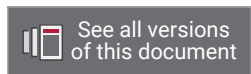


# SDAccel Environment Profiling and Optimization Guide

UG1207 (v2019.1) June 5, 2019



# Revision History

The following table shows the revision history for this document.

| Section                                    | Revision Summary  |
|--|---|
| 06/05/2019 Version 2019.1                  |   |
| <a href="#">Profile Summary Report</a>     | Updated <a href="#">Data Interpretation</a> to reflect current view.  |
| <a href="#">Waveform View</a>              | Updated images and text in <a href="#">Data Interpretation Waveform View</a> and <a href="#">Data Interpretation Live Waveform</a> to reflect current waveform views. |
| <a href="#">Memory Data Transfer Types</a> | Added information about data transfer types.  |
| <a href="#">Top Level Dataflow</a>         | Added discussion about the use of top level dataflow optimization.  |
| General updates                            | Updated figures and minor editorial changes.  |

# Table of Contents

|   |           |
|---|-----------|
| <b>Revision History</b> .....                                       | <b>2</b>  |
| <b>Chapter 1: Introduction</b> .....                                | <b>5</b>  |
| SDAccel Execution Model.....  | 5         |
| SDAccel Build Process.....  | 7         |
| SDAccel Optimization Flow Overview.....                             | 10        |
| <b>Chapter 2: SDAccel Profiling and Optimization Features</b> ..... | <b>16</b> |
| System Estimate.....  | 16        |
| HLS Report.....   | 21        |
| Profile Summary Report.....   | 24        |
| Application Timeline .....  | 30        |
| Waveform View.....  | 38        |
| Guidance.....   | 45        |
| Using Implementation Tools.....                                     | 48        |
| <b>Chapter 3: Kernel Optimization</b> .....                         | <b>51</b> |
| Interface Attributes (Detailed Kernel Trace).....                   | 51        |
| Optimizing Computational Parallelism.....                           | 61        |
| Optimizing Compute Units.....                                       | 73        |
| Optimizing Memory Architecture.....                                 | 76        |
| <b>Chapter 4: Host Optimization</b> .....                           | <b>82</b> |
| Reducing Overhead of Kernel Enqueing.....                           | 82        |
| Data Transfers.....   | 83        |
| Compute Unit Scheduling.....  | 88        |
| Using the clEnqueueMigrateMemObjects API to Transfer Data.....      | 91        |
| <b>Chapter 5: Topological Optimization</b> .....                    | <b>92</b> |
| Multiple Compute Units.....   | 92        |
| Using Multiple DDR Banks.....                                       | 92        |
| <b>Appendix A: Examples</b> .....                                   | <b>96</b> |

|  |           |
|--|-----------|
| <b>Appendix B: Additional Resources and Legal Notices.....</b> | <b>97</b> |
| Xilinx Resources.....  | 97        |
| Documentation Navigator and Design Hubs.....                   | 97        |
| References.....  | 97        |
| Please Read: Important Legal Notices.....                      | 98        |

# Introduction

This guide presents all SDx™ development environment features related to performance analysis of the design. It is also logically structured to assist in the actual performance improvement effort. Dedicated sections are available for the main components of the SDAccel™ environment performance bottlenecks, namely accelerator, PCIe® bus transfer, and host code. Each of these sections is structured to guide the developer from recognizing bottlenecks all the way to solution approaches to increase overall system performance.

**Note:** Performance optimization assumes, as a starting point, a working design intended for performance improvement. If erroneous behavior is encountered, look for guidance in *SDAccel Environment Debugging Guide* ([UG1281](#)).

Similarly, the general concepts regarding coding of host code or accelerator kernels are not explained here; these concepts are introduced in the *SDAccel Environment Programmers Guide* ([UG1277](#)).

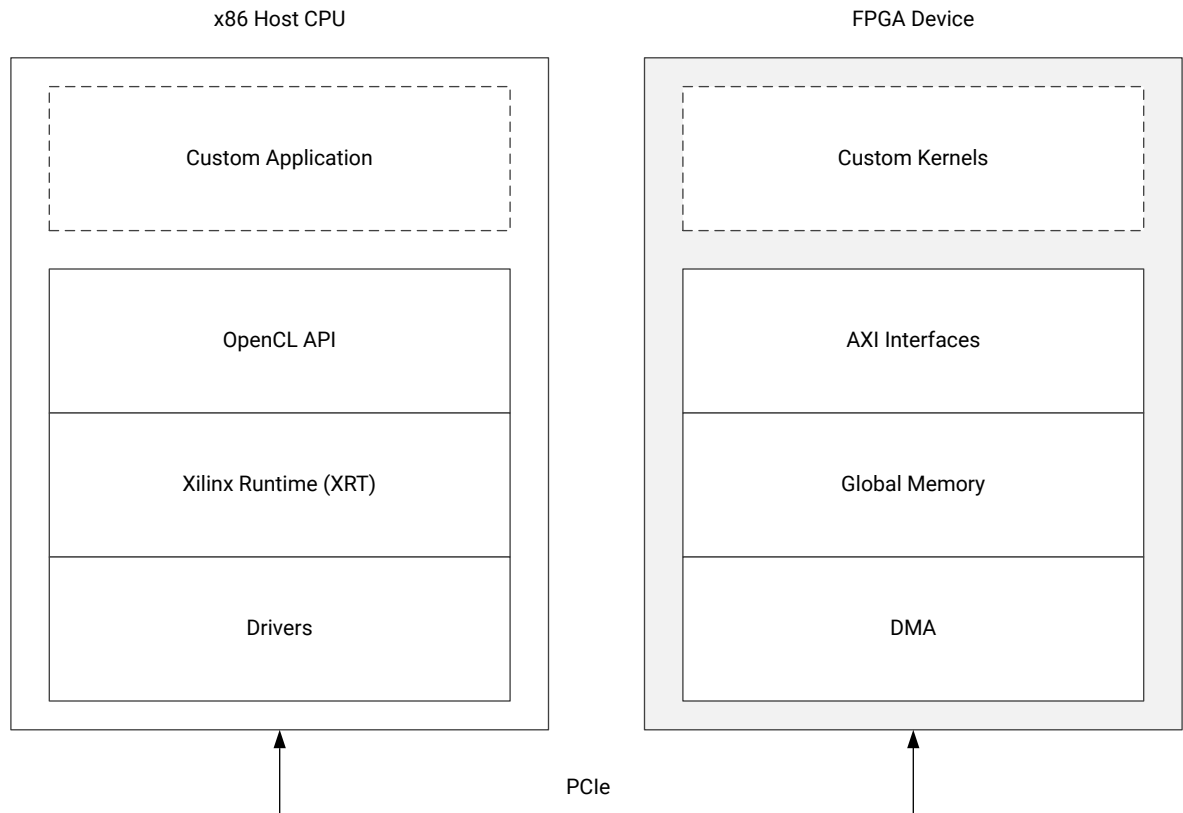
---

## SDAccel Execution Model

In the SDAccel framework, an application program is split between a host application and hardware accelerated kernels with a communication channel between them. The host application, written in C/C++ and using API abstractions like OpenCL, runs on an x86 server while hardware accelerated kernels run within the Xilinx FPGA. The API calls, managed by the Xilinx Runtime (XRT), are used to communicate with the hardware accelerators. Communication between the host x86 machine and the accelerator board, including control and data transfers, occurs across the PCIe bus. While control information is transferred between specific memory locations in hardware, global memory is used to transfer data between the host application and the kernels. Global memory is accessible by both the host processor and hardware accelerators, while host memory is only accessible by the host application.

For instance, in a typical application, the host will first transfer data, to be operated on by the kernel, from host memory into global memory. The kernel would subsequently operate on the data, storing results back to the global memory. Upon kernel completion, the host would transfer the results back into the host memory. Data transfers between the host and global memory introduce latency which can be costly to the overall acceleration. To achieve acceleration in a real system, the benefits achieved by hardware acceleration kernels must outweigh the extra latency of the data transfers. The general structure of this acceleration platform is shown in the following figure.

**Figure 1: Architecture of an SDAccel Application**



X21835-103118

The FPGA hardware platform, on the right-hand side, contains the hardware accelerated kernels, global memory along with the DMA for memory transfers. Kernels can have one or more global memory interfaces and are programmable. The SDAccel execution model can be broken down into these steps:

1. The host application writes the data needed by a kernel into the global memory of the attached device through the PCIe interface.
2. The host application sets up the kernel with its input parameters.
3. The host application triggers the execution of the kernel function on the FPGA.
4. The kernel performs the required computation while reading data from global memory, as necessary.
5. The kernel writes data back to global memory and notifies the host that it has completed its task.
6. The host application reads data back from global memory into the host memory and continues processing as needed.

The FPGA can accommodate multiple kernel instances at one time; this can occur between different types of kernels or multiple instances of the same kernel. The XRT transparently orchestrates the communication between the host application and the kernels in the accelerator. The number of instances of a kernel is determined by compilation options.

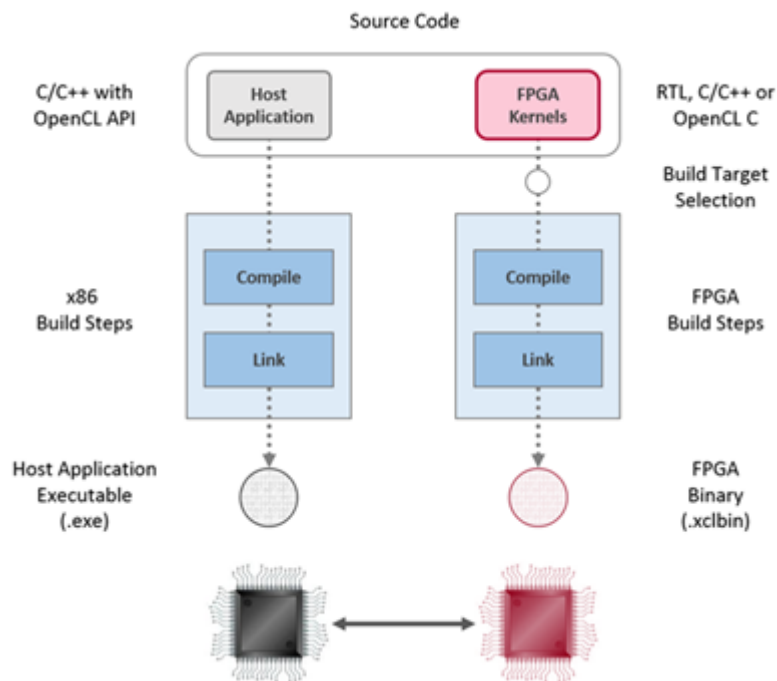
## SDAccel Build Process

The SDAccel environment offers all of the features of a standard software development environment:

- Optimized compiler for host applications
- Cross-compilers for the FPGA
- Robust debugging environment to help identify and resolve issues in the code
- Performance profilers to identify bottlenecks and optimize the code

Within this environment, the build process uses a standard compilation and linking process for both the software elements, and the hardware elements of the project. As shown in the following figure, the host application is built through one process using standard GCC compiler, and the FPGA binary is built through a separate process using the Xilinx `xocc` compiler.

Figure 2: Software/Hardware Build Process

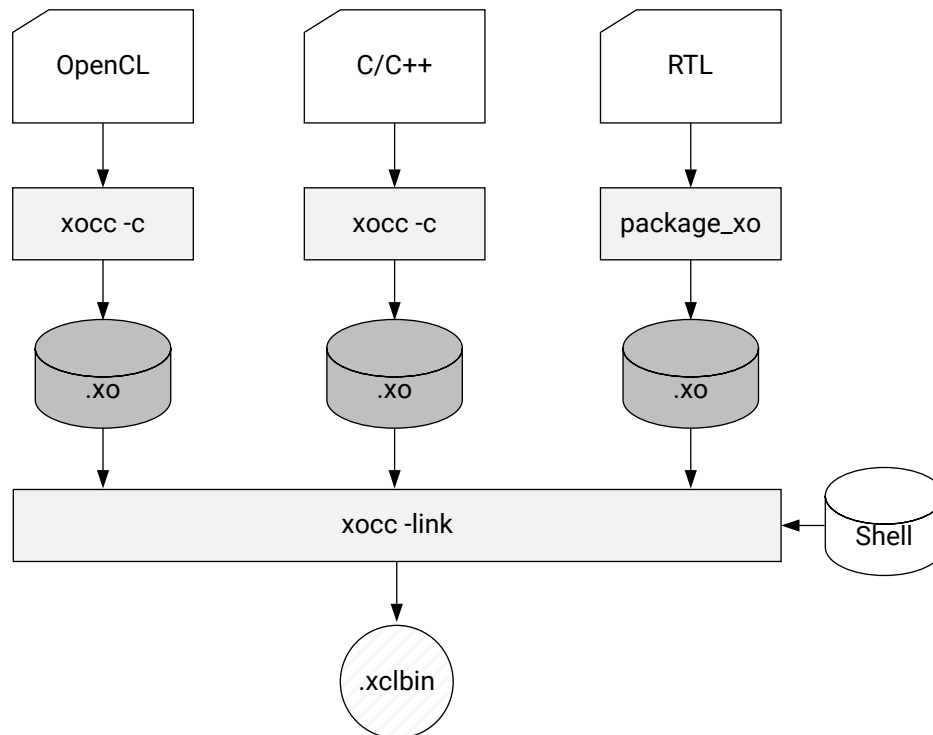


X22015-112618

1. Host application build process using GCC:
  - Each host application source file is compiled to an object file (.o).
  - The object files (.o) are linked with the Xilinx SDAccel runtime shared library to create the executable (.exe).
2. FPGA build process is highlighted in the following figure:
  - Each kernel is independently compiled to a Xilinx object (.xo) file.
    - C/C++ and OpenCL C kernels are compiled for implementation on an FPGA using the `xocc` compiler. This step leverages the Vivado® HLS compiler. Pragmas and attributes supported by Vivado HLS can be used in C/C++ and OpenCL C kernel source code to specify the desired kernel micro-architecture and control the result of the compilation process.
    - RTL kernels are compiled using the `package_xo` utility. The RTL kernel wizard in the SDAccel environment can be used to simplify this process.
  - The kernel .xo files are linked with the hardware platform (shell) to create the FPGA binary (.xclbin). Important architectural aspects are determined during the link step. In particular, this is where connections from kernel ports to global memory banks are established and where the number of instances for each kernel is specified.
    - When the build target is software or hardware emulation, as described below, `xocc` generates simulation models of the device contents.
    - When the build target is the system (actual hardware), `xocc` generates the FPGA binary for the device leveraging the Vivado Design Suite to run synthesis and implementation.



Figure 3: FPGA Build Process



X21155-111518

**Note:** The `xocc` compiler automatically uses the Vivado HLS and Vivado Design Suite tools to build the kernels to run on the FPGA platform. It uses these tools with predefined settings which have proven to provide good quality of results. Using the SDAccel environment and the `xocc` compiler does not require knowledge of these tools; however, hardware-savvy developers can fully leverage these tools and use all their available features to implement kernels.

## Build Targets

The SDAccel tool build process generates the host application executable (`.exe`) and the FPGA binary (`.xclbin`). The SDAccel build target defines the nature of FPGA binary generated by the build process.

The SDAccel tool provides three different build targets, two emulation targets used for debug and validation purposes, and the default hardware target used to generate the actual FPGA binary:

- Software Emulation (`sw_emu`):** Both the host application code and the kernel code are compiled to run on the x86 processor. This allows iterative algorithm refinement through fast build-and-run loops. This target is useful for identifying syntax errors, performing source-level debugging of the kernel code running together with application, and verifying the behavior of the system.

- **Hardware Emulation (hw\_emu):** The kernel code is compiled into a hardware model (RTL) which is run in a dedicated simulator. This build and run loop takes longer but provides a detailed, cycle-accurate, view of kernel activity. This target is useful for testing the functionality of the logic that will go in the FPGA and for getting initial performance estimates.
- **System (hw):** The kernel code is compiled into a hardware model (RTL) and is then implemented on the FPGA device, resulting in a binary that will run on the actual FPGA.

---

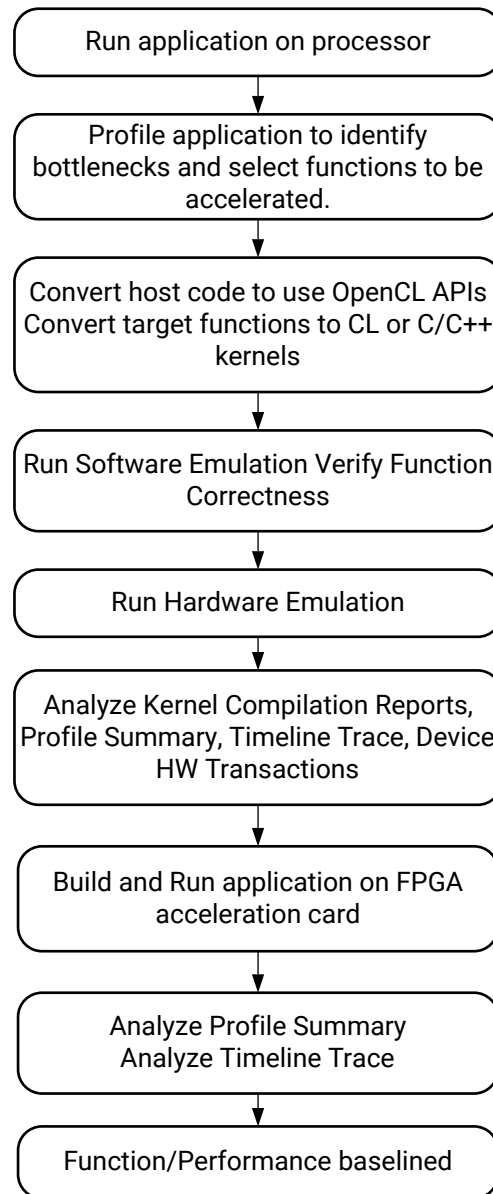
## SDAccel Optimization Flow Overview

SDAccel environment is a complete software development environment for creating, compiling, and optimizing C/C++/OpenCL applications to be accelerated on Xilinx FPGAs. The SDAccel environment includes the three recommended flows for optimizing an application. For information on each flow, refer to the following:

- [Baselining Functionalities and Performance](#)
- [Optimizing Data Movement](#)
- [Optimizing Kernel Computation](#)

### Baselining Functionalities and Performance

Before starting any optimization efforts, it is important to understand the performance of your application. This is achieved by baselining the application in terms of functionalities and performance.

Figure 4: **Baselining Functionalities and Performance Flow**


X22238-042419

## Identify Bottlenecks

The first step is to identify the bottlenecks of the current application running on your existing platform. The most effective way is to run the application with profiling tools, like `valgrind`, `callgrind`, and `GNU gprof`. The profiling data generated by these tools show the call graph with the number of calls to all functions and their execution time. The functions that consume the most execution time are good candidates to be offloaded and accelerated onto FPGAs.

## Convert Target Functions

After the target functions are selected, convert them to OpenCL C kernels or C/C++ kernels without any optimization. The application code calling these kernels will also need to be converted to use OpenCL APIs for data movement and task scheduling.



---

**TIP:** *Keep everything as simple as possible and minimize changes to the existing code, so you can quickly generate a working design on the FPGA and get the baselined performance and resource number.*

---

## Run Software and Hardware Emulation

Next, run software and hardware emulation to verify the function correctness and generate profiling data on the host code and the kernels. Analyze the kernel compilation reports, profile summary, timeline trace, and device hardware transactions to understand the baselined performance estimate such as timing, interval, and latency and resource utilization, such as DSP and block RAM.

## Build and Run the Application

The last step in baselining is building and running the application on an FPGA acceleration card. Analyze the reports from the system compilation and the profiling data from application execution to see the actual performance and resource utilization.



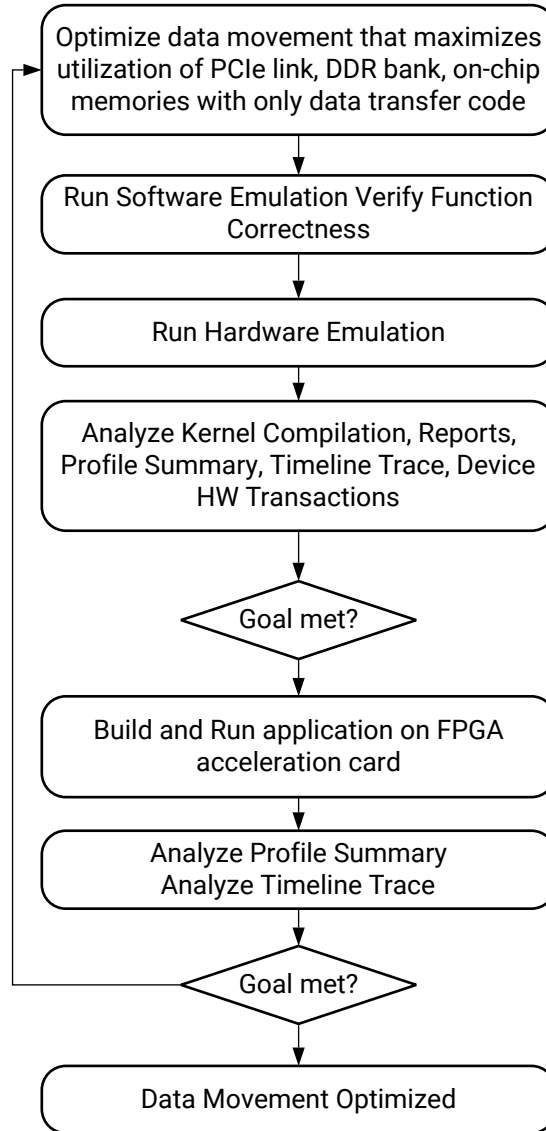
---

**TIP:** *Save all the reports during baselining, so you can reference and compare the results during optimization.*

---

# Optimizing Data Movement

Figure 5: Optimizing Data Movement Flow



X22239-042419

In the OpenCL API, all data is first transferred from the host memory to the global memory on the device and then from the global memory to the kernel for computation. The computation results are written back from the kernel to the global memory and then from the global memory to the host memory. A key factor in determining strategies for kernel optimization is understanding how data can be efficiently moved around.



**RECOMMENDED:** Before optimizing computation, first optimize the data movement in the application.

During data movement optimization, it is important to isolate data transfer code from computation code because inefficiency in computation might cause stalls in data movement.



---

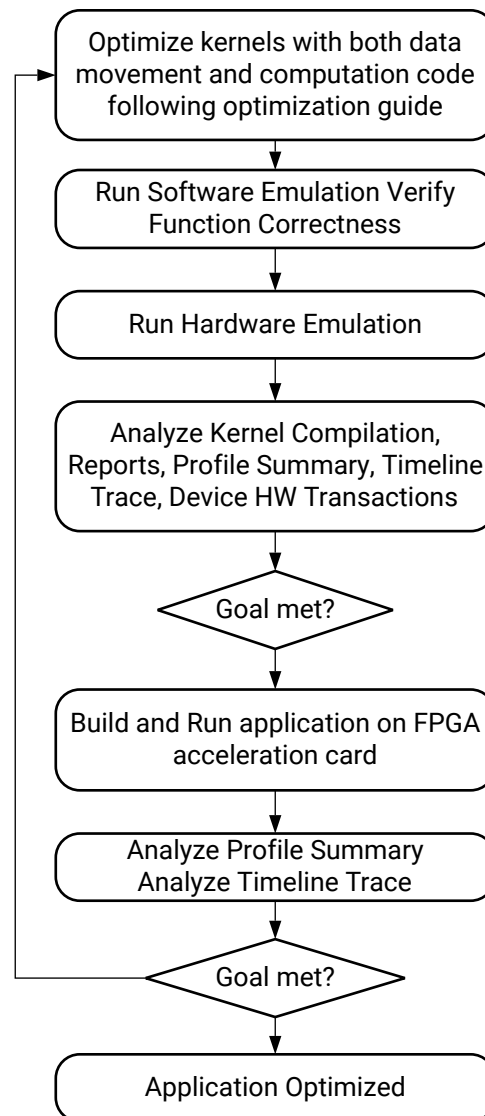
**RECOMMENDED:** *Xilinx recommends that you modify the host code and kernels with data transfer code only for this optimization step.*

---

The goal is to maximize the system level data throughput by maximizing PCIe bandwidth usage and DDR bandwidth usage. To achieve this goal, it usually takes multiple iterations of running software emulation, hardware emulation, as well as execution on FPGAs.

## Optimizing Kernel Computation

Figure 6: Optimizing Kernel Computation Flow



X22240-042419

One of the key benefits of an FPGA is that you can create custom logic for your specific application. The goal of kernel computation optimization is to create processing logic that can consume all the data as soon as they arrive at kernel interfaces. The key metric during this step is the initiation interval (II). This is generally achieved by expanding the processing code to match the data path with techniques, such as function pipelining, loop unrolling, array partitioning, data flowing, etc. The SDAccel environment produces various compilation reports and profiling data during hardware emulation and system run to assist your optimization effort. Refer to [Chapter 2: SDAccel Profiling and Optimization Features](#) for details on the compilation and profiling report.

