Emulating Your OpenCL Kernel</i>, updated command to run the host application and added instruction for overriding default temporary directory containing <code>&lt;process_ID&gt;-libkernel.so</code>.</li> <li>• Introduced the <code>--high-effort aoc</code> command flag in <i>Resolving Hardware Generation Fitting Errors during Kernel Compilation</i>.</li> <li>• In <i>Enabling Double Precision Floating-Point Operations</i>, introduced the <code>OPENCL_EXTENSION</code> pragma for enabling double precision floating-point operations.</li> <li>• Introduced OpenCL pipes support. Refer to <i>Implementing OpenCL Pipes</i> (and subsequent subtopics) and <i>Creating a Pipe Object in Your Host Application</i> for more information.</li> <li>• In <i>AOCL Channels Extension: Restrictions</i>, added code examples to demonstrate how to statically index into arrays of channel IDs.</li> <li>• In <i>Multiple Host Threads</i>, added recommendation for synchronizing OpenCL host function calls in a multi-threaded host application.</li> <li>• Introduced ICD and ACD support. Refer to <i>Linking Your Host Application to the Khronos ICD Loader Library</i> for more information.</li> <li>• Introduced <code>clGetBoardExtensionFunctionAddressAltera</code> for referencing user-accessible functions. Refer to <i>Accessing Custom Platform-Specific Functions</i> for more information.</li> <li>• In <i>Support Statuses of OpenCL Features</i> section, made the following updates: <ul style="list-style-type: none"> <li>– Listed the double precision floating-point functions that the Altera® SDK for OpenCL supports preliminarily.</li> <li>– Added <i>OpenCL C Programming Language Restrictions for Pipes</i>.</li> </ul> </li> </ul>
December 2014	14.1.0	<ul style="list-style-type: none"> <li>• Reorganized information flow. Information is now presented based on the tasks you might perform using the Altera SDK for OpenCL (AOCL) or the Altera RTE for OpenCL.</li> <li>• Removed information pertaining to the <code>--util &lt;N&gt;</code> and <code>-O3</code> Altera Offline Compiler (AOC) options.</li> <li>• Added the following information on PLDA QuickUDP IP core licensing in <i>Compiling Your OpenCL Kernel</i>: <ol style="list-style-type: none"> <li>1. A PLDA QuickUDP IP core license is required for the Stratix V Network Reference Platform or a Custom Platform that uses the QuickUDP IP core.</li> <li>2. Improper installation of the QuickUDP IP core license causes compilation to fail with an error message that refers to the QuickTCP IP core.</li> </ol> </li> <li>• Added reminder that conditionally shifting a large shift register is not recommended.</li> <li>• Removed the <i>Emulating Systems with Multiple Devices</i> section. A new <code>env CL_CONTEXT_EMULATOR_DEVICE_ALTERA=&lt;number_of_devices&gt;</code> command is now available for emulating multiple devices.</li> <li>• Removed language support limitation from the <i>Limitations of the AOCL Emulator</i> section.</li> <li>• In <i>AOCL Allocation Limits</i> under <i>Support Statuses of OpenCL Features</i> section, updated the maximum number of kernels per FPGA device from 32 to 64.</li> </ul>

*continued...*



Date	Version	Changes
June 2014	14.0.0	<ul style="list-style-type: none"> <li>• Removed the <code>--estimate-throughput</code> and <code>--sw-dimm-partition</code> AOC options</li> <li>• Added the <code>-march=emulator</code>, <code>-g</code>, <code>--big-endian</code>, and <code>--profile</code> AOC options</li> <li>• <code>--no-interleaving</code> needs <code>&lt;global_memory_type&gt;</code> argument</li> <li>• <code>-fp-relaxed=true</code> is now <code>--fp-relaxed</code></li> <li>• <code>-fpc=true</code> is now <code>--fpc</code></li> <li>• For non-SoC devices, <code>aocl diagnostic</code> is now <code>aocl diagnose</code> and <code>aocl diagnose &lt;device_name&gt;</code></li> <li>• <code>program</code> and <code>flash</code> need <code>&lt;device_name&gt;</code> arguments</li> <li>• Added <i>Identifying the Device Name of Your FPGA Board</i></li> <li>• Added <i>AOCL Profiler Utility</i></li> <li>• Added <i>AOCL Channels Extension</i> and associated subsections</li> <li>• Added <i>Attributes for Channels</i></li> <li>• Added <i>Match Data Layouts of Host and Kernel Structure Data Types</i></li> <li>• Added <i>Register Inference</i> and <i>Shift Register Inference</i></li> <li>• Added <i>Channels and Multiple Command Queues</i></li> <li>• Added <i>Shared Memory Accesses for OpenCL Kernels Running on SoCs</i></li> <li>• Added <i>Collecting Profile Data During Kernel Execution</i></li> <li>• Added <i>Emulate and Debug Your OpenCL Kernel</i> and associated subsections</li> <li>• Updated <i>AOC Kernel Compilation Flows</i></li> <li>• Updated <code>-v</code></li> <li>• Updated <i>Host Binary Requirement</i></li> <li>• Combined <i>Partitioning Global Memory Accesses</i> and <i>Partitioning Heterogeneous Global Memory Accesses</i> into the section <i>Partitioning Global Memory Accesses</i></li> <li>• Updated <i>AOC Allocation Limits</i> in <i>Appendix A</i></li> <li>• Removed <code>max_unroll_loops</code>, <code>max_share_resources</code>, <code>num_share_resources</code>, and <code>task</code> kernel attributes</li> <li>• Added <code>packed</code>, and <code>aligned(&lt;N&gt;)</code> kernel attributes</li> <li>• In <i>Support Statuses of OpenCL Features</i> section, updated the following AOCL allocation limits: <ul style="list-style-type: none"> <li>— Maximum number of contexts</li> <li>— Maximum number of queues</li> <li>— Maximum number of event objects per context</li> </ul> </li> </ul>
December 2013	13.1.1	<ul style="list-style-type: none"> <li>• Removed the section <code>-W</code> and <code>-Werror</code>, and replaced it with two sections: <code>-W</code> and <code>-Werror</code>.</li> <li>• Updated the following contents to reflect multiple devices support: <ul style="list-style-type: none"> <li>— The figure <i>The AOCL FPGA Programming Flow</i>.</li> <li>— <code>--list-boards</code> section.</li> <li>— <code>-board &lt;board_name&gt;</code> section.</li> <li>— Added the subsection <i>Programming Multiple FPGA Devices</i> under <i>FPGA Programming</i>.</li> </ul> </li> <li>• The following contents were added to reflect heterogeneous global memory support: <ul style="list-style-type: none"> <li>— <code>--no-interleaving</code> section.</li> <li>— <code>buffer_location</code> kernel attribute under <i>Kernel Pragmas and Attributes</i>.</li> <li>— <i>Partitioning Heterogeneous Global Memory Accesses</i> section.</li> </ul> </li> <li>• Modified support status designations in <i>Appendix: Support Statuses of OpenCL Features</i>.</li> <li>• Removed information on OpenCL programming language restrictions from the section <i>OpenCL Programming Language Implementation</i>, and presented the information in a new section titled <i>OpenCL Programming Language Restrictions</i>.</li> </ul>
<b>continued...</b>		

Date	Version	Changes
November 2013	13.1.0	<ul style="list-style-type: none"> <li>• Reorganized information flow.</li> <li>• Updated and renamed <i>Intel FPGA SDK for OpenCL Compilation Flow</i> to <i>AOCL FPGA Programming Flow</i>.</li> <li>• Added figures <i>One-Step AOC Compilation Flow</i> and <i>Two-Step AOC Compilation Flow</i>.</li> <li>• Updated the section <i>Contents of the AOCL Version 13.1</i>.</li> <li>• Removed the following sections:               <ul style="list-style-type: none"> <li>– <i>OpenCL Kernel Source File Compilation</i>.</li> <li>– <i>Using the Altera Offline Kernel Compiler</i>.</li> <li>– <i>Setting Up Your FPGA Board</i>.</li> <li>– <i>Targeting a Specific FPGA Board</i>.</li> <li>– <i>Running Your OpenCL Application</i>.</li> <li>– <i>Consolidating Your Kernel Source Files</i>.</li> <li>– <i>Aligned Memory Allocation</i>.</li> <li>– <i>Programming the FPGA Hardware</i>.</li> <li>– <i>Programming the Flash Memory of an FPGA</i>.</li> </ul> </li> <li>• Updated and renamed <i>Compiling the OpenCL Kernel Source File</i> to <i>AOCL Compilation Flows</i>.