# SDAccel Profiling and Optimization Features

The SDAccel™ environment generates various reports on the kernel resource and performance during compilation. It also collects profiling data during application execution in emulation mode and on the FPGA acceleration card. The reports and profiling data provide you with information on performance bottlenecks in the application and optimization techniques that can be used to improve performance. This chapter describes how to generate the reports and collect, display, and read the profiling results in the SDAccel environment.

---

## System Estimate

In the SDAccel development environment, generating FPGA binary files is the step with the longest execution time. The execution time is also most affected by the FPGA architecture and the number of compute units placed on the FPGA fabric. Therefore, it is essential for you to have a quicker way to understand the performance of the application before running it on the hardware, so you can spend more time iterating and optimizing your applications instead of waiting for the FPGA programming file to generate.

The system estimate in the SDAccel development environment takes into account the target hardware device and each compute unit in the application. Although an exact performance metric can only be measured by running the application on the FPGA, the estimation report in the development environment provides an accurate representation of the expected behavior.

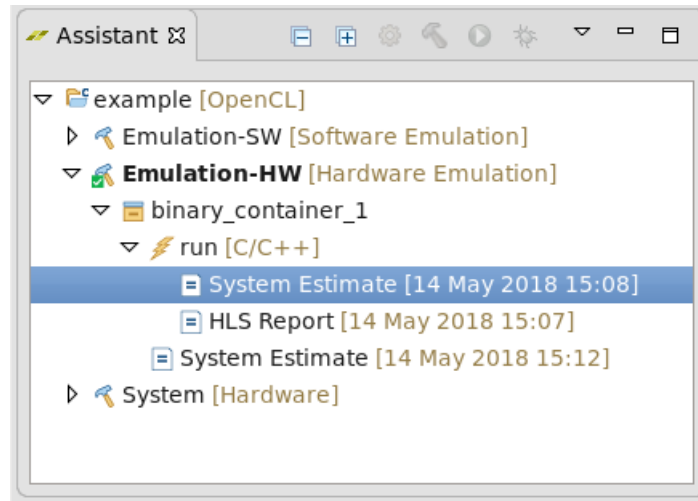
## GUI Flow

This report is automatically generated during the hardware emulation flow. There is one report generated for each kernel and a top report for the complete binary container. It is easy to access the reports from the Assistant window in the Emulation-HW folder.

The following figure shows the Assistant window with a System Estimate report for the `binary_container_1` and the kernel with the name `run`.



Figure 7: System Estimate Report in the Assistant Window



## Command Line

The following command generates the system performance estimate report `system_estimate.txt` for all kernels in `kernel.cl`:

```
xocc -c -t hw_emu --platform xilinx:adm-pcie-7v3:1ddr:3.0 --report estimate kernel.cl
```

The performance estimate report generated by the `xocc -report estimate` option provides information on every binary container in the application, as well as every compute unit in the design. The report is structured as follows:

- Target device information
- Summary of every kernel in the application
- Detailed information on every binary container in the solution

## Data Interpretation

The following example report file represents the information generated for the estimate report:

```
-----
Design Name:          _xocc_compile_kernel_bin.dir
Target Device:       xilinx:adm-pcie-ku3:2ddr-xpr:3.3
Target Clock:       200MHz
Total number of kernels: 1
-----

Kernel Summary
Kernel Name      Type  Target                OpenCL Library  Compute Units
-----
smithwaterman   clc   fpga0:OCL_REGION_0  xcl_xocc        1
-----
```

```

-----
OpenCL Binary:      xcl_xocc
Kernels mapped to: clc_region

Timing Information (MHz)
Compute Unit      Kernel Name      Module Name      Target Frequency
-----
smithwaterman_1  smithwaterman  smithwaterman    200

Estimated Frequency
-----
202.020203

Latency Information (clock cycles)
Compute Unit      Kernel Name      Module Name      Start Interval
-----
smithwaterman_1  smithwaterman  smithwaterman    29468

Best Case  Avg Case  Worst Case
-----
29467      29467      29467

Area Information
Compute Unit      Kernel Name      Module Name      FF      LUT      DSP      BRAM
-----
smithwaterman_1  smithwaterman  smithwaterman    2925    4304    1        10
-----
    
```

## Design and Target Device Summary

All design estimate reports begin with an application summary and information about the target device. The device information is provided in the following section of the report:

```

-----
Design Name:      _xocc_compile_kernel_bin.dir
Target Device:    xilinx:adm-pcie-ku3:2ddr-xpr:3.3
Target Clock:     200MHz
Total number of kernels: 1
-----
    
```

For the design summary, the only information that is provided is the design name and the selection of the target device. The other information provided in this section is the target board and the clock frequency.

- **Target Device:** The name of the board that runs the application compiled by the SDAccel development environment.
- **Target Clock:** Defines how fast the logic runs for compute units mapped to the FPGA fabric. Both of these parameters are fixed by the device developer.

These parameters cannot be modified from within the SDAccel environment.

## Kernel Summary

The Kernel Summary section lists all of the kernels defined for the current SDAccel solution. The following example shows the kernel summary:

```
Kernel Summary
Kernel Name      Type      Target                OpenCL Library  Compute Units
-----
smithwaterman    clc       fpga0:OCL_REGION_0   xcl_xocc        1
```

In addition to the kernel name, the summary also provides the execution target and type of the input source. Because there is a difference in compilation and optimization methodology for OpenCL™, C, and C/C++ source files, the type of kernel source file is specified.

The Kernel Summary section is the last summary information in the report. From here, detailed information on each compute unit binary container is presented.

## Timing Information

For each binary container, the detail section begins with the execution target of all compute units. It also provides timing information for every compute unit. As a general rule, if an estimated frequency is higher than that of the device target, the compute unit will be able to run in the device. If the estimated frequency is below the target frequency, the kernel code for the compute unit needs to be further optimized for the compute unit to run correctly on the FPGA fabric. This information is shown in the following example:

```
OpenCL Binary:      xcl_xocc
Kernels mapped to:  clc_region

Timing Information (MHz)
Compute Unit        Kernel Name      Module Name      Target Frequency
-----
smithwaterman_1    smithwaterman    smithwaterman    200

Estimated Frequency
-----
202.020203
```

It is important to understand the difference between the target and estimated frequencies. Compute units are not placed in isolation into the FPGA fabric. Compute units are placed as part of a valid FPGA design that can include other components defined by the device developer to support a class of applications.

Because the compute unit custom logic is generated one kernel at a time, an estimated frequency that is higher than the device target indicates to the developer using the SDAccel environment that there should not be any timing problems during the creation of the FPGA programming files.

## Latency Information

The latency information presents the execution profile of each compute unit in the binary container. When analyzing this data, it is important to keep in mind that all values are measured from the compute unit boundary through the custom logic. In-system latencies associated with data transfers to global memory are not reported as part of these values. Also, the latency numbers reported are only for compute units targeted at the FPGA fabric. The following is an example of the latency report:

```

Latency Information (clock cycles)
Compute Unit      Kernel Name      Module Name      Start Interval  Best Case
-----
smithwaterman_1  smithwaterman    smithwaterman    29468           29467

Avg Case  Worst Case
-----
29467     29467
    
```

The latency report is divided into the following fields:

- Start interval
- Best case latency
- Average case latency
- Worst case latency

The start interval defines the number of clock cycles that has to pass between invocations of a compute unit for a given kernel.

The best, average, and worst case latency numbers refer to how much time it takes the compute unit to generate the results of one ND Range data tile for the kernel. For cases where the kernel does not have data dependent computation loops, the latency values will be the same. Data dependent execution of loops introduces data specific latency variation that is captured by the latency report.

The interval or latency numbers will be reported as "undef" for kernels with one or more conditions listed below:

- OpenCL kernels that do not have explicit `reqd_work_group_size(x,y,z)`
- Kernels that have loops with variable bounds

**Note:** In case of undefined counters, consider using the TRIPCOUNT pragma.

**Note:** The latency information reflects estimates based on the analysis of the loop transformations and exploited parallelism of the model. These advanced transformations such as pipelining and data flow can heavily change the actual throughput numbers. Therefore, latency can only be used as relative guides between different runs.

## Area Information

There are a limited number of fundamental building blocks available in each FPGA. These fundamental blocks (FF, LUT, DSP, block RAM) are used by SDAccel development environment to generate the custom logic for each compute unit in the design. The number of each fundamental resource needed to implement the custom logic in a compute unit determines how many compute units can be simultaneously loaded into the FPGA fabric. The following example shows the area information reported for a compute unit:

| Area Information |               |               |      |      |     |      |
|------------------|---------------|---------------|------|------|-----|------|
| Compute Unit     | Kernel Name   | Module Name   | FF   | LUT  | DSP | BRAM |
| smithwaterman_1  | smithwaterman | smithwaterman | 2925 | 4304 | 1   | 10   |

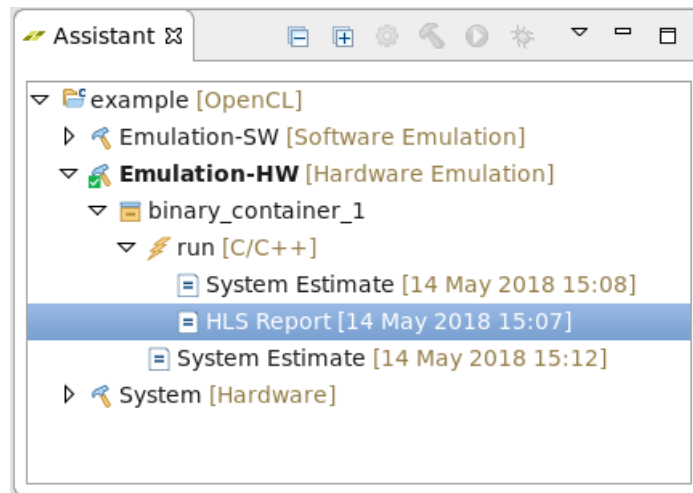
## HLS Report

After compiling a kernel using the SDx™ development environment GUI or the XOCC command line, the Vivado® High-Level Synthesis (HLS) tool HLS report is available. The HLS report includes details about the performance and logic usage of the custom-generated hardware logic from user kernel code. These details provide advanced users many insights into the kernel compilation results to guide kernel optimization.

## GUI Flow

After compiling a kernel using the SDx environment GUI, you can view the HLS Report in the Assistant window. The report is under the Emulation-HW or System build configuration, and has the <binary container> name, and the <kernel> name. This is illustrated in the following Assistant window:

Figure 8: Assistant Window



## Command Line

The HLS Report is designed to be viewed by the SDAccel environment GUI. However, for command line users, a textual representation of this report is also published. This report can be found inside the report directory situated under the kernel synthesis directory in the Vivado HLS tool solution directory.

Because the `xocc` command generates several additional levels of hierarchy above this synthesis directory, it is best to simply locate the file by name:

```
find . -name <module>_csynth.rpt
```

Where `<module>` is the name of the kernel.

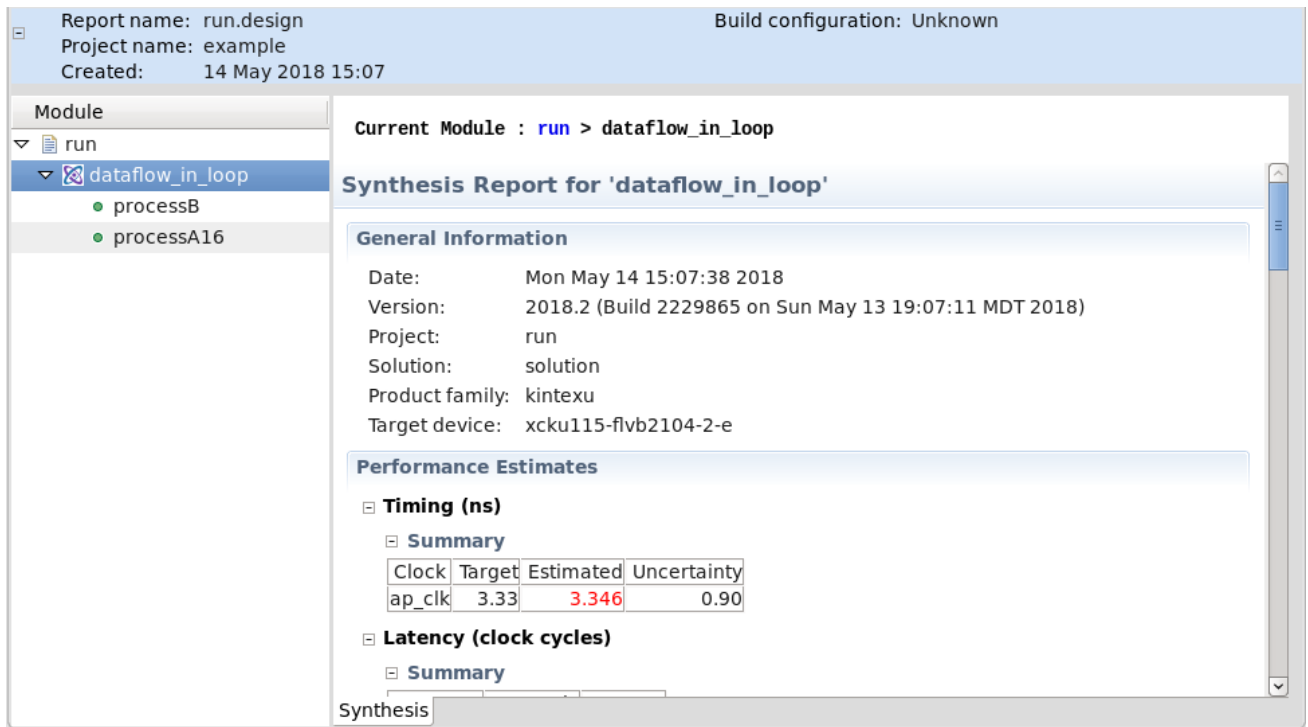
**Note:** The `find` command also supports the look up using wildcards such that the following command will look up all synthesis reports in any subdirectory:

```
find . -name "*_csynth.rpt"
```

## Data Interpretation

The left pane of the HLS Report shows the module hierarchy. Each module generated as part of the high level synthesis run is represented in this hierarchy. You can select any of these modules to present the synthesis details of the module in the right side of the Synthesis Report window.

Figure 9: HLS Report Window



Report name: run.design  
Project name: example  
Created: 14 May 2018 15:07

Build configuration: Unknown

Module: run > dataflow\_in\_loop

Current Module : run > dataflow\_in\_loop

### Synthesis Report for 'dataflow\_in\_loop'

**General Information**

Date: Mon May 14 15:07:38 2018  
Version: 2018.2 (Build 2229865 on Sun May 13 19:07:11 MDT 2018)  
Project: run  
Solution: solution  
Product family: kintexu  
Target device: xcku115-flvb2104-2-e

**Performance Estimates**

**Timing (ns)**

**Summary**

| Clock  | Target | Estimated | Uncertainty |
|--------|--------|-----------|-------------|
| ap_clk | 3.33   | 3.346     | 0.90        |

**Latency (clock cycles)**

**Summary**

Synthesis

The Synthesis Report is separated into several sections, namely:

- General Information
- Performance Estimates (timing and latency)
- Utilization Estimates
- Interface Information

If this information is part of a hierarchical block, it will sum up the information of the blocks contained in the hierarchy. Due to this fact, the hierarchy can also be navigated from within the report, when it is clear which instance contributes what to the overall design.



**CAUTION!** Regarding the absolute counts of cycles and latency, these numbers are based on estimates identified during synthesis, especially with advanced transformations, such as pipelining and dataflow; these numbers might not accurately reflect the final results. If you encounter question marks in the report, this might be due to variable bound loops, and you are encouraged to set trip counts for such loops to have some relative estimates presented in this report.

# Profile Summary Report

The SDAccel environment runtime automatically collects profiling data on host applications. After the application finishes execution, the profile summary is saved in HTML, CSV, and Google Protocol Buffer formats in the solution report directory or working directory. These reports can be reviewed in a web browser, spreadsheet viewer, or the integrated Profile Summary view in the SDAccel environment. The profile reports are generated in both SDAccel GUI and XOCC command line flows.

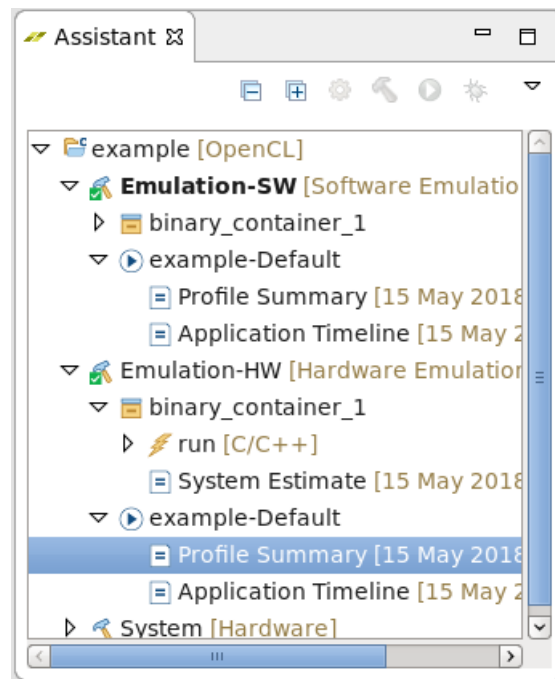
## GUI Flow

When you compile and execute an application from SDAccel environment, the profile summary is automatically generated.

To control the generation of profile information, simply edit the run configuration through the context menu of the build configuration, and select **Run → Run Configurations**.

After the configuration is run, the Assistant window enables easy access to the report from below the Run Configuration item. After the run configuration has executed, modifying the configuration can now be initiated directly through the context menu of the run configuration item in the Assistant window.

Figure 10: Profile Summary access in SDAccel GUI Flow



Double-click the report to open it.



## Command Line

Command line users execute standalone applications outside the SDAccel environment. To generate the profile summary data, you can compile your design without any additional options. However, linking the bitstream file (`xclbin`) requires the `--profile_kernel` option.

The argument provided through the `--profile_kernel` option can be used to limit data collection, which might be required in large systems. The general syntax for the `profile_kernel` option with respect to the profile summary report is:

```
--profile_kernel <[data]:<[kernel_name|all]:[compute_unit_name|all]:  
[interface_name|all]:[counters|all]>
```

The following three fields can be specified to determine which interface the performance monitor is applied to:

- `kernel_name`
- `compute_unit_name`
- `interface_name`

However, you can also specify the keyword `all` to apply the monitoring to all existing kernels, compute units, and interfaces with a single option. The last option, `<counters|all>` allows you to restrict the information gathering to just `counters` for large designs, while `all` (default) will include the collection of actual trace information.

**Note:** The `profile_kernel` option is additive and can be used multiple times on the link line.

If the `profile = true` option is specified in the `sdaccel.ini` file, when the program is executed, a `profile_summary.csv` file is created.

```
[Debug]  
profile = true
```

The `.csv` file needs to be manually converted to Google Protocol Buffer format (`.xprf`) before the profiling result can be viewed in the integrated Profile Summary view. The following command line example generates an `.xprf` file from the `.csv` input file:

```
sdx_analyze profile profile_summary.csv
```

### Display the Profile Summary

Use the following methods to display the SDAccel environment Profile Summary view created from the command line.

1. To display the report in your web browser of choice, do the following:
  - a. Run the following command:

```
sdx_analyze profile -i profile_summary.csv -f html
```

This creates an HTML file representing the data. that can be opened by the web browser of your choice. The file contains the same profiling result as presented in [GUI Flow](#).

- b. Navigate to the file location, and double-click the generated HTML file.
2. To display the report in the integrated SDAccel Profile Summary view, do the following:
  - a. Use the following command to convert the .csv data file into the protobuf format.

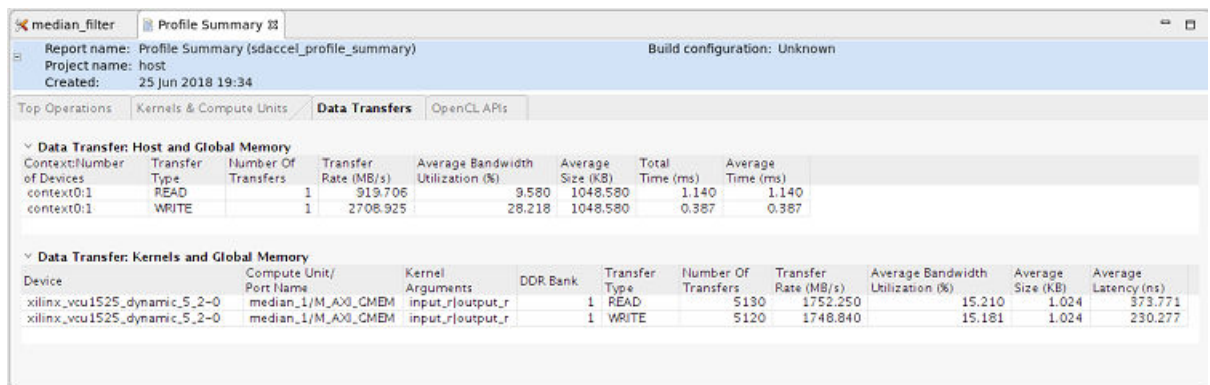
```
sdx_analyze profile -i profile_summary.csv -f protobuf
```

- b. Start SDAccel tool GUI by running the `sdx` command:

```
$ sdx
```

- c. Choose the default workspace when prompted.
- d. Select **File** → **Open File**.
- e. Browse to and then open the .xprf file created by the `sdx_analyze` command run in step a.

The following figure shows the Profile Summary view that displays OpenCL API calls, kernel executions, data transfers, and profile rule checks (PRCs).



## Data Interpretation

The profile summary includes a number of useful statistics for your OpenCL application. This can provide you with a general idea of the functional bottlenecks in your application. The profile summary consists of the following sections:

- **Top Operations:**
  - **Top Data Transfer: Kernels and Global Memory:** This table displays the profile data for top data transfers between FPGA and device memory.
    - **Device:** Name of device
    - **Compute Unit:** Name of compute unit

- **Number of Transfers:** Sum of write and read AXI transactions monitored on device
- **Average Bytes per Transfer:**  $(\text{Total Read Bytes} + \text{Total Write Bytes}) / (\text{Total Read AXI Transactions} + \text{Total Write AXI Transactions})$
- **Transfer Efficiency (%):**  $(\text{Average Bytes per Transfer}) / \min(4\text{K}, (\text{Memory Bit Width}/8 * 256))$   
 AXI4 specification limits the max burst length to 256 and max burst size to 4K bytes.
- **Total Data Transfer (MB):**  $(\text{Total Read Bytes} + \text{Total Write Bytes}) / 1.0\text{e}6$
- **Total Write (MB):**  $(\text{Total Write Bytes}) / 1.0\text{e}6$
- **Total Read (MB):**  $(\text{Total Read Bytes}) / 1.0\text{e}6$
- **Transfer Rate (MB/s):**  $(\text{Total Data Transfer}) / (\text{Compute Unit Total Time})$
- **Top Kernel Execution:**
  - **Kernel Instance Address:** Host address of kernel instance (in hex)
  - **Kernel:** Name of kernel
  - **Context ID:** Context ID on host
  - **Command Queue ID:** Command queue ID on host
  - **Device:** Name of device where kernel was executed (format: <device>-<ID>)
  - **Start Time (ms):** Start time of execution (in ms)
  - **Duration (ms):** Duration of execution (in ms)
  - **Global Work Size:** NDRange of kernel
  - **Local Work Size:** Work group size of kernel
- **Top Memory Writes: Host and Device Global Memory:**
  - **Buffer Address:** Host address of buffer (in hex)
  - **Context ID:** Context ID on host
  - **Command Queue ID:** Command queue ID on host
  - **Start Time (ms):** Start time of write transfer (in ms)
  - **Duration (ms):** Duration of write transfer (in ms)
  - **Buffer Size (KB):** Size of write transfer (in KB)
  - **Writing Rate (MB/s):**  $\text{Writing Rate} = (\text{Buffer Size}) / (\text{Duration})$

- **Top Memory Reads: Host and Device Global Memory:**
  - **Buffer Address:** Host address of buffer (in hex)
  - **Context ID:** Context ID on host
  - **Command Queue ID:** Command queue ID on host
  - **Start Time (ms):** Start time of read transfer (in ms)
  - **Duration (ms):** Duration of read transfer (in ms)
  - **Buffer Size (KB):** Size of read transfer (in KB)
  - **Reading Rate (MB/s):** Reading Rate = (Buffer Size) / (Duration)
- **Kernels & Compute Units:**
  - **Kernel Execution (includes estimated device times):** This table displays the profile data summary for all kernel functions scheduled and executed.
    - **Kernel:** Name of kernel
    - **Number of Enqueues:** Number of times kernel is enqueued
    - **Total Time (ms):** Sum of runtimes of all enqueues (measured from START to END in the OpenCL API execution model)
    - **Minimum Time (ms):** Minimum runtime of all enqueues
    - **Total Time (ms):** Sum of runtimes of all enqueues (measured from START to END in the OpenCL API execution model)
    - **Average Time (ms):** (Total Time) / (Number of Enqueues)
    - **Maximum Time (ms):** Maximum runtime of all enqueues
  - **Compute Unit Utilization (includes estimated device times):** This table displays the summary profile data for all compute units on the FPGA.
    - **Device:** Name of device (format: <device>-<ID>)
    - **Compute Unit:** Name of Compute Unit
    - **Kernel:** Kernel this Compute Unit is associated with
    - **Global Work Size:** NDRange of kernel (format is x:y:z)
    - **Local Work Size:** Local work group size (format is x:y:z)
    - **Number of Calls:** Number of times the Compute Unit is called
    - **Dataflow Execution:** Indicates if top level dataflow execution is enabled

- **Maximum Overlapping Executions:** How much executions were actually operating in parallel at some point during execution
- **Dataflow Acceleration:** Estimated improvement due to dataflow acceleration
- **Total Time (ms):** Sum of runtimes of all calls
- **Minimum Time (ms):** Minimum runtime of all calls
- **Average Time (ms):** (Total Time) / (Number of Work Groups)
- **Maximum Time (ms):** Maximum runtime of all calls
- **Clock Frequency (MHz):** Clock frequency used for a given accelerator (in MHz)
- **Data Transfers:**
  - **Data Transfer: Host and Global Memory:** This table displays the profile data for all read and write transfers between the host and device memory through the PCI Express® link.
  - **Context: Number of Devices:** Context ID and number of devices in context
  - **Transfer Type:** READ or WRITE
  - **Number of Transfers:** Number of host data transfers  
*Note:* Might contain `printf` transfers
  - **Transfer Rate (MB/s):** (Total Bytes Sent) / (Total Time in usec), where Total Time includes software overhead
  - **Average Bandwidth Utilization (%):** (Transfer Rate) / (Max. Transfer Rate), where Max. Transfer Rate = (256/8 bytes) \* (300 MHz) = 9.6 GBps
  - **Average Size (KB):** (Total KB sent) / (number of transfers)
  - **Total Time (ms):** Sum of transfer times
  - **Average Time (ms):** (Total Time) / (number of transfers)
  - **Data Transfer: Kernels and Global Memory:** This table displays the profile data for all read and write transfers between the FPGA and device memory.
    - **Device:** Name of device
    - **Compute Unit/Port Name:** <Name of Compute Unit>/<Name of Port>
    - **Kernel Arguments:** List of arguments connected to this port
    - **DDR Bank:** DDR bank number this port is connected to
    - **Transfer Type:** READ or WRITE
    - **Number of Transfers:** Number of AXI transactions monitored on device

**Note:** Might contain `printf` transfers)

- **Transfer Rate (MB/s):** (Total Bytes Sent) / (Compute Unit Total Time)
  - **Compute Unit Total Time:** Total execution time of compute unit
  - **Total Bytes Sent:** Sum of bytes across all transactions
- **Average Bandwidth Utilization (%):** (Transfer Rate) / (0.6 \* Max. Transfer Rate), where Max. Transfer Rate = (512/8 bytes) \* (300 MHz) = 19200 MBps
- **Average Size (KB):** (Total KB sent) / (number of AXI transactions)
- **Average Latency (ns):** (Total latency of all transaction) / (number of AXI transactions)
- **OpenCL API Calls:** This table displays the profile data for all OpenCL host API function calls executed in the host application.
  - **API Name:** Name of API function (for example, `clCreateProgramWithBinary`, `clEnqueueNDRangeKernel`)
  - **Number of Calls:** Number of calls to this API
  - **Total Time (ms):** Sum of runtimes of all calls
  - **Minimum Time (ms):** Minimum runtime of all calls
  - **Average Time (ms):** (Total Time) / (Number of Calls)
  - **Maximum Time (ms):** Maximum runtime of all calls

## Application Timeline

The Application Timeline window collects and displays host and device events on a common timeline to help you understand and visualize the overall health and performance of your systems. These events include:

- OpenCL API calls from the host code.
- Device trace data including AXI transaction start/stop, kernel start/stop, etc.