</li> <li>• Renamed <i>Passing File Scope Structures to OpenCL Kernels</i> to <i>Use Structure Arguments in OpenCL Kernels</i>.</li> <li>• Updated and renamed <i>Augmenting Your OpenCL Kernel by Specifying Kernel Attributes and Pragmas</i> to <i>Kernel Pragmas and Attributes</i>.</li> <li>• Renamed <i>Loading Kernels onto an FPGA</i> to <i>FPGA Programming</i>.</li> <li>• Consolidated <i>Compiling and Linking Your Host Program</i>, <i>Host Program Compilation Settings</i>, and <i>Library Paths and Links</i> into a single section.</li> <li>• Inserted the section <i>Preprocessor Macros</i>.</li> <li>• Renamed <i>Optimizing Global Memory Accesses</i> to <i>Partitioning Global Memory Accesses</i>.</li> </ul>
June 2013	13.0 SP1.0	<ul style="list-style-type: none"> <li>• Added the section <i>Setting Up Your FPGA Board</i>.</li> <li>• Removed the subsection <i>Specifying a Target FPGA Board</i> under <i>Kernel Programming Considerations</i>.</li> <li>• Inserted the subsections <i>Targeting a Specific FPGA Board</i> and <i>Generating Compilation Reports</i> under <i>Compiling the OpenCL Kernel Source File</i>.</li> <li>• Renamed <i>File Scope __constant Address Space Qualifier</i> to <i>__constant Address Space Qualifiers</i>, and inserted the following subsections:               <ul style="list-style-type: none"> <li>– <i>Function Scope __constant Variables</i>.</li> <li>– <i>File Scope __constant Variables</i>.</li> <li>– <i>Points to __constant Parameters from the Host</i>.</li> </ul> </li> <li>• Inserted the subsection <i>Passing File Scope Structures to OpenCL Kernels</i> under <i>Kernel Programming Considerations</i>.</li> <li>• Renamed <i>Modifying Your OpenCL Kernel by Specifying Kernel Attributes and Pragmas</i> to <i>Augmenting Your OpenCL Kernel by Specifying Kernel Attributes and Pragmas</i>.</li> <li>• Updated content for the <code>unroll</code> pragma directive in the section <i>Augmenting Your OpenCL Kernel by Specifying Kernel Attributes and Pragmas</i>.</li> <li>• Inserted the subsections <i>Out-of-Order Command Queues</i> and <i>Modifying Host Program for Structure Parameter Conversion</i> under <i>Host Programming Considerations</i>.</li> <li>• Updated the sections <i>Loading Kernels onto an FPGA Using clCreateProgramWithBinary</i> and <i>Aligned Memory Allocation</i>.</li> <li>• Updated flash programming instructions.</li> <li>• Renamed <i>Optional Extensions</i> in <i>Appendix B</i> to <i>Atomic Functions</i>, and updated its content.</li> <li>• Removed <i>Platform Layer and Runtime Implementation</i> from <i>Appendix B</i>.</li> </ul>
<b>continued...</b>		



Date	Version	Changes
May 2013	13.0.1	<ul style="list-style-type: none"> <li>• Explicit memory fence functions are now supported; the entry is removed from the table OpenCL Programming Language Implementation.</li> <li>• Updated the section Programming the Flash Memory of an FPGA.</li> <li>• Added the section <i>Modifying Your OpenCL Kernel by Specifying Kernel Attributes and Pragmas</i> to introduce kernel attributes and pragmas that can be implemented to optimize kernel performance.</li> <li>• Added the section Optimizing Global Memory Accesses to discuss data partitioning.</li> <li>• Removed the section <i>Programming the FPGA with the aocl program Command</i> from Appendix A.</li> </ul>
May 2013	13.0.0	<ul style="list-style-type: none"> <li>• Updated compilation flow.</li> <li>• Updated kernel compiler commands.</li> <li>• Included AlteraOpenCL Utility commands. SDK for</li> <li>• Added the section <i>OpenCL Programming Considerations</i>.</li> <li>• Updated flash programming procedure and moved it to <i>Appendix A</i>.</li> <li>• Included a new <code>clCreateProgramWithBinary</code> SDK for FPGA hardware programming flow.</li> <li>• Moved the hostless <code>clCreateProgramWithBinary</code> hardware programming flow to <i>Appendix A</i> under the title <i>Programming the FPGA with the aocl program Command</i>.</li> <li>• Moved updated information on allocation limits and OpenCL language support to <i>Appendix B</i>.</li> </ul>
November 2012	12.1.0	Initial release.