While useful for debugging and profiling the application, timeline and device trace data are not collected by default because the runtime needs to periodically unload the trace data from the FPGA, which can add additional time to the overall application execution. However, the device data is collected with dedicated hardware inside the FPGA, so the data collection does not affect kernel functionality on the FPGA. The following sections describe setups required to enable time and device data collection.

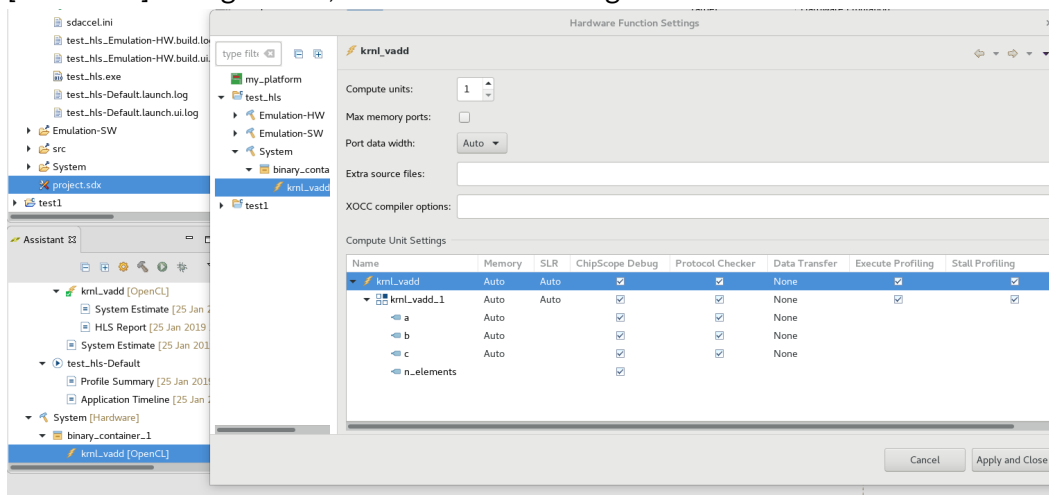
Turning on device profiling is intrusive and can negatively affect overall performance. This feature should be used for system performance debugging only.

**Note:** Device profiling can be used in Emulation-HW without negative impact.

## GUI Flow

Timeline and device trace data collection is part of run configuration for an SDAccel™ project created from the integrated SDAccel environment. Follow the steps below to enable it:

1. Instrumenting the code is required for System execution. This is done through the Hardware Function Settings dialog box. In the Assistant window, right-click the kernel under the System [Hardware] configuration, and select the Settings Command.



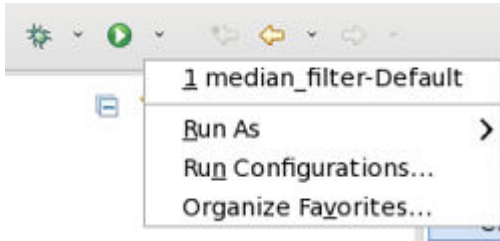
With respect to application timeline functionality, you can enable Data Transfer, Execute Profile, and Stall Profiling. These options are instrumenting all ports of each instance of any kernel. As these options insert additional hardware, instrumenting all ports might be too much. Towards that end, more control is available through command line options as detailed in the [Command Line](#) section. These options are only valid for system runs. During hardware emulation, this data is generated by default.

- **Data Transfer:** This option enables monitoring of data ports.
- **Execute Profiling:** This option provides minimum port data collection during system run. This option records the execution times of the compute units. Execute profiling is enabled by default for data and stall profiling.
- **Stall Profiling:** This option includes the stall monitoring logic in the bitstream.

2. Specify what information is actually going to be reported during a run.

**Note:** Only information actually exposed from the hardware during system execution is reported.

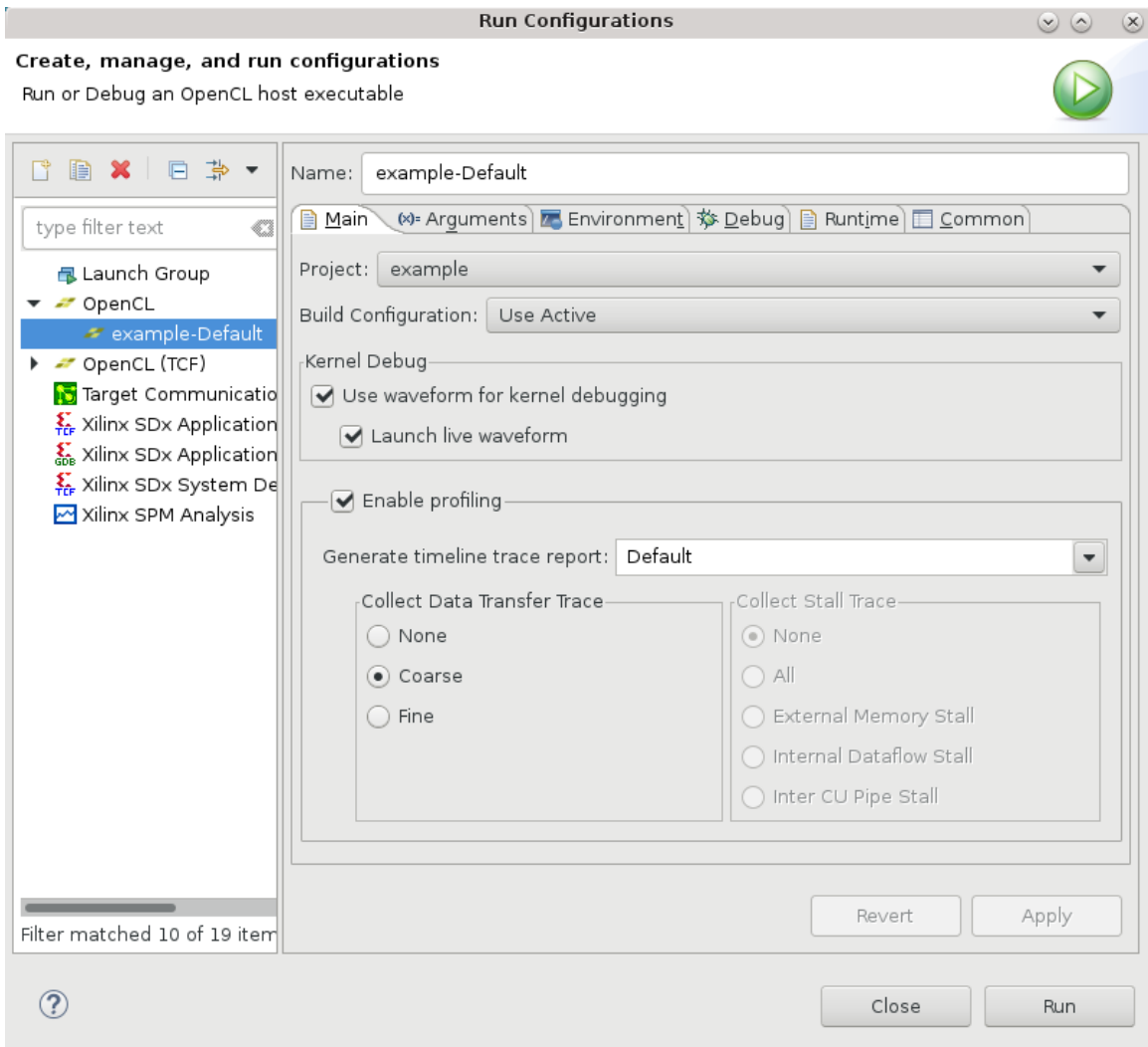
To configure reporting, click the down arrow next to the Debug or Run button, and then select **Run Configurations** to open the Run Configurations window.



3. In the Run Configurations window, click the Profile tab.

Ensure the Enable profiling check box is selected. This enables basic profiling support. With respect to trace data, ensure that Generate timeline trace report actually gathers the information in the build config you are running.

Default implies that no trace data capturing is supported in system execution, but enabled by default in hardware emulation.





You can also select the amount of information to gather during runtime. Select the granularity for trace data collection independently for Data Transfer Trace and Stall Trace.

The Data Transfer Trace options are as follows:

- **Coarse:** Show compute unit transfer activity from beginning of first transfer to end of last transfer (before compute unit transfer ends).
- **Fine:** Show all AXI-level burst data transfers.
- **None:** Turn off reading and reporting of device-level trace during runtime.

The Stall Trace Options are as follows:

- **None:** Turn off any stall trace information gathering.
- **All:** Record all stall trace information.
- **External Memory Stall:** Memory stalls to DDR (for example, AXI4 read from DDR).
- **Internal Dataflow Stall:** Intra-kernel streams (for example, writing to a full FIFO between data flow blocks).
- **Inter CU Pipe Stall:** Inter-kernel pipe (for example, writing to a full OpenCL™ pipe between kernels).

If you have multiple run configurations for the same project, you must change the profile settings for each run configuration.

4. After running configurations, in the Assistant window, double-click **Application Timeline** to open the Application Timeline window.

## Command Line

Complete the following steps to enable timeline and device trace data collection in the command line flow:

1. This step is responsible for the FPGA bitstream instrumentation with SDx Accel Monitors (SAM) and SDx Performance Monitors (SPMs). The instrumentation is performed through the `--profile_kernel`, which has three distinct instrumentation options (`data`, `stall`, and `exec`).

**Note:** The `--profile_kernel` option is ignored except for system compilation and linking. During hardware emulation, this data is generated by default.

The `--profile_kernel` option has three fields that are required to determine the specific kernel interface to which the monitors are applied. However, if resource usage is not an issue, the keyword `all` enables you to apply the monitoring to all existing kernels, compute units, and interfaces with a single option. Otherwise, you can specify the `kernel_name`, `compute_unit_name`, and `interface_name` explicitly to limit instrumentation. The last option, `<counters|all>` allows you to restrict the information gathering to just `counters` for large designs, while `all` (default) includes the collection of actual trace information.

**Note:** The `--profile_kernel` option is additive and can be used multiple times on the link line.

- **data:** This option enables monitoring of data ports through SAM and SPM IPs. This option needs to be set only during linking.

```
-l --profile_kernel <[data]:<[kernel_name|all]:[compute_unit_name|all]:[interface_name|all]:[counters|all]>
```

- **stall:** This option needs to be applied during compilation:

```
-c --profile_kernel <[stall]:<[kernel_name|all]:[compute_unit_name|all]:[counters|all]>
```

and during linking:

```
-l --profile_kernel <[stall]:<[kernel_name|all]:[compute_unit_name|all]:[counters|all]>
```

This option includes the stall monitoring logic (using SAM IP) in the bitstream. However, it does require that stall ports are present on the kernel interface. To facilitate this, the option is required for compilation of the C/C++/OpenCL API kernel modules.

- **exec:** This option provides minimum port data collection during system run. It simply records the execution times of the kernel through the use of SAM IP. This feature is by default enabled on any port that uses the data or stall data collection. This option needs to be provided only during linking.

```
-l --profile_kernel <[exec]:<[kernel_name|all]:[compute_unit_name|all]>:[counters|all]
```

2. After the kernels are instrumented, data gathering must be enabled during runtime execution. Do this through the use of the `sdaccel.ini` file that is in the same directory as the host executable. The following `sdaccel.ini` file will enable maximum information gathering during runtime:

```
[Debug]
profile=true
timeline_trace=true
data_transfer_trace=coarse
stall_trace=all
```

- **profile=<true|false>:** When this option is specified as true, basic profile monitoring is enabled. Without any additional options, this implies that the host runtime logging profile summary is enabled. However, without this option enabled, no monitoring is performed at all.
- **timeline\_trace=<true|false>:** This option will enable timeline trace information gathering of the data. Without adding profile IP into the FPGA (data), it will only show host information. At a minimum, to get more compute unit start and end execution times in the timeline trace, the compute unit needs to be linked with `--profile_kernel exec`.
- **data\_transfer\_trace=<coarse|fine|off>:** This option enables device-level AXI data transfers trace:
  - **coarse:** Show compute unit transfer activity from beginning of first transfer to end of last transfer (before compute unit transfer ends).

- `fine`: Show all AXI-level burst data transfers.
  - `off`: Turn off reading and reporting of device-level trace during runtime.
  - `stall_trace=<dataflow|memory|pipe|all|off>`: Specify what types of stalls to capture and report in timeline trace. The default is `off`.
    - `off`: Turn off any stall trace information gathering.
 

**Note:** Enabling stall tracing can often fill the trace buffer, which results in incomplete and potentially corrupt timeline traces. This can be avoided by setting `trace_stall=off`.
    - `all`: Record all stall trace information.
    - `dataflow`: Intra-kernel streams (for example, writing to full FIFO between dataflow blocks).
    - `memory`: External memory stalls (for example, AXI4 read from the DDR).
    - `pipe`: Inter-kernel pipe (for example, writing to full OpenCL pipe between kernels).
3. To allow User Function analysis in the Application Timeline it is necessary to run the waveform tool as part of the Emulation flow. This requires to launch the waveform tool through the `sdaccel.ini` file.

```
[Emulation]
launch_waveform=batch
```

- `launch_waveform=<gui|batch>`: This option automatically starts the waveform tool during emulation. See [Waveform View](#) for more details.
    - `gui`: Start the graphical user interface for the live waveform view
    - `batch`: Start the waveform processing as a background process.
4. In command line mode, CSV files are generated to capture the trace data. These CSV reports need to be converted to the Application Timeline format using the `sdx_analyze` utility before they can be opened and displayed in the SDAccel environment GUI.

```
sdx_analyze trace timeline_trace.csv -k timeline_kernels.csv -f wdb
```

This creates the `timeline_trace.wdb` file by default, which can be opened from the GUI. The `timeline_kernels.csv` file contains specific kernel trace data which might not always be available. In this case, the option `-k timeline_kernels.csv` should be omitted from the `sdx_analyze` command.

5. To view the timeline report host and device waveforms, do the following:
- a. Start the SDx environment by running the command:

```
$sdx
```

- b. Choose a workspace when prompted.
- c. Select **File** → **Open File**, browse to the `.wdb` file generated during hardware emulation or system run, and open it.

## Data Interpretation

The following figure shows the Application Timeline window that displays host and device events on a common timeline. This information helps you to understand details of application execution and identify potential areas for improvements.

Figure 11: Application Timeline Window



The Application Timeline window trace includes two main sections:

- **Host:** Shows the trace of all the activity originating from the host side.
- **Device:** Shows the activity of the compute-units on the FPGA.

Under the host, different activities are categorized as OpenCL™ API calls, Data Transfer, and Kernels.

The complete tree has the following structure:

- **Host:**
  - **OpenCL API Calls:** All OpenCL API calls are traced here. The activity time is measured from the host perspective.
    - **General:** All general OpenCL API calls, such as `clCreateProgramWithBinary()`, `clCreateContext()`, and `clCreateCommandQueue`, are traced here.

- **Queue:** OpenCL API calls that are associated with a specific command queue are traced here. This includes commands, such as `clEnqueueMigrateMemObjects`, `clEnqueueNDRangeKernel`, and so on. If the user application creates multiple command queues, then this section shows as many queues and activities under it.
- **Data Transfer:** In this section, the DMA transfers are traced. The data transfer from the host to the device appear under **Write**, and the transfers from device to host appear under **Read**. The additional section **Copy** traces direct communication between kernels.
- **Kernel Enqueues:** The active kernel executions are shown here. The kernels here should not be confused with your kernels/compute-unit on the device. In this instance, kernels refers to the NDRangeKernels and the tasks created by the `clEnqueueNDRangeKernels()` and `clEnqueueTask()` APIs, and these are plotted against the time measured from the host's perspective. Multiple kernels can be scheduled to be executed at the same time and they are traced from the point they are scheduled to run until the end of kernel execution. This is the reason for multiple entries. The number of rows depend on the number of overlapping kernel executions.

**Note:** Overlapping of the kernels should not be mistaken for actual real parallel execution on the device as the process might not be ready to actually execute right away.

- **Device "name":**
- **Binary Container "name":**
- **Accelerator "name":** This is the name of the compute unit (also known as accelerator) on the FPGA.
- **User Functions:** In the case of the Vivado HLS tool kernels, functions that are implemented as data flow processes are traced here. The trace for these functions show the number of active instances of these functions that are currently executing in parallel. These names are generated in hw emulation when waveform is enabled.

**Note:** Function level activity is only possible in Hardware Emulation.

- **Function: "name a"**
- **Function: "name b"**
- **Read:** A compute unit reads from the DDR over AXI-MM ports. The trace of data a read by a compute unit is shown here. The activity is shown as transaction and the tooltip for each transaction shows more details of the AXI transaction. These names are generated when `--profile_kernel data` is used.
  - **m\_axi\_<bundle name>(port)**
- **Write:** A compute unit writes to the DDR over AXI-MM ports. The trace of data written by a compute unit is shown here. The activity is shown as transactions and the tool-tip for each transaction shows more details of the AXI transaction. This is generated when `--profile_kernel data` is used.

- m\_axi\_<bundle name>(port)

## Waveform View

The SDx development environment can generate a waveform view and launch a live waveform viewer when running hardware emulation. It displays in-depth details on the emulation results at system level, compute unit level, and at function level. The details include data transfers between the kernel and global memory, data flow through inter-kernel pipes as well as data flow through intra-kernel pipes. They provide many insights into the performance bottleneck from the system level down to individual function call to help developers optimize their applications.

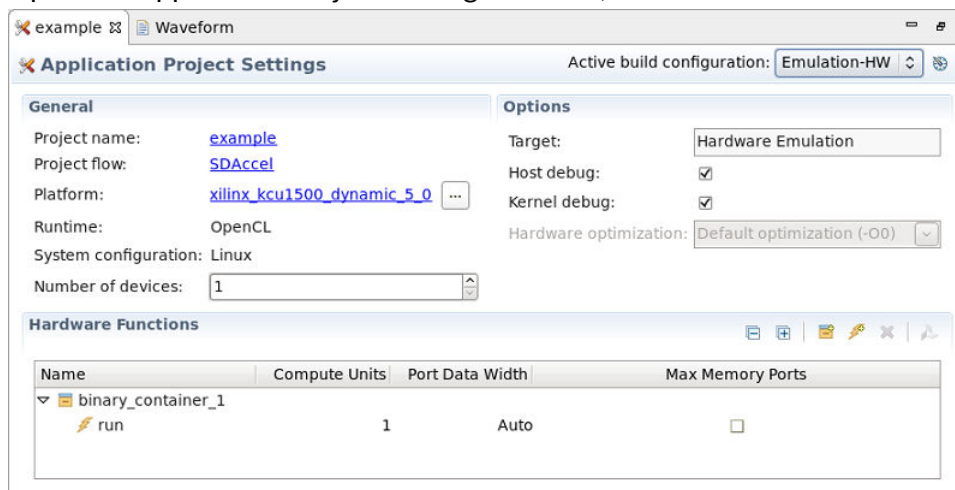
By default, the waveform view and live waveform viewer are not enabled. This is because the views require that the runtime generates a simulation waveform during hardware emulation, which consumes more time and disk space. The following sections describe the setup required to enable data collection.

**Note:** The waveform view allows you to look directly at the device transactions from within the SDx development environment. In contrast, the live waveform capability actually spawns the simulation waveform view that visualizes the hardware transactions in addition to potential user-selected internal signals.

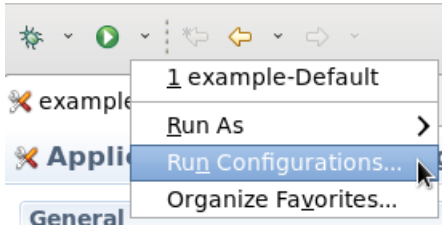
## GUI Flow

Follow the steps below to enable waveform data collection and to open the viewer:

1. Open the Application Project Settings window, and select the **Kernel debug** check box.



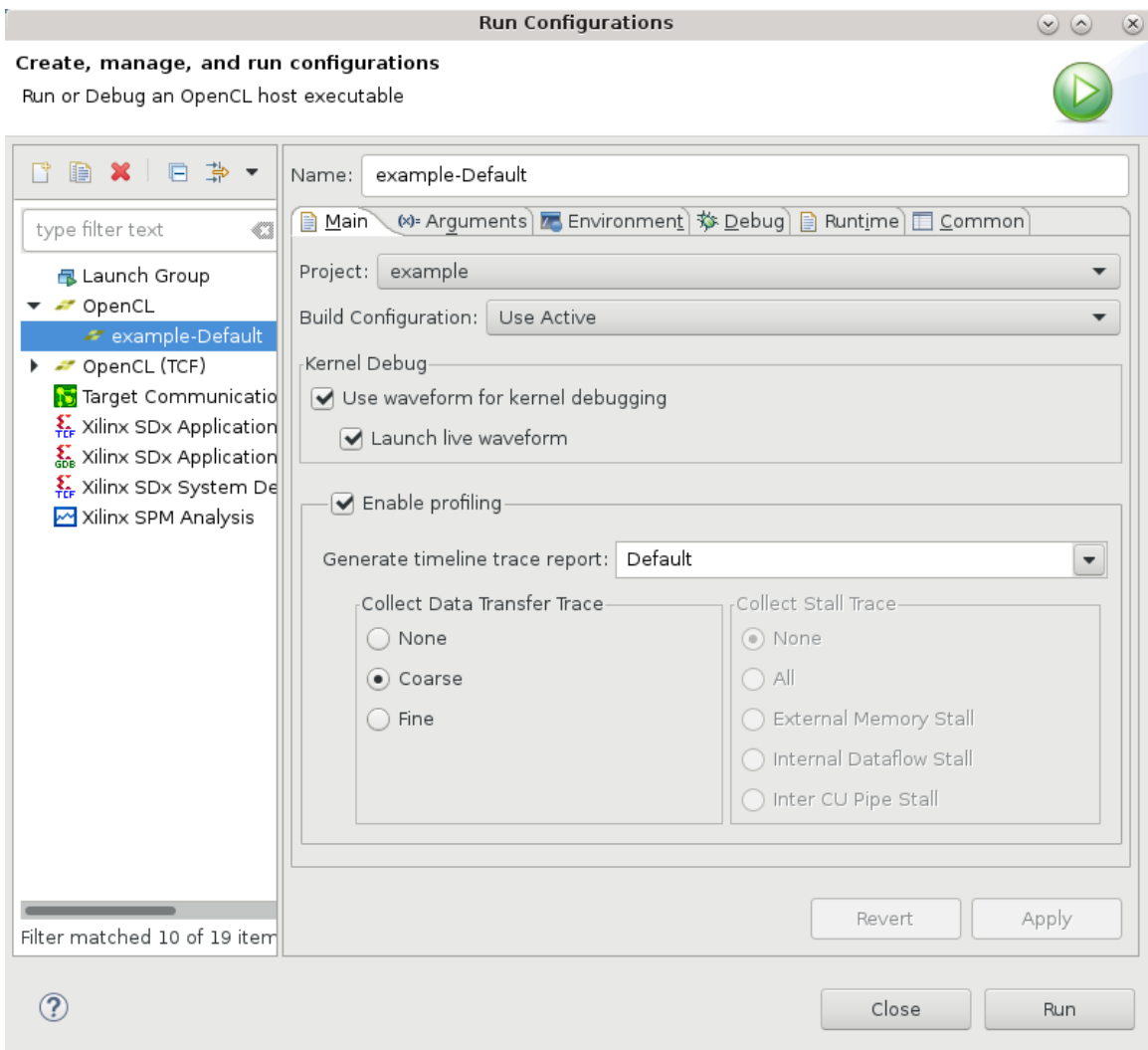
2. Click the down arrow next to the Run button, and select **Run Configurations** to open the Run Configurations window.



3. On the Run Configurations window, click **Main**, and select the **Use waveform for kernel debugging** check box.

Optionally:

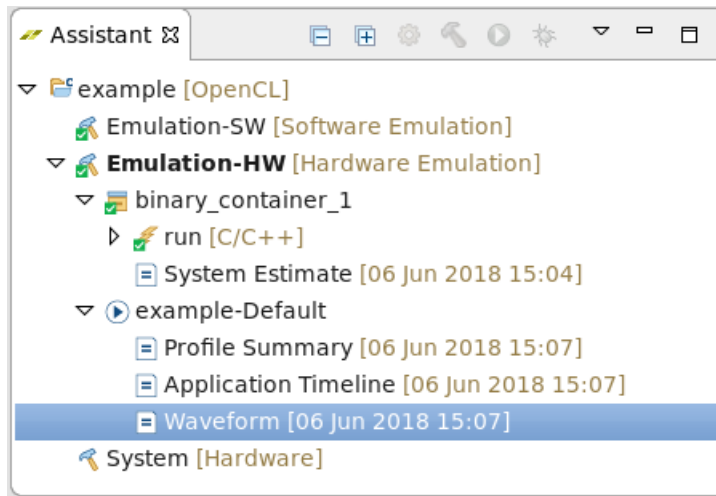
- To bring up the Simulation window to view the Live Waveform while the hardware emulation is running, deselect **Launch live waveform**.
- To enable basic profiling, select **Enable profiling**.



4. If you have multiple run configurations for the same project, change the profile settings for each run configuration.

- If you have not selected the Live Waveform viewer to be launched automatically, open the Waveform view from the SDx Development Environment.

In the SDx Development Environment, double-click **Waveform** in the Assistant window to open the Waveform view window.



## Command Line

Use the following instructions to enable waveform data collection from the Command Line during hardware emulation and open the viewer:

- Turn on debug code generation during kernel compilation.

```
xocc -g -t hw_emu ...
```

- Create an `sdaccel.ini` file in the same directory as the host executable with the contents below:

```
[Debug]
profile=true
timeline_trace=true
```

This enables maximum observability. The options in detail are:

- profile=<true|false>:** Setting this option to true, enables profile monitoring. Without any additional options, this implies that the host runtime logging profile summary is enabled. However, without this option enabled, no monitoring is performed at all.
  - timeline\_trace=<true|false>:** This option enables timeline trace information gathering of the data.
- To see the live waveform and additional simulation waveforms, add the following to the emulation section in the `sdaccel.ini`:

```
[Emulation]
launch_waveform=gui
```



- **launch\_waveform=<batch|gui>**: The `gui` option enables the Live Waveform viewer, while the `batch` option will record the waveform activity for post-processing.

A Live Waveform viewer is spawned during the execution of the hardware emulation, which allows you to examine the waveforms in detail.

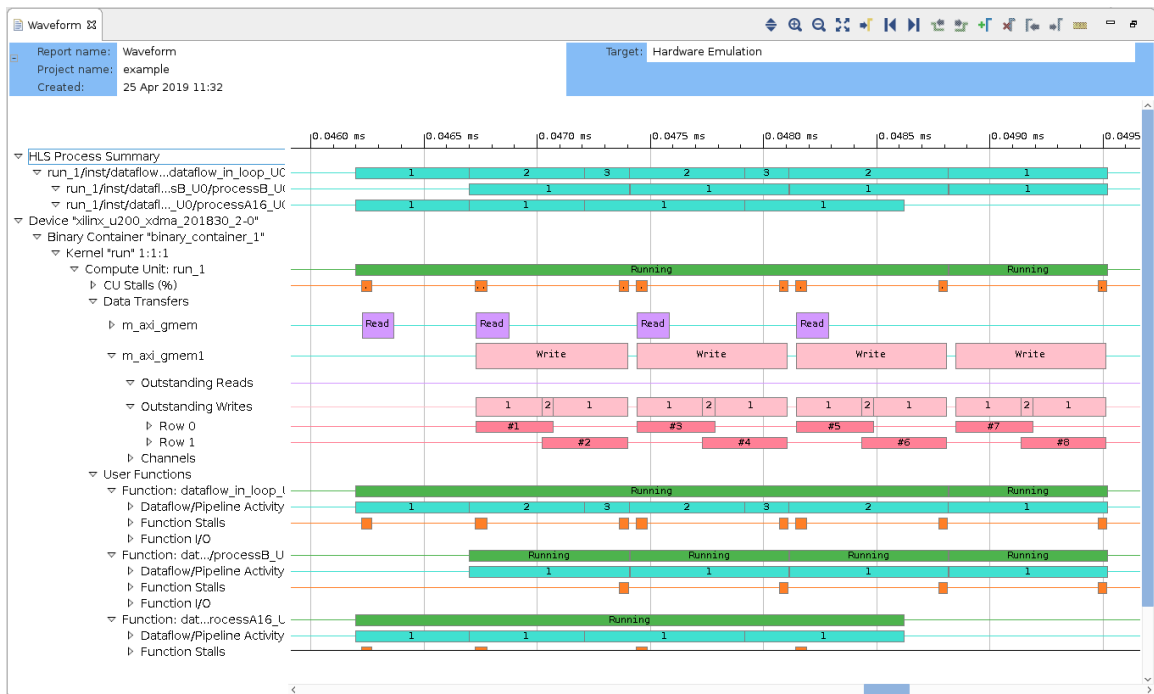
4. Execute hardware emulation. The hardware transaction data is collected in the file `<hardware_platform>-<device_id>-<xclbin_name>.wdb`.
5. If no Live Waveform viewer is requested, follow the steps below to open the Waveform view:
  - a. Start the SDx IDE by running the following command: `$sdx`.
  - b. Choose a workspace when prompted.
  - c. Select **File** → **Open File**, browse to the `.wdb` file generated during hardware emulation.

Alternatively, `xsim` can be used to open the `.wdb` file using the following command: `xsim --gui <file>.wdb`. For more details about `xsim`, refer to *Vivado Design Suite User Guide: Logic Simulation (UG900)*.

## Data Interpretation Waveform View

The following image shows the Waveform view:

Figure 12: Waveform View



The waveform view is organized hierarchically for easy navigation.

**Note:** This view is based on the actual waveforms generated during hardware emulation (Kernel Trace). This allows this view to descend all the way down to the individual signals responsible for the abstracted data. However, as it is post processing the data, no additional signals can be added, and some of the runtime analysis such as DATAFLOW transactions cannot be visualized.

The hierarchy tree and descriptions are:

- **HLS Process Summary:** This summary section contains a hierarchical representation of the activity report of each sequential process contained within the generated RTL. Visualizing the active processes within the HLS design allows to profile in detail which process is active for how long within each activation of the top module. Therefore, this view enables the analysis with respect to individual process performance as well as the overall concurrent execution of independent processes. According to Amdahl's Law, processes dominating the overall execution have the highest potential to improve performance, if process execution time can be reduced.
- **Device "name":** Target device name.
  - **Binary Container "name":** Binary container name.
    - **Kernel "name" 1:1:1:** For each kernel and for each compute unit of that kernel, this section breaks down the activities originating from the compute unit.
      - **Compute Unit: "name":** Compute unit name.
        - **CU Stalls (%):** Stall signals are provided by the HLS tool to inform you when a portion of their circuit is stalling because of external memory accesses, internal streams (that is, dataflow), or external streams (that is, OpenCL pipes). The stall bus, shown in detailed kernel trace, compiles all of the lowest level stall signals and reports the percentage that are stalling at any point in time. This provides a factor of how much of the kernel is stalling at any point in the simulation.
 

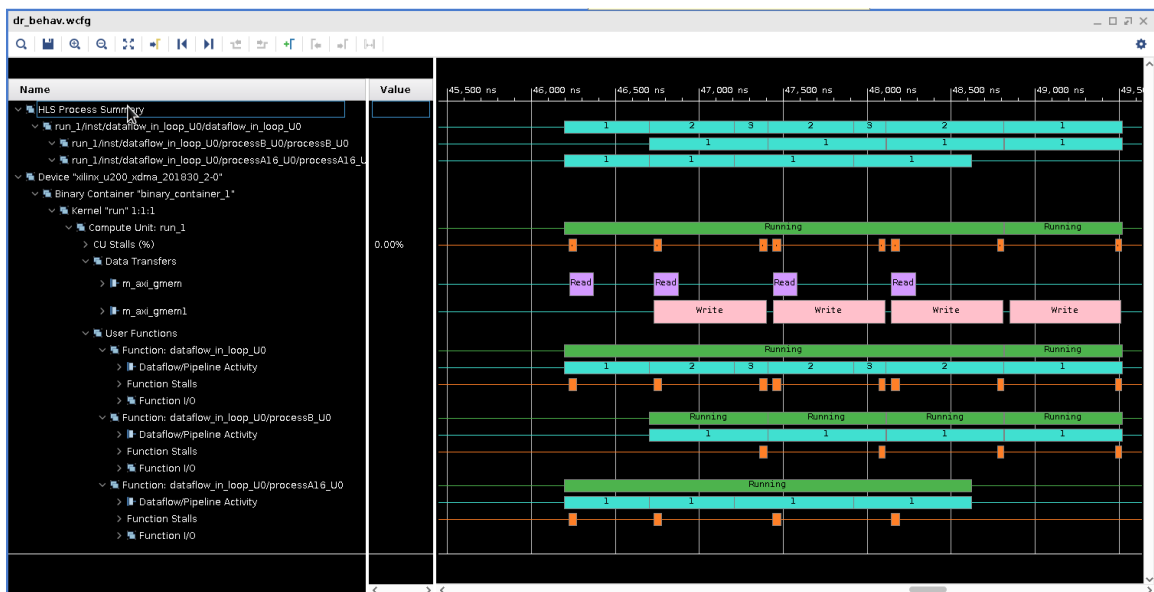
For example, if there are 100 lowest level stall signals, and 10 are active on a given clock cycle, then the CU Stall percentage is 10%. If one goes inactive, then it would be 9%.
- **Data Transfers:** This shows the read/write data transfer accesses originating from each Master AXI port of the compute unit to the DDR.
- **User Functions:** This information is available for the HLS tool kernels and shows the user functions.
  - **Function: "name":**
    - **Dataflow/Pipeline Activity:** This shows the number of parallel executions of the function if the function is implemented as a dataflow process.
    - **Active Iterations:** This shows the currently active iterations of the dataflow. The number of rows is dynamically incremented to accommodate the visualization of any concurrent execution.

- **StallNoContinue:** This is a stall signal that tells if there were any output stalls experienced by the dataflow processes (function is done, but it has not received a continue from the adjacent dataflow process).
- **RTL Signals:** These are the underlying RTL control signals that were used to interpret the above transaction view of the dataflow process.
- **Function Stalls:** Shows the different types of stalls experienced by the process.
  - **External Memory:** Stalls experienced while accessing the DDR memory.
  - **External Stream:** Stalls triggered by streaming.
- **Function I/O:** Actual interface signals grouped according to their associated block interfaces.
  - **Function: "name":** Function name.
  - **Function: "name":** Function name.

## Data Interpretation Live Waveform

The following figure shows the live waveform viewer while running hardware emulation.

Figure 13: Live Waveform Viewer



The live waveform viewer is organized hierarchically for easy navigation. Below are the hierarchy tree and descriptions.

**Note:** As the live waveform viewer is presented only as part of the actual hardware simulation run (xsim), you can annotate extra signals and internals of the register transfer (RTL) to the same view. Also, all grouped and combined groups can be fully expanded to the actual contributing signals. For more information on working with xsim, refer to *Vivado Design Suite User Guide: Logic Simulation (UG900)*.

- **HLS Process Summary:** This summary section contains a hierarchical representation of the activity report of each sequential process contained within the generated RTL. Visualizing the active processes within the HLS design allows to profile in detail which process is active for how long within each activation of the top module. Therefore, this view enables the analysis with respect to individual process performance as well as the overall concurrent execution of independent processes. According to Amdahl's Law, processes dominating the overall execution have the highest potential to improve performance if process execution time can be reduced.
- **Device "name":** Target device name.
  - **Binary Container "name":** Binary container name.
  - **Kernel "name" 1:1:1:** For each kernel and for each compute unit of that kernel this section breaks down the activities originating from the compute unit.
    - **Compute Unit: "name":** Compute unit name.
      - **CU Stalls (%):** Stall signals are provided by the Vivado HLS tool to inform you when a portion of the circuit is stalling because of external memory accesses, internal streams (that is, dataflow), or external streams (that is, OpenCL™ pipes). The stall bus shown in detailed kernel trace compiles all of the lowest level stall signals and reports the percentage that are stalling at any point in time. This provides a factor of how much of the kernel is stalling at any point in the simulation.
 

For example: If there are 100 lowest level stall signals, and 10 are active on a given clock cycle, then the CU Stall percentage is 10%. If one goes inactive, then it would be 9%.
  - **Data Transfers:** This shows the read/write data transfer accesses originating from each Master AXI port of the compute unit to the DDR.
  - **User Functions:** This information is available for the HLS kernels and shows the user functions.
    - **Function: "name":**
      - **Dataflow/Pipeline Activity:** This shows the number of parallel executions of the function if the function is implemented as a dataflow process
      - **Active Iterations:** This shows the currently active iterations of the dataflow. The number of rows is dynamically incremented to accommodate the visualization of any concurrent execution.

- **StallNoContinue:** This is a stall signal that tells if there were any output stalls experienced by the dataflow processes (function is done, but it has not received a continue from the adjacent dataflow process).
- **RTL Signals:** These are the underlying RTL control signals that were used to interpret the above transaction view of the dataflow process.
- **Function Stalls:** Shows the different types of stalls experienced by the process.
  - **External Memory:** Stalls experienced while accessing the DDR memory.
  - **External Stream:** Stalls triggered by streaming.
- **Function I/O:** Actual interface signals grouped according to their associated block interfaces.
- **Function: "name":** Function name
- **Function: "name":** Function name

## Guidance

The Guidance view is designed to provide you with feedback throughout the development process. It presents in a single location all issues encountered from building the actual design all the way through runtime analysis.

It is crucial to understand that the Guidance view is intended to help you to identify potential issues in the design. These issues might be source code related or due to missed tool optimizations. Also, the rules are generic rules based on experiences on a vast set of reference designs. Nevertheless, these rules might not be applicable for a specific design. Therefore, it is up to you to understand the specific guidance rules, and take appropriate action based on your specific algorithm and requirements.

## GUI Flow

The Guidance view is automatically populated and displayed in the lower central tab view. After running hardware emulation, the Guidance view might look like the following:

Figure 14: Guidance View

| Name                          | Threshold | Actual  | Details   |
|-------------------------------|-----------|---------|---|
| Emulation-HW (27)             |           |         |   |
| example-Default (23)          |           |         |   |
| Host Data Transfer (3)        |           |         |   |
| HOST_WRITE_TRANSFER_SIZE (1)  | > 4.096   |         |   |
| ✓ HOST_WRITE_TRANSFER_SIZE #1 | > 4.096   | 32.768  | Host write average size was 32.768 KB across 4 tr |
| HOST_MIGRATE_MEM (1)          | > 0       |         |   |
| ✓ HOST_MIGRATE_MEM #1         | > 0       | 8       | Migrate Memory OpenCL APIs were used 8 time(s)    |
| HOST_READ_TRANSFER_SIZE (1)   | > 4.096   |         |   |
| ✓ HOST_READ_TRANSFER_SIZE #1  | > 4.096   | 32.768  | Host read average size was 32.768 KB across 4 tr  |
| Resource Usage (6)            |           |         |   |
| KERNEL_UTIL (1)               | = 100.000 |         |   |
| ✓ KERNEL_UTIL #1              | = 100.000 | 100.000 | Kernel run - global size: 1, local size: 1.       |
| KERNEL_COUNT (1)              | > 1       |         |   |
| ⚠ KERNEL_COUNT #1             | > 1       | 1       | Kernel run was executed 4 time(s) with 1 comput   |
| OVERUSED_CUS (1)              | < 16      |         |   |
| ✓ OVERUSED_CUS #1             | < 16      | 1       | Kernel run required 1 compute unit call(s).       |
| DEVICE_UTIL (1)               | > 0       |         |   |
| ✓ DEVICE_UTIL #1              | > 0       | 0.339   | Device xilinx_kcu1500_dynamic_5_0-0 was utilize   |
| UNUSED_CUS (1)                | > 0       |         |   |

**Note:** You can produce the Guidance view through the Vivado HLS tool post-compilation, but you will not get Profile Rule Checks.

To simplify visualizing the guidance information, the GUI flow allows you to search, and filter the Guidance view to locate specific guidance rule entries. It is also possible to collapse or expand the tree view or even suppress the hierarchical tree representation and visualize a condensed representation of the guidance rules. Finally, you can select what is shown in the Guidance view. You can enable or disable the visualization of warnings, as well as met rules, and restrict the specific content based on the source of the messages such as build and emulation.

By default, the Guidance view shows all guidance information for the project selected in the drop-down.

To restrict the content to an individual build or run step, do the following:

1. Use the command **Window → Preferences**
2. Select the category **Xilinx Sdx → Guidance**.
3. Deselect **Group guidance rule checks by project**.

## Command Line

The Guidance data is best analyzed through the GUI, which consolidates all guidance information for the flow. Nevertheless, the tool automatically generates HTML files containing the guidance information. As guidance information is generated throughout the tool flow, several guidance files are generated. The simplest way to locate the guidance reports is to search for the `guidance.html` files.

```
find . -name "*guidance.html" -print
```

This command lists all guidance files generated, which can be opened with any web browser.

## Data Interpretation

The Guidance view places each entry in a separate row. Each row might contain the name of the guidance rule, threshold value, actual value, and a brief but specific description of the rule. The last field provides a link to reference material intended to assist in understanding and resolving any of the rule violations.

In the GUI Guidance view, guidance rules are grouped by categories and unique IDs in the Name column and annotated with symbols representing the severity. These are listed individually in the HTML report. In addition, as the HTML report does not show tooltips, a full Name column is included in the HTML report as well.

The following list describes all fields and their purpose as included in the HTML guidance reports.

- **Id:** Each guidance rule is assigned a unique id. Use this id to uniquely identify a specific message from the guidance report.
- **Name:** The Name column displays a mnemonic name uniquely identifying the guidance rule. These names are designed to assist in memorizing specific guidance rules in the view.
- **Severity:** The Severity column allows the easy identification of the importance of a guidance rule.
- **Full Name:** The Full Name provides a less cryptic name compared to the mnemonic name in the Name column.
- **Categories:** Most messages are grouped within different categories. This allows the GUI to display groups of messages within logical categories under common tree nodes in the Guidance view.
- **Threshold:** The Threshold column displays an expected threshold value, which determines whether or not a rule is met. The threshold values are determined from many applications that follow good design and coding practices.
- **Actual:** The Actual column displays the values actually encountered on the specific design. This value is compared against the expected value to see if the rule is met.

- **Details:** The Details column provides a brief but specific message describing the specifics of the current rule.
- **Resolution:** The Resolution column provides a pointer to common ways the model source code or tool transformations can be modified to meet the current rule. Clicking the link brings up a pop-up window or the documentation with tips and code snippets that you can apply to the specific issue.

---

## Using Implementation Tools

### Exploring Kernel Optimizations Using Vivado HLS

All kernel optimizations using OpenCL or C/C++ can be performed from within the SDAccel environment. The primary performance optimizations, such as those discussed in this chapter (pipelining function and loops, applying dataflow to enable greater concurrency between functions and loops, unrolling loops, etc.), are performed by the Xilinx® FPGA design tool, Vivado® HLS tool.

The SDAccel environment automatically calls the HLS tool. However, to use the GUI analysis capabilities, you must launch the HLS tool directly from within the SDAccel environment. Using the HLS tool in standalone mode enables the following enhancements to the optimization methodology:

- Focusing solely on the kernel optimization, there is no requirement to execute emulation.
- The ability to create multiple solutions, compare their results, and explore the solution space to find the most optimum design.
- The ability to use the interactive Analysis Perspective to analyze the design performance.



---

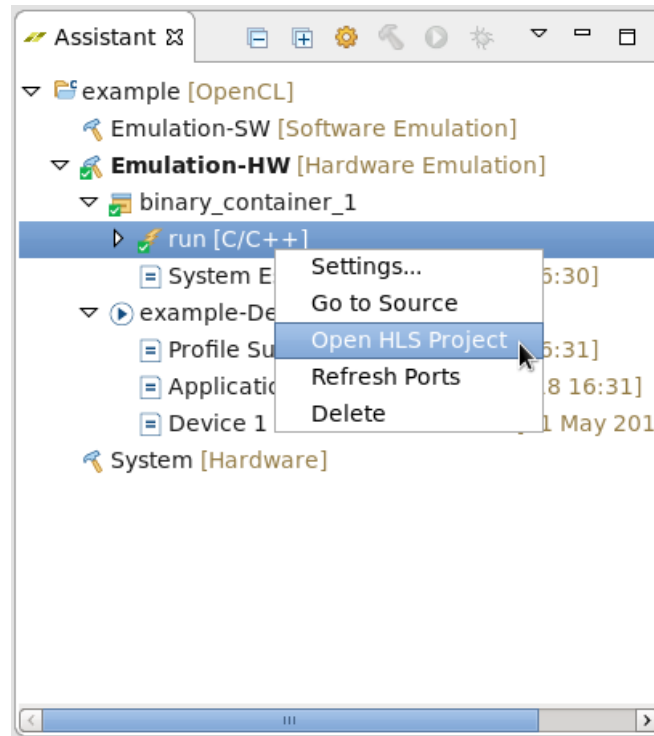
**IMPORTANT!** Only the kernel source code is incorporated back into the SDAccel environment. After exploring the optimization space, ensure that all optimizations are applied to the kernel source code as OpenCL attributes or C/C++ pragmas.

---

To open the HLS tool in standalone mode, from the Assistant window, right-click the hardware function object, and select **Open HLS Project**, as shown in the following figure.



Figure 15: Open HLS Project



## Controlling FPGA Implementation with the Vivado Design Suite

The SDAccel development environment provides a smooth flow from an OpenCL/C/C++ model all the way to an FPGA accelerated implementation. In most cases, this flow completely abstracts away the underlying fact that the programmable region in the FPGA is configured to implement the kernel functionality. This fully isolates the developer from typical hardware constraints such as routing delays and kernel placement. However, in some cases these concerns will have to be looked at especially when large designs are to be implemented. Towards this end, the development environment allows you to fully control the Vivado Design Suite backend tool.

The SDAccel environment calls the Vivado Design Suite to automatically run RTL synthesis and implementation. You also have the option of launching the design suite directly from within the SDAccel environment. When invoking the Vivado Integrated Design Environment (IDE) in standalone mode in the SDAccel environment, you can open the Vivado synthesis project or the Vivado implementation project to edit, manage, and control the project.

The Vivado project can be opened in the SDAccel environment after the build targeting the system configuration has completed.

To open Vivado IDE in standalone mode, from the Xilinx drop-down menu, select **Vivado Integration** and **Open Vivado Project**. Choose between the Vivado synthesis and implementation projects, and click **OK**.

Using the Vivado IDE in standalone mode enables the exploration of various synthesis and implementation options for further optimizing the kernel for performance and area. Familiarity with the design suite is recommended to make the most use of these features.



**IMPORTANT!** *The optimization switches applied in the standalone project are not automatically incorporated back into the SDAccel environment. After exploring the optimization space, ensure that all optimization parameters are passed to the SDAccel environment using the `--xp` option for `xocc`. For example:*

```
--xp "vivado_prop:run.impl_1.{STEPS.PLACE_DESIGN.ARGS.TCL.POST}={<File and path>}"
```

This optimization flow is supported in the command line flow by calling `xocc -interactive` to bring up the Vivado IDE, on the current project. In the IDE, generate a DCP, which can be saved and reused during linking with `xocc`. The specific options are:

- `--interactive` allows the Vivado IDE to be launched from within the `xocc` environment, with the right project loaded.
- `--reuse_impl` allows a pre-implemented and timing closed Vivado tool design checkpoint (.dcp) file to be brought in and used directly in SDx environment flow for xclbin generation.

# Kernel Optimization

One of the key advantages of an FPGA is its flexibility and capacity to create customized designs specifically for your algorithm. This enables various implementation choices to trade off algorithm throughput vs. power consumption. The downside of creating custom logic is that the design needs to go through the traditional FPGA design flow.

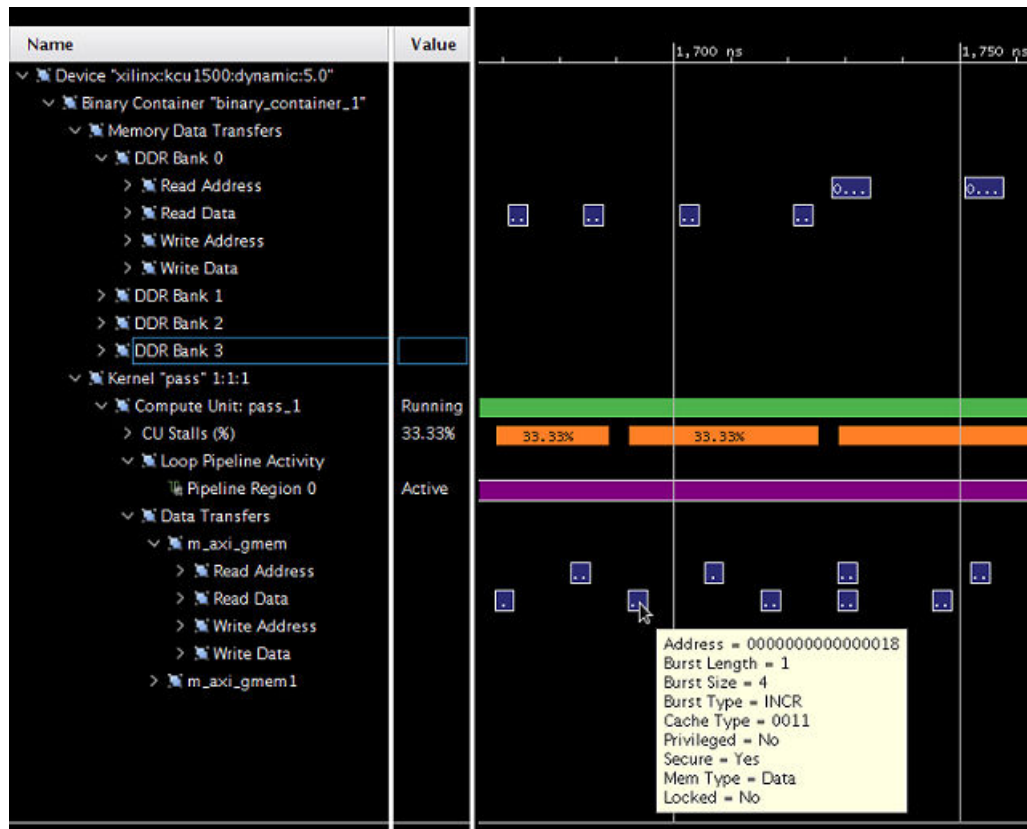
The following guidelines help manage the design complexity and achieve the desired design goals.

---

## Interface Attributes (Detailed Kernel Trace)

The detailed kernel trace provides easy access to the AXI transactions and their properties. The AXI transactions are presented for the global memory, as well as the Kernel side (Kernel "pass" 1:1:1) of the AXI interconnect. The following figure illustrates a typical kernel trace of a newly accelerated algorithm.

Figure 16: Accelerated Algorithm Kernel Trace



The following fields are the most important in respect to performance:

- **Burst Length:** Describes how many beats are sent within one transaction
- **Burst Size:** Describes the number of bytes being transferred as part of one beat

Given a burst length of 1 and just 4 Bytes per package, it will require many individual AXI transactions to transfer any reasonable amount of data.

**Note:** The SDAccel™ environment never creates burst sizes less than 4 bytes, even if smaller data is transmitted. In this case, if consecutive items are accessed without AXI bursts enabled, you can observe multiple AXI reads to the same address.

Therefore, small burst lengths, as well as burst sizes, considerably less than 512-bits are good opportunities to optimize interface performance. The following sections show improved implementations:

- [Using Burst Data Transfers](#)
- [Using AXI4 Data Width](#)

## Top Level Dataflow

Top level dataflow optimization provides the capability to the kernels to be executed in a parallel processing architecture by overlapping multiple kernel execution. In cases of long kernel latency, this optimization is useful and will improve the performance of the overall application. Xilinx recommends enabling this optimization only to kernels which have dataflow at the top function by applying the interface pragma `ap_ctrl_chain` to the return port (along with `s_axilite`).

```
#pragma HLS INTERFACE ap_ctrl_chain port=return bundle=control
```

The following example is provided to illustrate that both the `ap_ctrl_chain` and `s_axilite` option is required for the return port to enable top level dataflow.

```
void N_stage_Adders(int *input, int *output, int incr, int size)
{
    ...
    #pragma HLS INTERFACE s_axilite port=return bundle=control
    #pragma HLS INTERFACE ap_ctrl_chain port=return bundle=control
    ...
    ...
    #pragma HLS dataflow
    read_input(input, streamArray[0], size);
    compute_loop: for (int i = 0 ; i < STAGES ; i++)
    {
        #pragma HLS UNROLL
        adder(streamArray[i], streamArray[i+1], incr, size);
    }
    write_result(output, streamArray[STAGES], size);
}
```

**Note:** The profiling capability for this optimization is currently only available in hardware.

## Using Burst Data Transfers

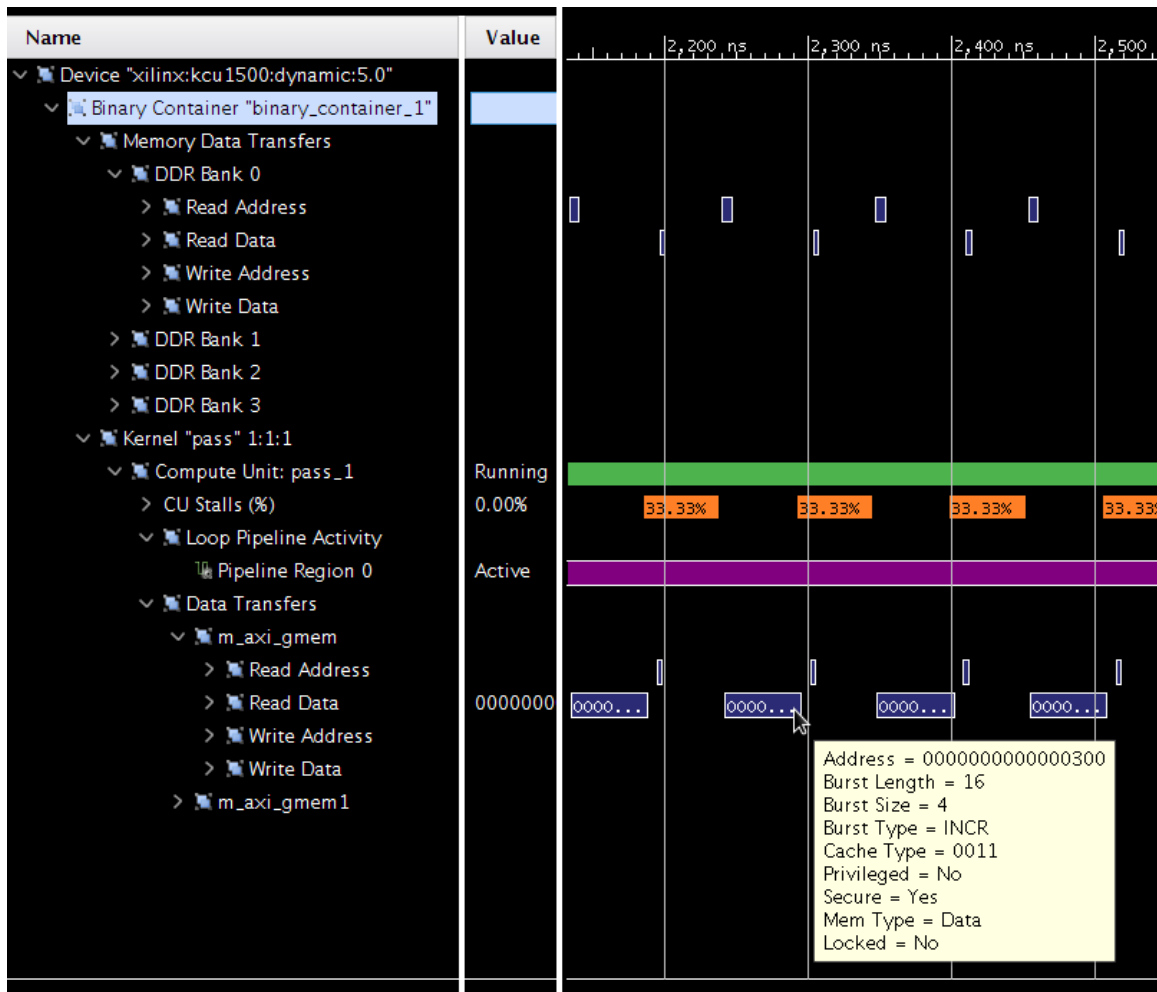
Transferring data in bursts hides the memory access latency and improves bandwidth usage and efficiency of the memory controller.



**RECOMMENDED:** *Infer burst transfers from successive requests of data from consecutive address locations. Refer to "Inferring Burst Transfer from/to Global Memory" in SDAccel Environment Programmers Guide (UG1277) for more details.*

If burst data transfers occur, the detailed kernel trace will reflect the higher burst rate as a larger burst length number:

Figure 17: Burst Data Transfer with Detailed Kernel Trace



In the previous figure, it is also possible to observe that the memory data transfers following the AXI interconnect are actually implemented rather differently (shorter transaction time). If you hover over these transactions, you would see that the AXI interconnect has packed the 16x4 Byte transaction into a single burst transaction of 1x64 Bytes. This effectively uses the AXI4 bandwidth which is even more favorable. The next section focuses on this optimization technique in more detail.

Burst inference is heavily dependent on coding style and access pattern. To avoid potential modeling pitfalls, refer to the *SDAccel Environment Programmers Guide* (UG1277). However, you can ease burst detection and improve performance by isolating data transfer and computation, as shown in the following code snippet:

```
void kernel(T in[1024], T out[1024]) {
    T tmpIn[1024];
    T tmpOu[1024];
    read(in, tmpIn);
    process(tmpIn, tmpOut);
    write(tmpOut, out);
}
```

In short, the function `read` is responsible for reading from the AXI input to an internal variable (`tmpIn`). The computation is implemented by the function `process` working on the internal variables `tmpIn` and `tmpOut`. The function `write` takes the produced output and writes to the AXI output.

The isolation of the read and write function from the computation results in:

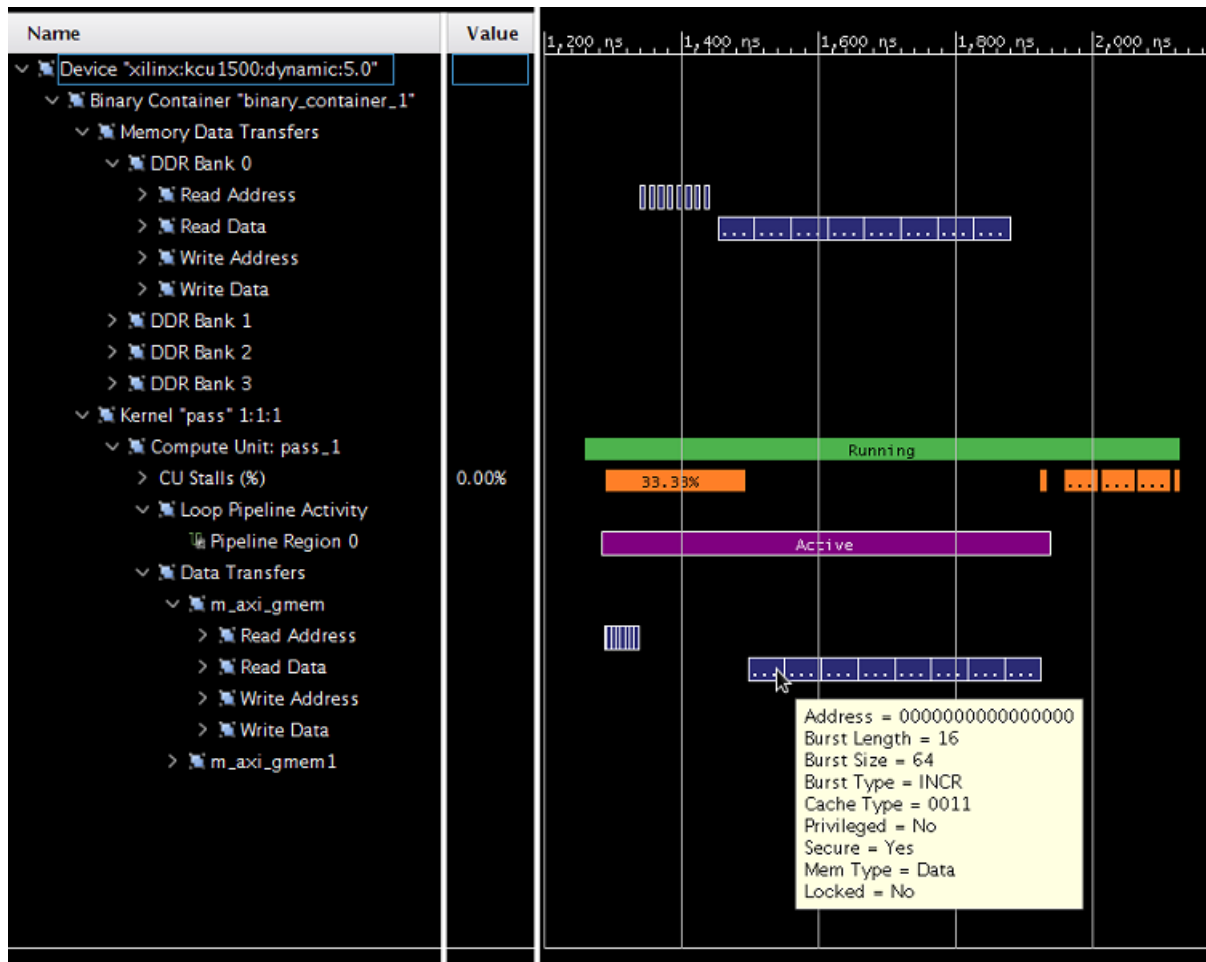
- Simple control structures (loops) in the read/write function which makes burst detection simpler.
- Isolation of the computational function away from the AXI interfaces, simplifies potential kernel optimization. See the [Kernel Optimization](#) chapter for more information.
- Internal variables are mapped to on-chip memory, which allow faster access compared to AXI transactions. Acceleration platforms supported in SDAccel environment can have as much as 10 MB on-chip memories that can be used as pipes, local memories, and private memories. Using these resources effectively can greatly improve the efficiency and performance of your applications.

## Using AXI4 Data Width

The user data width between the kernel and the memory controller can be configured by the SDAccel environment compiler based on the data types of the kernel arguments. To maximize the data throughput, Xilinx recommends that you choose data types mapping to the full data width on the memory controller. The memory controller in all supported acceleration cards supports 512-bit user interface, which can be mapped to OpenCL™ vector data types, such as `int16` or C/C++ arbitrary precision data type `ap_int<512>`.

As shown on the following figure, you can observe burst AXI transactions (Burst Length 16) and a 512-bit beat size (Burst Size 64 Bytes).

Figure 18: Burst AXI Transactions



This example shows good interface configuration as it maximizes AXI data width and also shows actual burst transactions.

Complex structs or classes, used to declare interfaces, can lead to very complex hardware interfaces due to memory layout and data packing differences. This can introduce potential issues that are very difficult to debug in a complex system.



**RECOMMENDED:** Use simple structs for kernel arguments that can be packed to a power of two boundary. Refer to the Custom Data Type Example in kernel\_to\_gmem category at [Xilinx On-boarding Example GitHub](#) for the recommended way to use structs.

## OpenCL API Attributes

The OpenCL API provides attributes to support a more automatic approach to incrementing AXI data width usage. The change of the interface data types, as stated above, is supported in the API as well but will require the same code changes as C/C++ to the algorithm to accommodate the larger input vector.



To eliminate manual code modifications, the following OpenCL attributes are interpreted to perform data path widening and vectorization of the algorithm. A detailed description can be found in the *SDx Pragma Reference Guide* ([UG1253](#)).

- `vec_type_hint`
- `reqd_work_group_size`
- `xcl_zero_global_work_offset`

In the following example, examine the combined functionality on the following case:

```
__attribute__((reqd_work_group_size(64, 1, 1)))
__attribute__((vec_type_hint(int)))
__attribute__((xcl_zero_global_work_offset))
__kernel void vector_add(__global int* c, __global const int* a, __global
const int* b) {
    size_t idx = get_global_id(0);
    c[idx] = a[idx] + b[idx];
}
```

In this case, the hard coded interface is a 32-bit wide data path (`int *c, int* a, int *b`), which drastically limits the memory throughput if implemented directly. However, the automatic widening and transformation is applied, based on the values of the three attributes.

- **`__attribute__((vec_type_hint(int)))`**: Declares that `int` is the main type used for computation and memory transfer (32-bit). This knowledge is used to calculate the vectorization/widening factor based on the target bandwidth of the AXI interface (512-bits). In this example the factor would be  $16 = 512\text{-bits} / 32\text{-bit}$ . This implies that in theory, 16 values could be processed if vectorization can be applied.
- **`__attribute__((reqd_work_group_size(X, Y, Z)))`**: Defines the total number of work items (where `X`, `Y`, and `Z` are positive constants).  $X * Y * Z$  is the maximum number of work items therefore defining the maximum possible vectorization factor which would saturate the memory bandwidth. In this example, the total number of work items is  $64 * 1 * 1 = 64$ .

The actual vectorization factor to be applied will be the greatest common divider of the vectorization factor defined by the actual coded type or the `vec_type_hint`, and the maximum possible vectorization factor defined through `reqd_work_group_size`.

The quotient of maximum possible vectorization factor divided by the actual vectorization factor provides the remaining loop count of the OpenCL description. As this loop is pipelined, it can be advantageous to have several remaining loop iterations to take advantage of a pipelined implementation. This is especially true if the vectorized OpenCL code has long latency.

There is one optional parameter that is highly recommended to be specified for performance optimization on OpenCL interfaces.

- The `__attribute__((xcl_zero_global_work_offset))` instructs the compiler that no global offset parameter is used at runtime, and all accesses are aligned. This gives the compiler valuable information with regard to alignment of the work groups, which in turn usually propagate to the alignment of the memory accesses (less hardware).

It should be noted, that the application of these transformations changes the actual design to be synthesized. Partially unrolled loops require reshaping of local arrays in which data is stored. This usually behaves nicely, but can interact poorly in rare situations.

For example:

- For partitioned arrays, when the partition factor is not divisible by the unrolling/vectorization factor.
  - The resulting access requires a lot of multiplexers and will create a difficult problem for the scheduler (might severely increase memory usage and compilation time), Xilinx recommends that you use partitioning factors that are powers of two (as the vectorization factor is always a power of two).
- If the loop being vectorized has an unrelated resource constraint, the scheduler complains about `II` not being met.
  - This is not necessarily correlated with a loss of performance (usually it is still performing better) because the `II` is computed on the unrolled loop (which has therefore a multiplied throughput for each iteration).
  - The scheduler informs you of the possible resources constraints and resolving those will further improve the performance.
  - Note that a common occurrence is that a local array does not get automatically reshaped (usually because it is accessed in a later section of the code in non-vectorizable way).

## Reducing Kernel-to-Kernel Communication Latency with OpenCL Pipes

This section specifically applies to OpenCL kernels. For C++ kernels, kernel-to-kernel streaming is provided. This is discussed in [Memory Data Transfer Types](#).

The OpenCL API 2.0 specification introduces a new memory object called a pipe. A pipe stores data organized as a FIFO. Pipe objects can only be accessed using built-in functions that read from and write to a pipe. Pipe objects are not accessible from the host. Pipes can be used to stream data from one kernel to another inside the FPGA without having to use the external memory, which greatly improves the overall system latency.

In the SDAccel development environment, pipes must be statically defined outside of all kernel functions; dynamic pipe allocation using the OpenCL 2.x `clCreatePipe` API is not currently supported. The depth of a pipe must be specified by using the `xcl_reqd_pipe_depth` attribute in the pipe declaration.

The valid depth values are as follows:

- 16
- 32
- 64
- 128
- 256
- 512
- 1024
- 2048
- 4096
- 8192
- 16384
- 32768

A given pipe can have one and only one producer and consumer in different kernels.

```
pipe int p0 __attribute__((xcl_reqd_pipe_depth(32)));
```

Pipes can be accessed using standard OpenCL `read_pipe()` and `write_pipe()` built-in functions in non-blocking mode or using the Xilinx extended `read_pipe_block()` and `write_pipe_block()` functions in blocking mode. The status of pipes can be queried using OpenCL `get_pipe_num_packets()` and `get_pipe_max_packets()` built-in functions. See the OpenCL C Specification, Version 2.0 from Khronos Group for more details on these built-in functions.

The following function signatures are the currently supported pipe functions, where `gentype` indicates the built-in OpenCL C scalar integer or floating-point data types.

```
int read_pipe_block (pipe gentype p, gentype *ptr)
int write_pipe_block (pipe gentype p, const gentype *ptr)
```

The following “Blocking Pipes Example” from [SDAccel Getting Started Examples](#) on GitHub uses pipes to pass data from one processing stage to the next using blocking `read_pipe_block()` and `write_pipe_block()` functions:

```
pipe int p0 __attribute__((xcl_reqd_pipe_depth(32)));
pipe int p1 __attribute__((xcl_reqd_pipe_depth(32)));
// Input Stage Kernel : Read Data from Global Memory and write into Pipe P0
kernel __attribute__((reqd_work_group_size(1, 1, 1)))
void input_stage(__global int *input, int size)
{
    __attribute__((xcl_pipeline_loop))
    mem_rd: for (int i = 0 ; i < size ; i++)
    {
        //blocking Write command to pipe P0
        write_pipe_block(p0, &input[i]);
    }
}
```

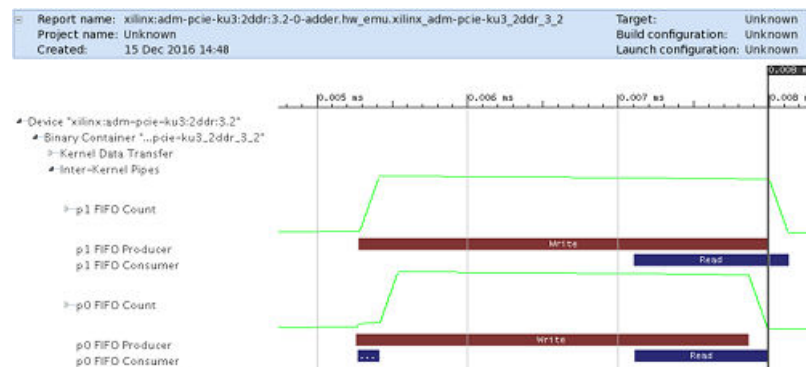
```

// Adder Stage Kernel: Read Input data from Pipe P0 and write the result
// into Pipe P1
kernel __attribute__((reqd_work_group_size(1, 1, 1)))
void adder_stage(int inc, int size)
{
    __attribute__((xcl_pipeline_loop))
    execute: for(int i = 0 ; i < size ; i++)
    {
        int input_data, output_data;
        //blocking read command to Pipe P0
        read_pipe_block(p0, &input_data);
        output_data = input_data + inc;
        //blocking write command to Pipe P1
        write_pipe_block(p1, &output_data);
    }
}

// Output Stage Kernel: Read result from Pipe P1 and write the result to
Global
// Memory
kernel __attribute__((reqd_work_group_size(1, 1, 1)))
void output_stage(__global int *output, int size)
{
    __attribute__((xcl_pipeline_loop))
    mem_wr: for (int i = 0 ; i < size ; i++)
    {
        //blocking read command to Pipe P1
        read_pipe_block(p1, &output[i]);
    }
}
    
```

The Device Traceline view shows the detailed activities and stalls on the OpenCL pipes after hardware emulation is run. This information can be used to choose the correct FIFO sizes to achieve the optimal application area and performance.

**Figure 19: Device Traceline View**



# Optimizing Computational Parallelism

By default, C/C++ does not model computational parallelism, as it always executes any algorithm sequentially. Conversely, the OpenCL API does model computational parallelism with respect to work groups, but it does not use any additional parallelism within the algorithm description. However, fully configurable computational engines, like FPGAs, allow more freedom to exploit computational parallelism.

## Coding Data Parallelism

To leverage computational parallelism during the implementation of an algorithm on the FPGA, it should be mentioned that the synthesis tool must be able to recognize computational parallelism from the source code first. Loops and functions are prime candidates for reflecting computational parallelism and compute units in the source description. However, even in this case, it is key to verify that the implementation takes advantage of the computational parallelism as in some cases the SDAccel tool might not be able to apply the desired transformation due to the structure of the source code.

It is quite common, that some computational parallelism might not be reflected in the source code to begin with. In this case, it will need to be added. A typical example is a kernel that might be described to operate on a single input value, while the FPGA implementation might execute computations more efficiently in parallel on multiple values. This kind of parallel modeling is described in the [Using Full AXI Data Width](#) section. A 512-bit interface can be created using OpenCL vector data types such as `int16` or C/C++ arbitrary precision data type `ap_int<512>`.

**Note:** These vector types can also be used as a powerful way to model data parallelism within a kernel, with up to 16 data paths operating in parallel in case of `int16`. Refer to the *Median Filter Example* in the *vision* category at [SDAccel Getting Started Examples](#) on GitHub for the recommended way to use vectors.

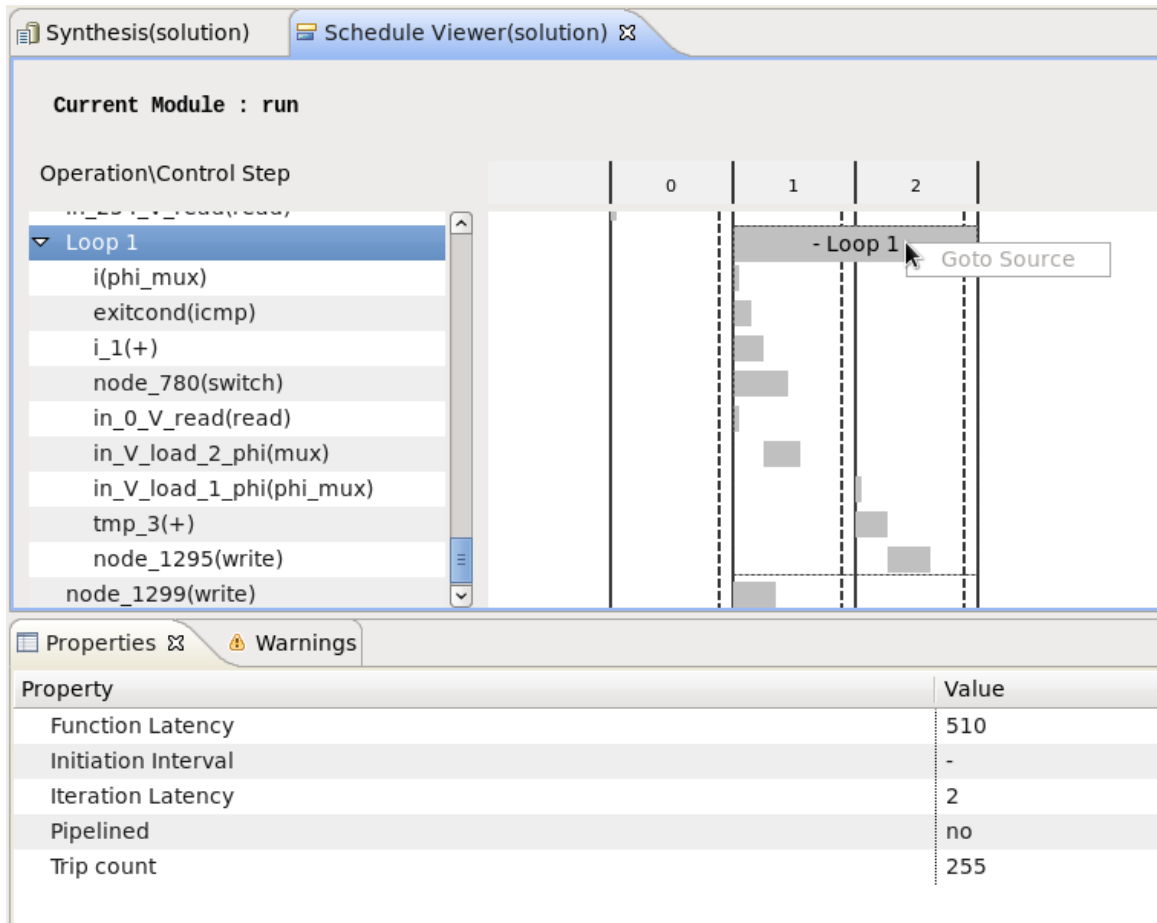
## Loop Parallelism

Loops are the basic C/C++/OpenCL™ API method of representing repetitive algorithmic code. The following example illustrates various implementation aspects of a loop structure:

```
for(int i = 0; i<255; i++) {
    out[i] = in[i]+in[i+1];
}
out[255] = in[255];
```

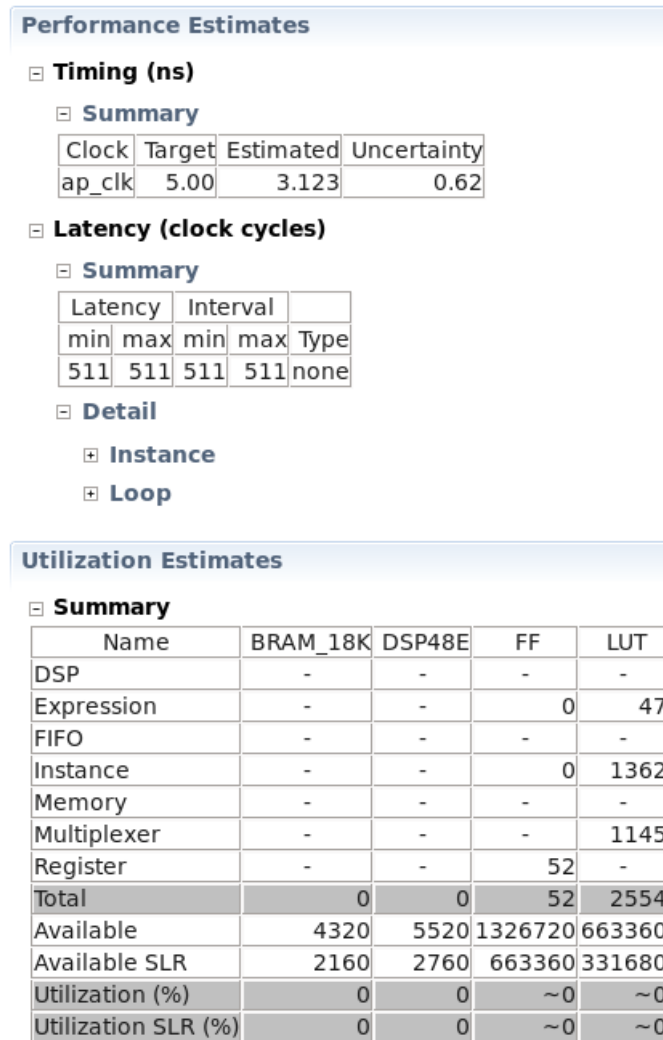
This code iterates over an array of values and adds consecutive values, except the last value. If this loop is implemented as written, each loop iteration requires two cycles for implementation, which results in a total of 510 cycles for implementation. This can be analyzed in detail in the Schedule Viewer in the HLS Project:

Figure 20: Implemented Loop Structure in Schedule Viewer



This can also be analyzed in terms of total numbers and latency through the Vivado synthesis results:

Figure 21: Synthesis Results Performance Estimates



The key numbers here are the latency numbers and total LUT usage. For example, depending on the configuration, you could get latency of 511 and total LUT usage of 47. As you will see, these values can widely vary based on the implementation choices. While this implementation will require very little area, it results in significant latency.

## Unrolling Loops

Unrolling a loop enables the full parallelism of the model to be exploited. To do this, you can simply mark a loop to be unrolled and the tool will create the implementation with the most parallelism possible. To mark a loop to unroll, an OpenCL API loop can be marked with the UNROLL attribute:

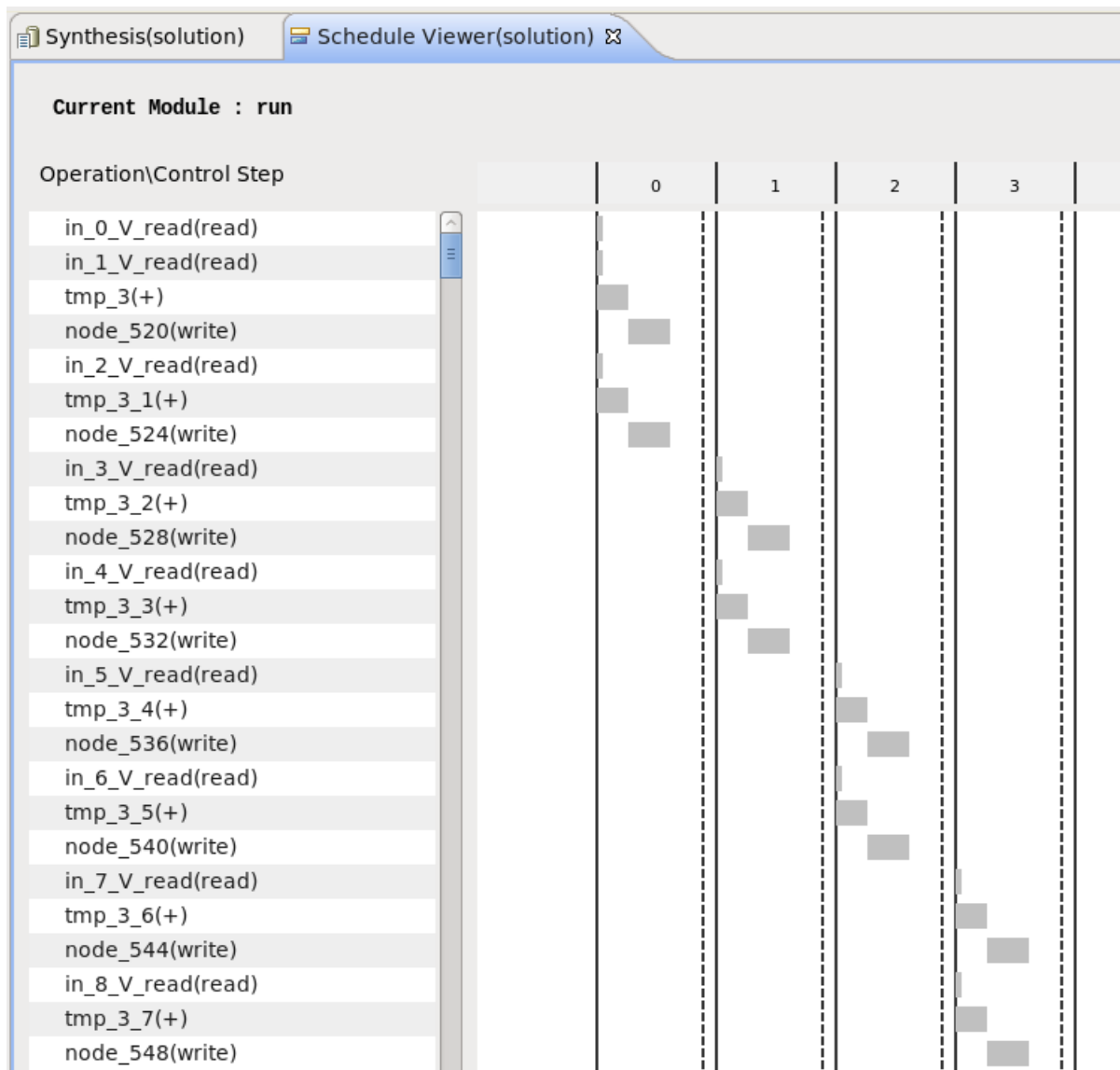
```
__attribute__((opencl_unroll_hint))
```

or a C/C++ loop can use the UNROLL pragma:

```
#pragma HLS UNROLL
```

When applied to this specific example, the Schedule Viewer in the HLS Project will be:

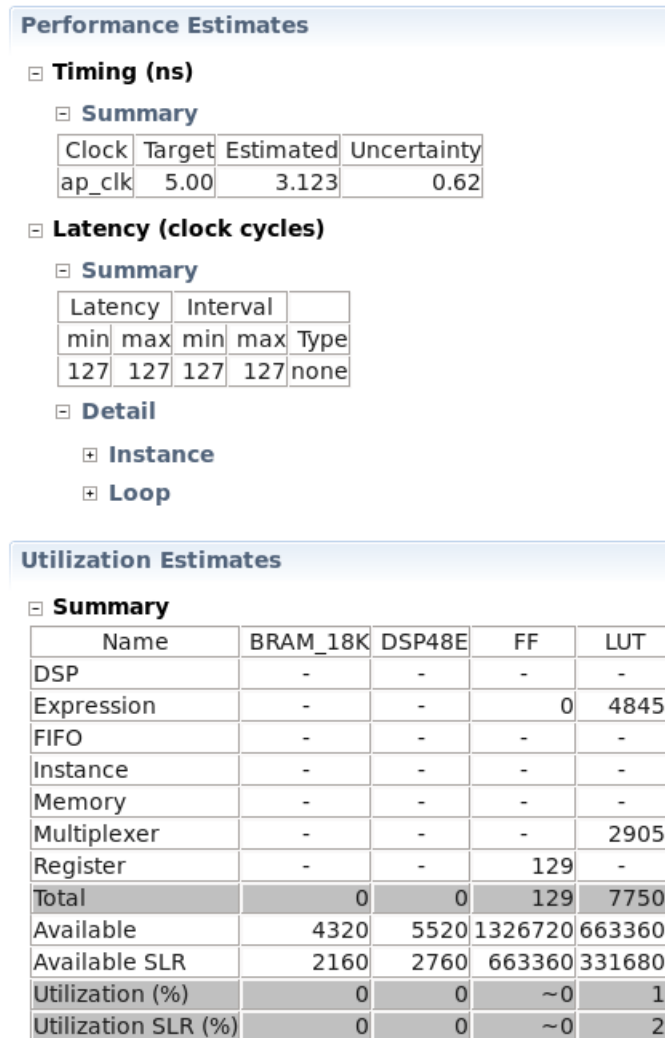
Figure 22: Schedule Viewer



With an estimated performance of:



Figure 23: Performance Estimates



As you can see, the total latency was considerably improved to 127 cycles, and as expected, the computational hardware was increased to 4845 LUTs to perform the same computation in parallel.

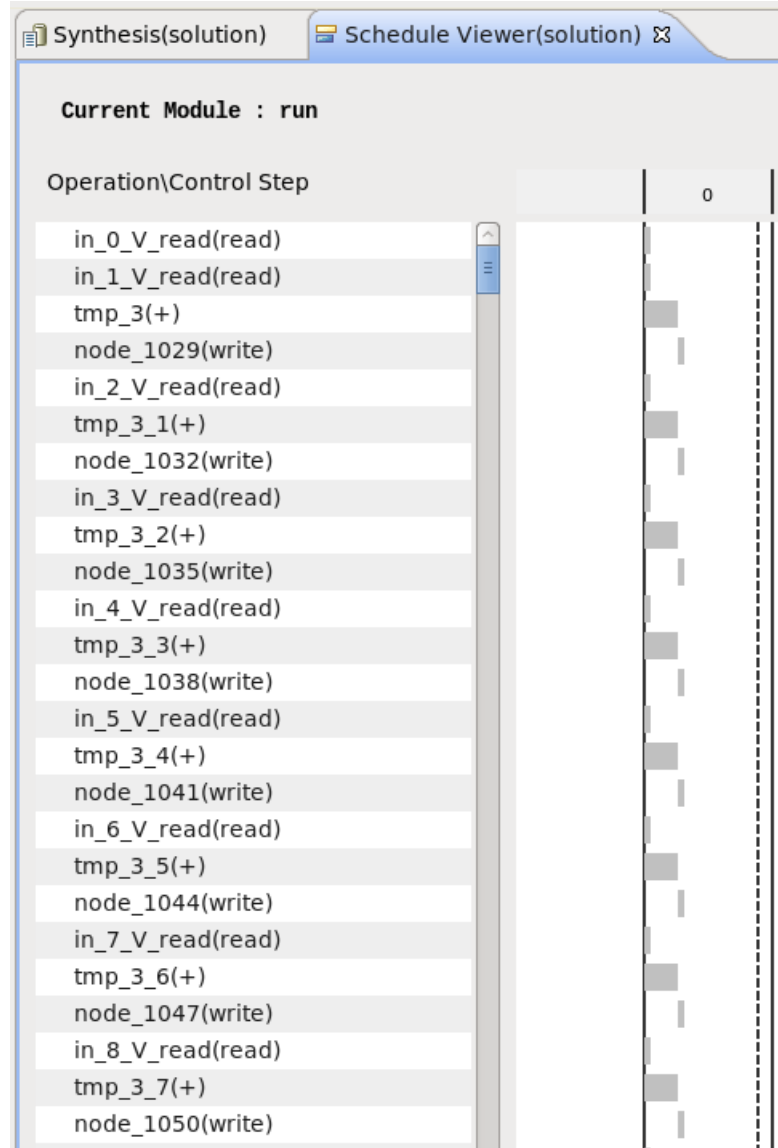
However, if you analyze the for-loop, you might ask why this algorithm cannot be implemented in a single cycle, as each addition is completely independent of the previous loop iteration. The reason is the memory interface to be used for the `out` variable. The SDAccel environment uses dual port memory by default for an array. However, this implies that at most two values can be written to the memory per cycle. Therefore, to see a fully parallel implementation, you must specify that the `out` variable is kept in registers, as shown in the following example:

```
#pragma HLS array_partition variable=out complete dim= 0
```

For more information see the *pragma HLS array\_partition* section in *SDx Pragma Reference Guide (UG1253)*.

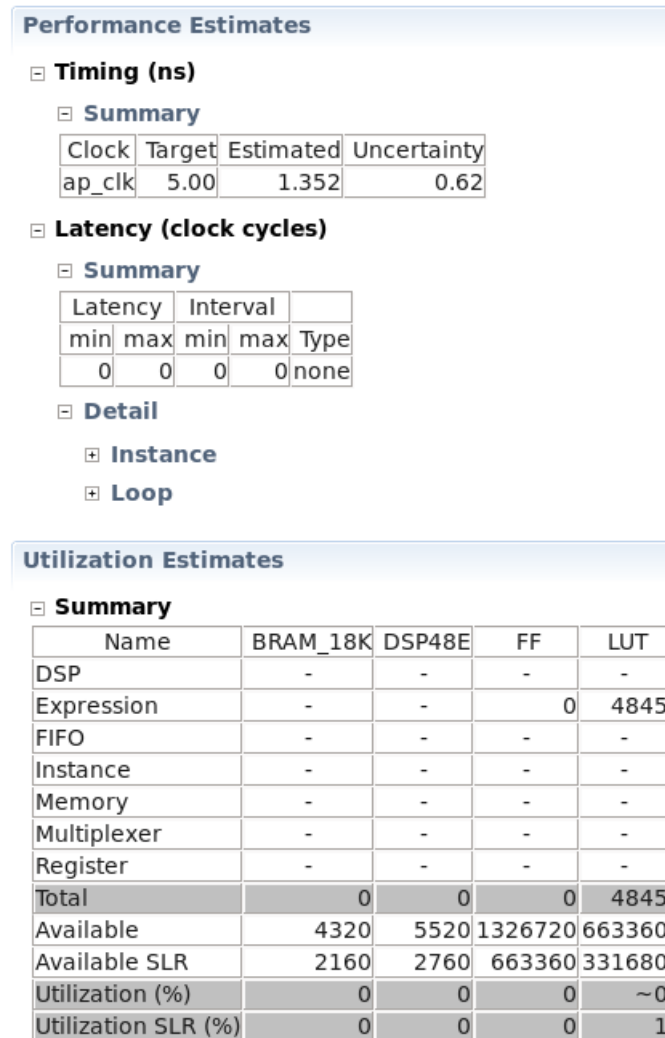
The results of this transformation can be observed in the following Schedule Viewer:

Figure 24: Transformation Results in Schedule Viewer



The associated estimates are:

Figure 25: Transformation Results Performance Estimates



As you can see, this code can be implemented as a combinatorial function requiring only a fraction of the cycle to complete.

## Pipelining Loops

Pipelining loops allows you to overlap iterations of a loop in time. Allowing iterations to operate concurrently is often a good compromise, as resources can be shared between iterations (less resource usage), while requiring less execution time compared to loops that are not unrolled.

Pipelining is enabled in C/C++ through the following pragma:

```
#pragma HLS PIPELINE
```

While the OpenCL API uses the following attribute:

```
__attribute__((xcl_pipeline_loop))
```

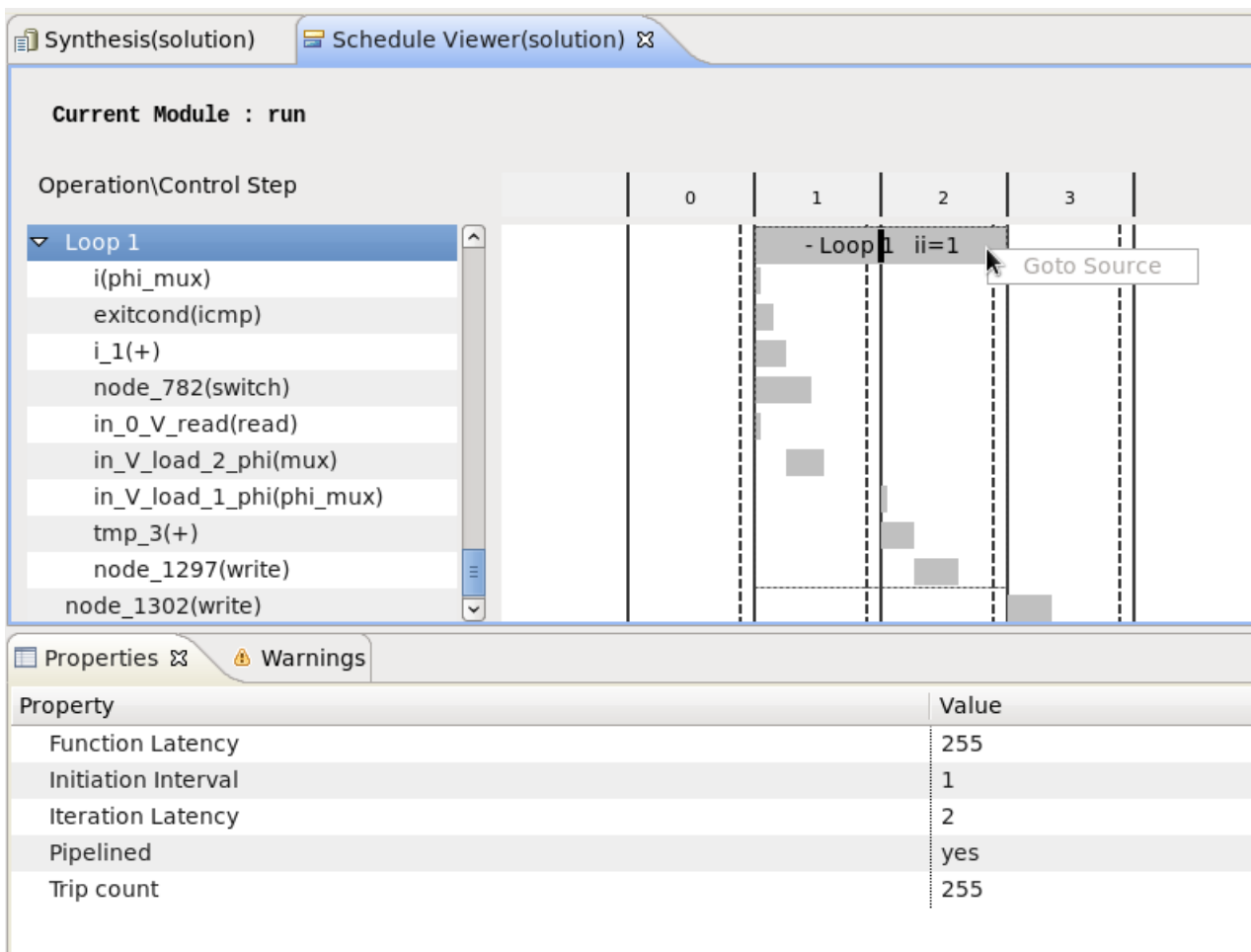
**Note:** The OpenCL API has an additional way of specifying loop pipelining. This has to do with the fact that work item loops are not explicitly stated and pipelining these loops requires the attribute:

```
__attribute__((xcl_pipeline_workitems))
```

More details to any of these specifications are provided in the *SDx Pragma Reference Guide* (UG1253) and the *SDAccel Environment Programmers Guide* (UG1277).

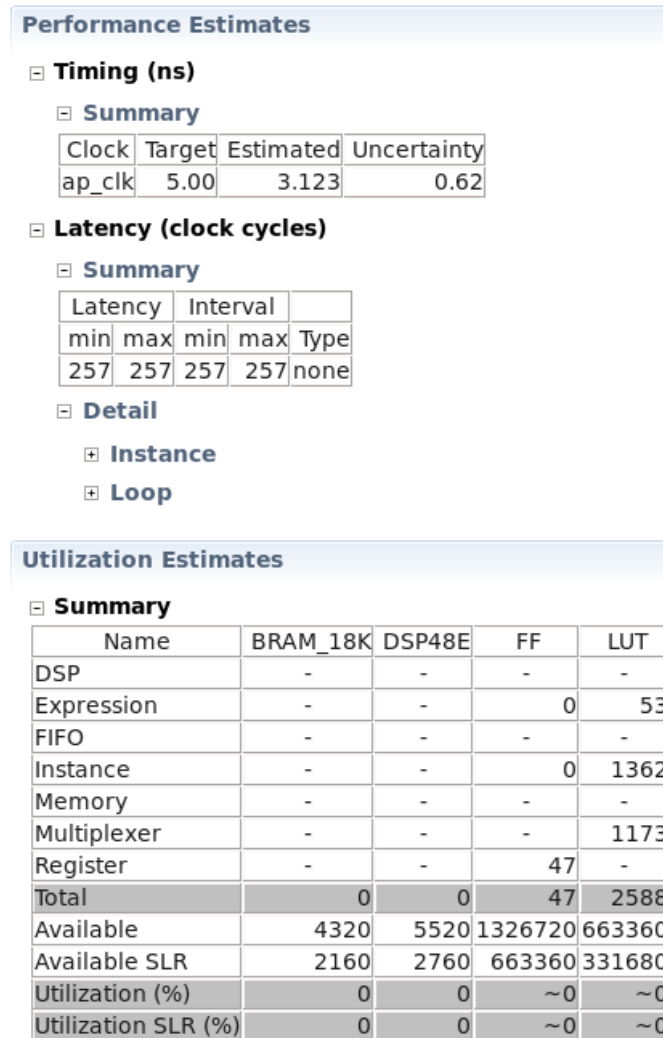
In this example, the Schedule Viewer in the HLS Project produces the following information:

Figure 26: Pipelining Loops in Schedule Viewer



With the overall estimates being:

Figure 27: Performance Estimates



Because each iteration of a loop only consumes two cycles of latency, there can only be a single iteration overlap. This enables the total latency to be cut into half compared to the original, resulting in 257 cycles of total latency. However, when compared to unrolling, this reduction in latency was achieved using fewer resources.

In most cases, loop pipelining by itself can improve overall performance. However, the effectiveness of the pipelining will depend on the structure of the loop. Some common limitations are:

- Resources with limited availability such as memory ports or process channels can limit the overlap of the iterations (II).
- Similarly, loop-carried dependencies, such as those created by variables conditions computed in one iteration affecting the next, might increase the initial interval of the pipeline.

These are reported by the tool during high-level synthesis and can be observed and examined in the Schedule Viewer. For the best possible performance, the code might have to be modified to eliminate these limiting factors, or the tool needs to be instructed to eliminate some dependency by restructuring the memory implementation of an array or breaking the dependencies all together.

## Task Parallelism

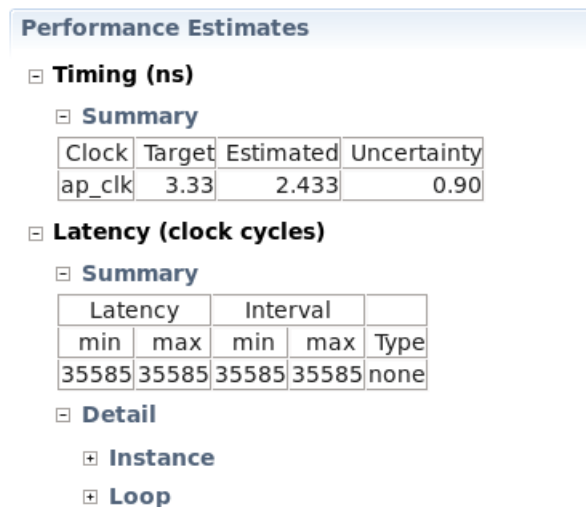
Task parallelism allows you to take advantage of data flow parallelism. In contrast to loop parallelism, when task parallelism is deployed, full execution units (tasks) are allowed to operate in parallel taking advantage of extra buffering introduced between the tasks.

Look at the following example:

```
void run (ap_uint<16> in[1024],
         ap_uint<16> out[1024]
        ) {
    ap_uint<16> tmp[128];
    for(int i = 0; i<8; i++) {
        processA(&(in[i*128]), tmp);
        processB(tmp, &(out[i*128]));
    }
}
```

When this code is executed, the function processA and processB are executed sequentially 128 times in a row. Given the combined latency for processA and processB in the loop is 278, the total latency can be estimated as:

Figure 28: Performance Estimates



The extra cycle is due to loop setup and can be observed in the Schedule Viewer.

For C/C++ code, task parallelism is performed by adding the DATAFLOW pragma into the for-loop:

```
#pragma HLS DATAFLOW
```

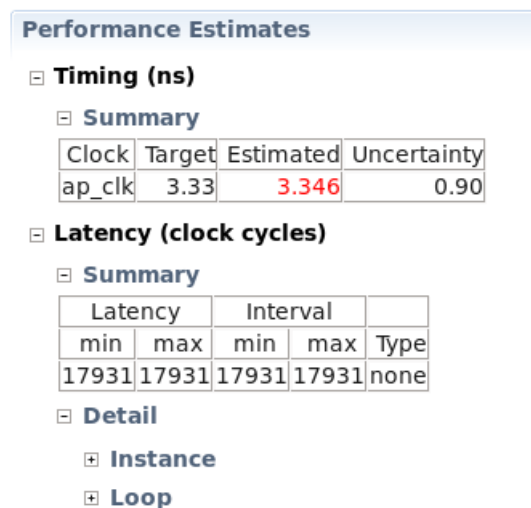
For OpenCL API code, add the attribute before the for-loop:

```
__attribute__((xcl_dataflow))
```

Refer to *SDx Pragma Reference Guide (UG1253)* and *SDAccel Environment Programmers Guide (UG1277)* for more details regarding the specifics and limitations of these modifiers.

As illustrated by the estimates in the HLS Report, applying the transformation will considerably improve the overall performance effectively using a double (ping pong) buffer scheme between the tasks:

Figure 29: Performances Estimates

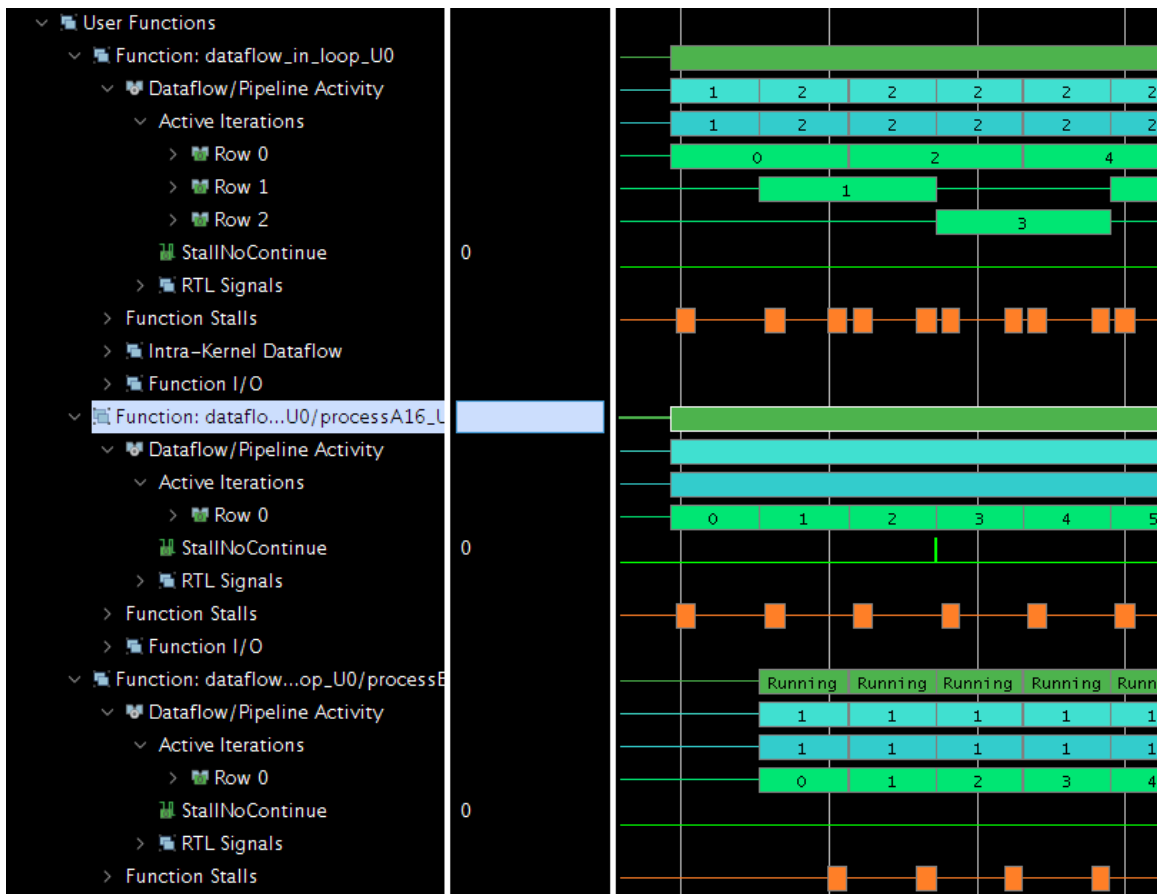


The overall latency of the design has almost halved in this case due to concurrent execution of the different tasks of the different iterations. Given the 139 cycles per processing function and the full overlap of the 128 iterations, this allows the total latency to be:

```
(1x only processA + 127x both processes + 1x only processB) * 139 cycles = 17931 cycles
```

Using task parallelism is a very powerful way improve performance when it comes to implementation. However, the effectiveness of applying the DATAFLOW pragma to a specific and arbitrary piece of code might vary vastly. The coding guidelines for applying DATAFLOW effectively are provided in *SDx Pragma Reference Guide (UG1253)* and *SDAccel Environment Programmers Guide (UG1277)*. However, to understand the final implementation of the DATAFLOW pragma, it is often necessary to actually look at the execution pattern of the individual tasks. Towards that end, the SDAccel environment provides the Detailed Kernel Trace, which nicely illustrates concurrent execution.

Figure 30: Detailed Kernel Trace



For this detailed kernel trace, the tool displays the start of the dataflowed loop, as shown in the previous figure. It illustrates how processA is starting up right away with the beginning of the loop, while processB waits until the completion of the processA before it can start up its first iteration. However, while processB completes the first iteration of the loop, processA begins operating on the second iteration and so forth.

A more abstract representation of this information is presented in the Application Timeline (Host & Device) and Device Hardware Transaction View (device-only during hardware emulation).



---

# Optimizing Compute Units

## Data Width

One aspect for performance is the data width required for the implementation. The tool propagates port widths throughout the algorithm. In some cases, especially when starting out with an algorithmic description, the C/C++/OpenCL™ API code might only use large data types such as integers even at the ports of the design. However, as the algorithm gets mapped to a fully configurable implementation, smaller data types such as 10- or 12-bit might often suffice. Towards that end it is beneficial to check the size of basic operations in the HLS Synthesis report during optimization. In general, when the SDAccel environment maps an algorithm onto the FPGA, much processing is required to comprehend the C/C++/OpenCL API structure and extract operational dependencies. Therefore, to perform this mapping the SDAccel environment generally partitions the source code into operational units which are then mapped onto the FPGA. Several aspects influence the number and size of these operational units (ops) as seen by the tool.

In the following figure, the basic operations and their bitwidth are reported.

Figure 31: Operations Utilization Estimates

| Utilization Estimates                       |           |          |           |            |             |             |
|---|-----------|----------|-----------|------------|-------------|-------------|
| <input type="checkbox"/> <b>Summary</b>     |           |          |           |            |             |             |
| Name  | BRAM_18K  | DSP48E   | FF        | LUT        |             |             |
| DSP   | -         | -        | -         | -          |             |             |
| Expression                                  | -         | -        | 0         | 102        |             |             |
| FIFO  | -         | -        | -         | -          |             |             |
| Instance                                    | -         | -        | -         | -          |             |             |
| Memory                                      | 0         | -        | 24        | 12         |             |             |
| Multiplexer                                 | -         | -        | -         | 80         |             |             |
| Register                                    | -         | -        | 51        | -          |             |             |
| <b>Total</b>                                | <b>0</b>  | <b>0</b> | <b>75</b> | <b>194</b> |             |             |
| Available                                   | 4320      | 5520     | 1326720   | 663360     |             |             |
| Available SLR                               | 2160      | 2760     | 663360    | 331680     |             |             |
| Utilization (%)                             | 0         | 0        | ~0        | ~0         |             |             |
| Utilization SLR (%)                         | 0         | 0        | ~0        | ~0         |             |             |
| <input type="checkbox"/> <b>Detail</b>      |           |          |           |            |             |             |
| <input type="checkbox"/> <b>Instance</b>    |           |          |           |            |             |             |
| <input type="checkbox"/> <b>DSP48</b>       |           |          |           |            |             |             |
| <input type="checkbox"/> <b>Memory</b>      |           |          |           |            |             |             |
| <input type="checkbox"/> <b>FIFO</b>        |           |          |           |            |             |             |
| <input type="checkbox"/> <b>Expression</b>  |           |          |           |            |             |             |
| Variable Name                               | Operation | DSP48E   | FF        | LUT        | Bitwidth P0 | Bitwidth P1 |
| i_1_fu_124_p2                               | +         | 0        | 0         | 11         | 3           | 1           |
| i_2_fu_148_p2                               | +         | 0        | 0         | 15         | 7           | 1           |
| i_3_fu_179_p2                               | +         | 0        | 0         | 15         | 7           | 1           |
| sum_i9_fu_194_p2                            | +         | 0        | 0         | 15         | 8           | 8           |
| sum_i_fu_158_p2                             | +         | 0        | 0         | 15         | 8           | 8           |
| exitcond_fu_118_p2                          | icmp      | 0        | 0         | 9          | 3           | 4           |
| exitcond_i6_fu_173_p2                       | icmp      | 0        | 0         | 11         | 7           | 8           |
| exitcond_i_fu_142_p2                        | icmp      | 0        | 0         | 11         | 7           | 8           |
| <b>Total</b>                                |           | <b>8</b> | <b>0</b>  | <b>0</b>   | <b>102</b>  | <b>39</b>   |
| <input type="checkbox"/> <b>Multiplexer</b> |           |          |           |            |             |             |
| <input type="checkbox"/> <b>Register</b>    |           |          |           |            |             |             |

Look for bit widths of 16, 32, and 64 bits commonly used in algorithmic descriptions, and verify that the associated operation from the C/C++/OpenCL API source actually requires the bit width to be this large. This can considerably improve the implementation of the algorithm, as smaller operations require less computation time.

## Fixed Point Arithmetic

Some applications use floating point computation only because they are optimized for other hardware architectures. As explained in [Deep Learning with INT8 Optimization on Xilinx Devices](#), using fixed point arithmetic for applications, like deep learning, can save the power efficiency and area significantly while keeping the same level of accuracy.



---

**RECOMMENDED:** Explore fixed point arithmetic for your application before committing to using floating point operations.

---

## Macro Operations

It is sometimes advantageous to think about larger computational elements. The tool will operate on the source code independently of the remaining source code, effectively mapping the algorithm without consideration of surrounding operations onto the FPGA. When applied, SDAccel tool keeps operational boundaries, effectively creating macro operations for specific code. This uses the following principles:

- Operational locality to the mapping process.
- Reduction in complexity for the heuristics.

This might create vastly different results when applied. In C/C++ macro operations are created with the help of

```
#pragma HLS inline off
```

While in the OpenCL API, the same kind of macro operation can be generated by *not* specifying the following attribute, when defining a function.:

```
__attribute__((always_inline))
```

## Using Optimized Libraries

The OpenCL specification provides many math built-in functions. All math built-in functions with the `native_` prefix are mapped to one or more native device instructions and will typically have better performance compared to the corresponding functions (without the `native_` prefix). The accuracy and in some cases the input ranges of these functions is implementation-defined. In SDAccel™ environment these `native_` built-in functions use the equivalent functions in the Vivado® HLS tool Math library, which are already optimized for Xilinx® FPGAs in terms of area and performance.



---

**RECOMMENDED:** Xilinx recommends that you use `native_` built-in functions or the HLS tool Math library if the accuracy meets the application requirement.

---

## Optimizing Memory Architecture

A key aspect of implementation is memory architecture. Because of the limited access bandwidth, it can heavily impact the overall performance, as shown in the following example:

```
void run (ap_uint<16> in[256][4],
         ap_uint<16> out[256]
        ) {
    ...
    ap_uint<16> inMem[256][4];
    ap_uint<16> outMem[256];

    ... Preprocess input to local memory

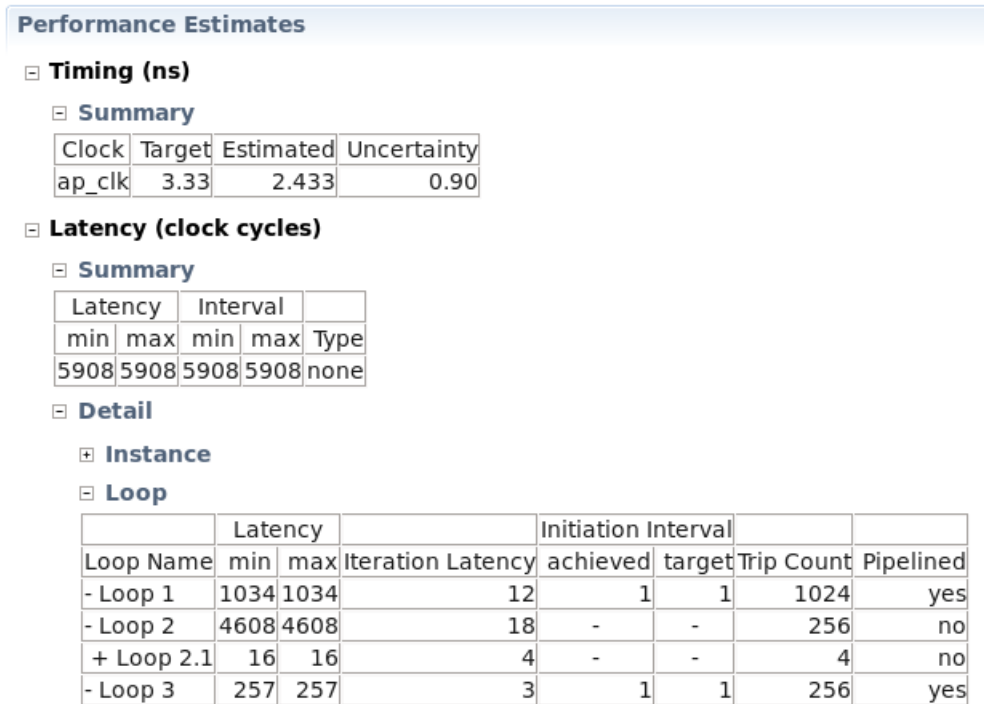
    for( int j=0; j<256; j++) {
        #pragma HLS PIPELINE OFF
        ap_uint<16> sum = 0;
        for( int i = 0; i<4; i++) {

            sum += inMem[j][i];
        }
        outMem[j] = sum;
    }

    ... Postprocess write local memory to output
}
```

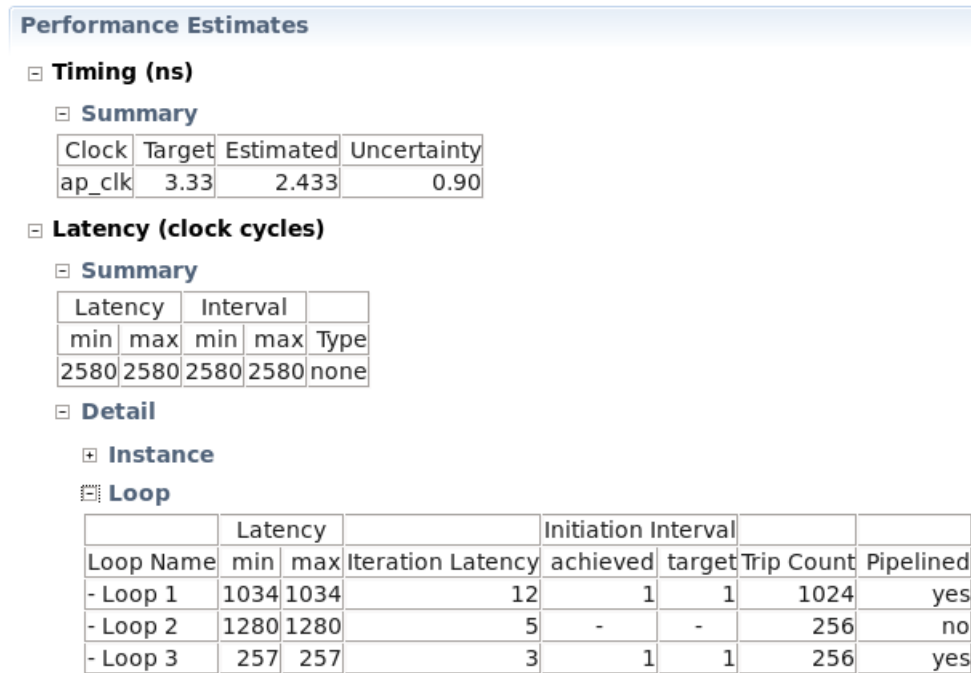
This code adds the four values associated with the inner dimension of the two dimensional input array. If implemented without any additional modifications, it results in the following estimates:

Figure 32: Performance Estimates



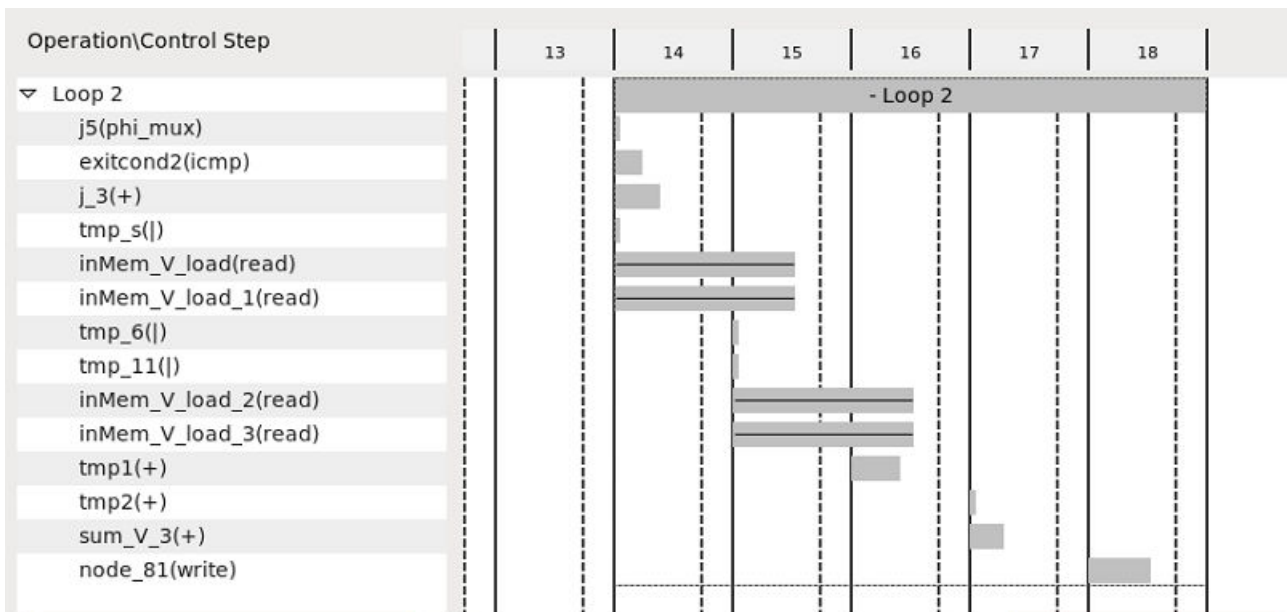
The overall latency of 4608 (Loop 2) is due to 256 iterations of 18 cycles (16 cycles spent in the inner loop, plus the reset of sum, plus the output being written). This can be observed in the Schedule Viewer in the HLS Project. The estimates become considerably better when unrolling the inner loop.

Figure 33: Performance Estimates



However, this improvement is largely due to the fact that this process uses both ports of a dual port memory. This can be seen from the Schedule Viewer in the HLS Project:

Figure 34: Schedule Viewer

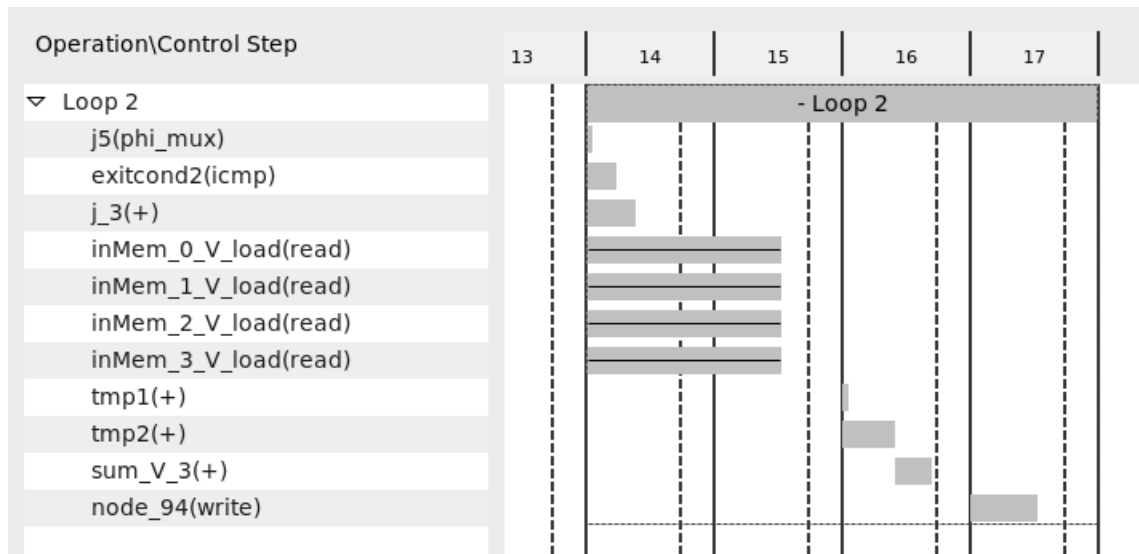


As you can see, two read operations are performed per cycle to access all the values from the memory to calculate the sum. This is often an undesired result as this completely blocks the access to the memory. To further improve the results, the memory can be split into four smaller memories along the second dimension:

```
#pragma HLS ARRAY_PARTITION variable=inMem complete dim=2
```

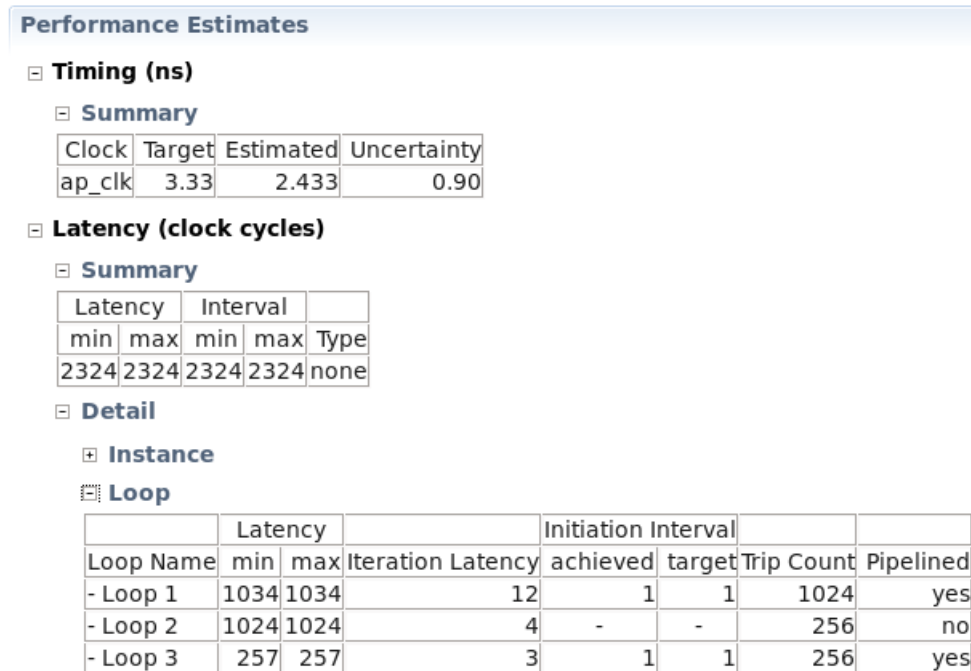
This results in four array reads, all executed on different memories using a single port:

Figure 35: Executed Four Arrays Results



Using a total of  $256 * 4$  cycles = 1024 cycles for loop 2.

Figure 36: Performance Estimates

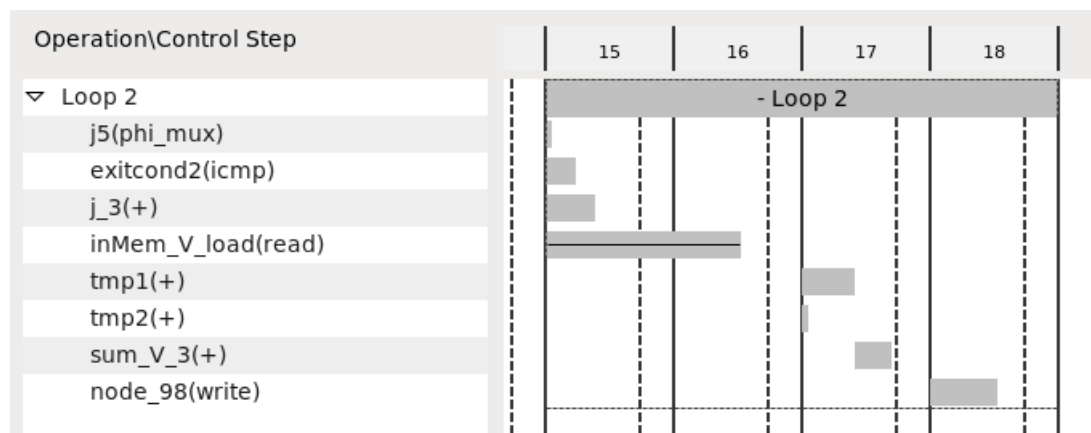


Alternatively, the memory can be reshaped into to a single memory with four words in parallel. This is performed through the following pragma:

```
#pragma HLS array_reshape variable=inMem complete dim=2
```

This results in the same latency as the array partitioning, but with a single memory using a single port:

Figure 37: Latency Result





Although, either solution creates comparable results with respect to overall latency and utilization, reshaping the array results in cleaner interfaces and less routing congestion making this the preferred solution.

**Note:** This completes array optimization, in a real design the latency could be further improved by exploiting loop parallelism (see the [Loop Parallelism](#) section).

```
void run (ap_uint<16> in[256][4],
         ap_uint<16> out[256]
        ) {
    ...

    ap_uint<16> inMem[256][4];
    ap_uint<16> outMem[256];
    #pragma HLS array_reshape variable=inMem complete dim=2

    ... Preprocess input to local memory

    for( int j=0; j<256; j++) {
        #pragma HLS PIPELINE OFF
        ap_uint<16> sum = 0;
        for( int i = 0; i<4; i++) {
            #pragma HLS UNROLL
            sum += inMem[j][i];
        }
        outMem[j] = sum;
    }

    ... Postprocess write local memory to output
}
```

# Host Optimization

This chapter focuses on host code optimization. The host code uses the OpenCL™ API to schedule the individual compute unit executions and data transfers from and to the FPGA board. As a result, you need to think about concurrent execution through the OpenCL queue(s). This section discusses in detail common pitfalls and how to recognize and address them.

---

## Reducing Overhead of Kernel Enqueuing

The OpenCL API execution model supports data parallel and task parallel programming models. Kernels are usually enqueued by the OpenCL runtime multiple times and then scheduled for execution on the device. You must send the command to start the kernel in one of two ways:

- For the data parallel case, use the `clEnqueueNDRange` API.
- For the task parallel case, use the `clEnqueueTask` API.

The dispatching process is executed on the host processor, the actual commands, and kernel arguments must to be sent to the FPGA through the PCIe® link. In the current Xilinx runtime (XRT), the overhead of dispatching the command and arguments to the FPGA is between 30us and 60us, depending the number of arguments on the kernel. You can reduce the impact of this overhead by minimizing the number of times the kernel needs to be executed.

For the data parallel case, Xilinx recommends that you carefully choose the global and local work sizes for your host code and kernel so that the global work size is a small multiple of the local work size. Ideally, the global work size is the same as the local work size as shown in the following code snippet:

```
size_t global = 1;
size_t local = 1;
clEnqueueNDRangeKernel(world.command_queue, kernel, 1, nullptr,
                        &global, &local, 2, write_events.data(),
                        &kernel_events[0]);
```

For the task parallel case, Xilinx recommends that you minimize the calls to `clEnqueueTask`. Ideally, you should finish all the work load in a single call to `clEnqueueTask`.

# Data Transfers

## Memory Data Transfer Types

The FPGA memory hierarchy, as well as the PCIe-based host to compute unit data transfers, provide a vast set of different data transfer options. To achieve maximum performance, these options should be reviewed. The basic trade-offs are between:

- Kind of memory to be deployed
- Streaming data transfer
- No host memory data transfers between compute units

**Note:** For more information about using the `--sp` option, the `--sc` option, `xclbin --info` command, and the `platforminfo` command, described in the sections below, refer to *SDx Command and Utility Reference Guide (UG1279)*.

### Memory Hierarchy

FPGA accelerator cards include several different memory hierarchies available for communication between host and kernels (CUs). For each memory hierarchy, there are different advantages and disadvantages to consider for optimization purposes.

- **DDR Memory:** This memory is external to the FPGA. It is the largest of the memory options but as a result has the longest latency to access.
  - **Usage (link flag):** `-sp [kernel|cu].[arg|port]:sptag`  
Values for `sptag` can be looked up through `xclbin --info` or the `platforminfo` command.
- **PL Memory:** This memory is internal to the FPGA. It is usually smaller than the external DDR memory, however, if supported by the platform, it has lower latency than the DDR memory.
  - **Usage (link flag):** `-sp [kernel|cu].[arg|port]:sptag`  
Values for `sptag` can be looked up through `xclbin.info` and `platform info`.

### Streaming Data Transfer

When streams are deployed between compute units or a compute unit and the host, a sequential data path is established. As a result, the transmitted data can only be received in the order sent. The tool distinguishes two types of streaming data transfers, one between compute units and the one between compute units and host.

- **Kernel-to-Kernel:** In case of kernel to kernel streaming communication, the ports are connected through the `-sc` option:

Usage (link flag): `-sc src.port:dst:port`

This is supported in the SDAccel environment through the use of `hls::streams` based on the `ap_axiu` data type on the kernel interfaces. This requires the `ap_axi_sdata.h` header file to be included.

- **Compute Unit to Host:** Streaming data between compute units and host require QDMA queues. As a result, this flow requires a QDMA platform.

This is supported in the SDAccel environment through the use of `hls::streams` based on the `qdma_axis` data type on the kernel interfaces. The member field “last” is used to indicate when a specified workload is completed. This allows streaming communication of various length. For more details on host code modeling, refer to the *SDAccel Environment Programmers Guide* ([UG1277](#)).

## No Host Memory Data Transfers between Compute Units

In memory mapped data transfers, blocks of data are written to and from the memory associated with the port. This enables the reader and writer to perform random access within the memory buffer used for communication. The SDAccel environment allows efficient buffer-based communication between kernels by using shared buffers between processes, which requires the same memory, or with the help of the `clEnqueueCopyBuffer` API, which implements efficient copying between buffers without host memory interaction.

## Peer-to-Peer

The SDAccel development environment can support systems with multiple accelerator cards on a single host to work together to accelerate large software systems. Specifically, for communication between different accelerator cards, the development environment supports direct peer-to-peer communication. This can be achieved by enabling the direct access of one of the cards DDR memory space which enables direct communication (without host memory).

## Enable Direct Access to DDR Memory Space

The following steps describe the setup from a local source buffer to an exported destination buffer. However, before actually modifying the host code, it is important to enable peer-to-peer communication. This is performed through the use of `xutil "xutil p2p --enable"` for the receiving device. The whole DDR address space of the device will be mapped to the host I/O memory space. For more information, refer to the *SDx Command and Utility Reference Guide* ([UG1279](#)).

After that it is necessary to modify the host code to prepare for direct peer-to-peer communication:

1. Create `buf_dst` with `XCL_MEM_EXT_P2P_BUFFER` flag.

There should not be any associated user space buffer (host buffer) for `buf_dst`.

2. Export and import the `buf_dst`, using the `xclGetMemObjectFd` and `xclGetMemObjectFromFd` APIs, in context: `buf_dst_exported`.
3. Use the regular `clEnqueueCopyBuffer` API (`src_command_queue, buf_src, buf_dst_exported, 0,0, buffer_size,,,`) command to copy the buffers between the devices.

**Note:** As mentioned in the Step 1 and Step 2, the exported buffer should be a p2p buffer. If a non-p2p buffer is exported and a copy operation is executed, the buffer would still be copied but through the host memory and not directly from card-to-card. The card-to-card transfers will be reported separately in the Profile Summary report.

More information about peer-to-peer connectivity can be found in the [XRT](#) documentation on GitHub.

## Overlapping Data Transfers with Kernel Computation

Applications, such as database analytics, have a much larger data set than the available memory on the acceleration device. They require the complete data to be transferred and processed in blocks. Techniques that overlap the data transfers with the computation are critical to achieve high performance for these applications.

The following code snippet shows the vector add kernel from the OpenCL Overlap Data Transfers with Kernel Computation Example in the [host](#) category from [Xilinx On-boarding Example GitHub](#).

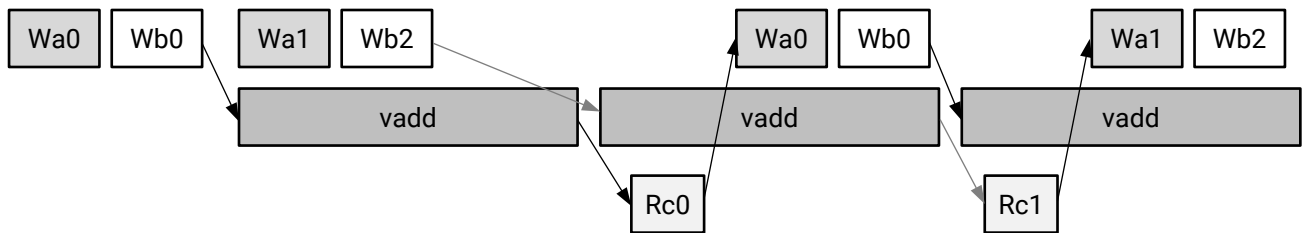
```
kernel __attribute__((reqd_work_group_size(1, 1, 1)))
void vadd(global int* c,
          global const int* a,
          global const int* b,
          const int offset,
          const int elements)
{
    int end = offset + elements;
    vadd_loop: for (int x=offset; x<end; ++x) {
        c[x] = a[x] + b[x];
    }
}
```

For this example, there are four tasks to perform in the host application:

1. Write buffer a (Wa)
2. Write buffer b (Wb)
3. Execute vadd kernel
4. Read buffer c (Rc)

The asynchronous nature of OpenCL data transfer and kernel execution APIs allows overlap of data transfers and kernel execution, as shown in the following figure. In this example, double buffering is used for all buffers so that the compute unit can process one set of buffers while the host can operate on the other set of buffers. The OpenCL event object provides an easy way to set up complex operation dependencies and synchronize host threads and device operations. The arrows in the following figure show how event triggering can be set up to achieve optimal performance.

Figure 38: Event Triggering Set Up



X22780-042519

The following host code snippet enqueues the four tasks in a loop. It also sets up event synchronization between different tasks to ensure that data dependencies are met for each task. The double buffering is set up by passing different memory objects values to the `clEnqueueMigrateMemObjects` API. The event synchronization is achieved by having each API call wait for other event as well as trigger its own event when the API completes.

```
for (size_t iteration_idx = 0;
     iteration_idx < num_iterations;
     iteration_idx++) {
    int flag = iteration_idx % 2;

    if (iteration_idx >= 2) {
        clWaitForEvents(1, &map_events[flag]);
        OCL_CHECK(clReleaseMemObject(buffer_a[flag]));
        OCL_CHECK(clReleaseMemObject(buffer_b[flag]));
        OCL_CHECK(clReleaseMemObject(buffer_c[flag]));
        OCL_CHECK(clReleaseEvent(read_events[flag]));
        OCL_CHECK(clReleaseEvent(kernel_events[flag]));
    }

    buffer_a[flag] = clCreateBuffer(world.context,
                                    CL_MEM_READ_ONLY | CL_MEM_USE_HOST_PTR,
                                    bytes_per_iteration,
                                    &A[iteration_idx * elements_per_iteration],
                                    NULL);
    buffer_b[flag] = clCreateBuffer(world.context,
                                    CL_MEM_READ_ONLY | CL_MEM_USE_HOST_PTR,
                                    bytes_per_iteration,
                                    &B[iteration_idx * elements_per_iteration],
                                    NULL);
    buffer_c[flag] = clCreateBuffer(world.context,
                                    CL_MEM_WRITE_ONLY | CL_MEM_USE_HOST_PTR,
                                    bytes_per_iteration,
                                    &device_result[iteration_idx * elements_per_iteration],
                                    NULL);

    array<cl_event, 2> write_events;
```

```

printf("Enqueueing Migrate Mem Object (Host to Device) calls\n");
// These calls are asynchronous with respect to the main thread
// because are passing the CL_FALSE as the third parameter.
// Because we are passing the events from the previous kernel call
// into the wait list, it will wait for the previous operations
// to complete before continuing
OCL_CHECK(clEnqueueMigrateMemObjects(
    world.command_queue, 1, &buffer_a[iteration_idx % 2],
    0 /* flags, 0 means from host */,
    0, NULL,
    &write_events[0]));
set_callback(write_events[0], "ooo_queue");

OCL_CHECK(clEnqueueMigrateMemObjects(
    world.command_queue, 1, &buffer_b[iteration_idx % 2],
    0 /* flags, 0 means from host */,
    0, NULL,
    &write_events[1]));
set_callback(write_events[1], "ooo_queue");

OCL_CHECK(clSetKernelArg(kernel, 0, sizeof(cl_mem),
    &buffer_c[iteration_idx % 2]));
OCL_CHECK(clSetKernelArg(kernel, 1, sizeof(cl_mem),
    &buffer_a[iteration_idx % 2]));
OCL_CHECK(clSetKernelArg(kernel, 2, sizeof(cl_mem),
    &buffer_b[iteration_idx % 2]));
OCL_CHECK(clSetKernelArg(kernel, 3, sizeof(int),
    &elements_per_iteration));

printf("Enqueueing NDRange kernel.\n");
// This event needs to wait for the write buffer operations to complete
// before executing. We are sending the write_events into its wait list
// to ensure that the order of operations is correct.
OCL_CHECK(clEnqueueNDRangeKernel(world.command_queue, kernel, 1,
    nullptr, &global, &local, 2,
    write_events.data(),
    &kernel_events[flag]));
set_callback(kernel_events[flag], "ooo_queue");

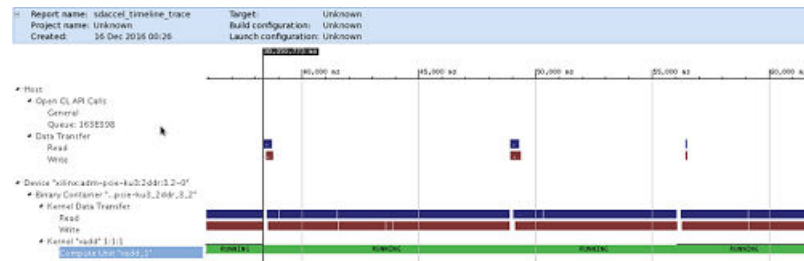
printf("Enqueueing Migrate Mem Object (Device to Host) calls\n");
// This operation only needs to wait for the kernel call. This call will
// potentially overlap the next kernel call as well as the next read
// operations
OCL_CHECK(clEnqueueMigrateMemObjects(world.command_queue, 1,
    &buffer_c[iteration_idx % 2],
    CL_MIGRATE_MEM_OBJECT_HOST, 1,
    &kernel_events[flag],
    &read_events[flag]));

set_callback(read_events[flag], "ooo_queue");
clEnqueueMapBuffer(world.command_queue, buffer_c[flag], CL_FALSE,
    CL_MAP_READ, 0, bytes_per_iteration, 1,
    &read_events[flag], &map_events[flag], 0);
set_callback(map_events[flag], "ooo_queue");

OCL_CHECK(clReleaseEvent(write_events[0]));
OCL_CHECK(clReleaseEvent(write_events[1]));
}
    
```

The Application Timeline view clearly shows that the data transfer time is completely hidden, while the `vadd_1` compute unit is running constantly.

Figure 39: Data Transfer Time Hidden in Application Timeline View



## Buffer Memory Segmentation

Allocation and deallocation of memory buffers can lead to memory segmentation in the DDRs. This might result in sub-optimal performance of compute units, even if they could theoretically execute in parallel.

This problem occurs most often when multiple pthreads for different compute units are used, and the threads allocate and release many device buffers with different sizes every time they enqueue the kernels. In this case, the timeline trace will exhibit gaps between kernel executions and it just seems the processes are sleeping.

Each buffer allocated by runtime should be continuous in hardware. For large memory, it might take a lot of time to wait for that space to be freed, when many buffers are allocated and deallocated. This can be resolved by allocating device buffer, and reusing it between different enqueues of a kernel.

## Compute Unit Scheduling

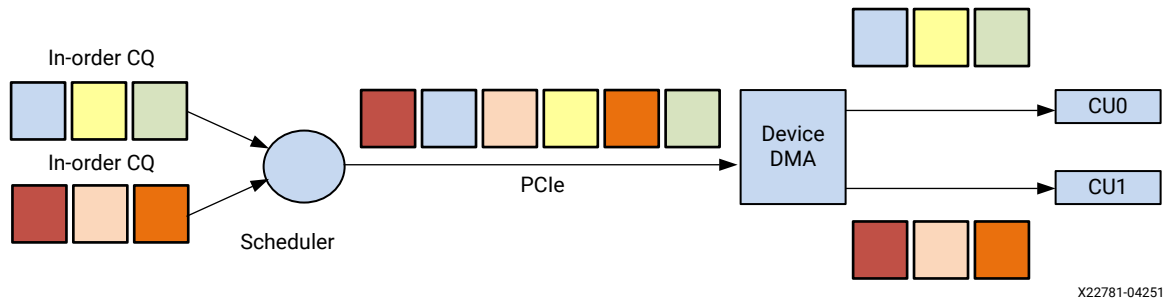
Scheduling kernel operations is key to overall system performance. This becomes even more important when implementing multiple compute units (of the same kernel or of different kernels). This section examines the different command queues responsible for scheduling the kernels.

### Multiple In-Order Command Queues

The following figure shows an example with two in-order command queues, CQ0 and CQ1. The scheduler dispatches commands from each queue in order, but commands from CQ0 and CQ1 can be pulled out by the scheduler in any order. You must manage synchronization between CQ0 and CQ1, if required.



Figure 40: Example with Two In-Order Command Queues



X22781-042519

The following code snippet from the *Concurrent Kernel Execution Example* in `host` category from [SDAccel Getting Started Examples](#) on GitHub sets up multiple in-order command queues and enqueues commands into each queue:

```

cl_command_queue ordered_queue1 = clCreateCommandQueue(
    world.context, world.device_id, CL_QUEUE_PROFILING_ENABLE, &err)

cl_command_queue ordered_queue2 = clCreateCommandQueue(
    world.context, world.device_id, CL_QUEUE_PROFILING_ENABLE, &err);

clEnqueueNDRangeKernel(ordered_queue1, kernel_mscale, 1, offset,
    global, local, 0, nullptr,
    &kernel_events[0]);

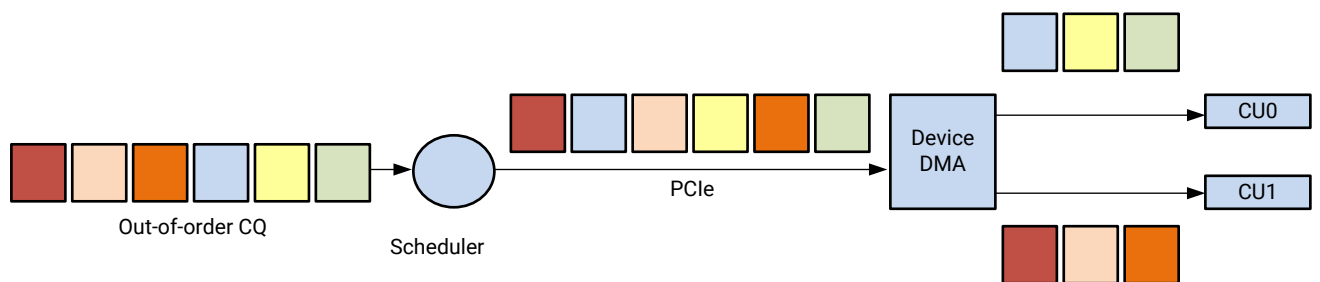
clEnqueueNDRangeKernel(ordered_queue1, kernel_madd, 1, offset,
    global, local, 0, nullptr,
    &kernel_events[1]);

clEnqueueNDRangeKernel(ordered_queue2, kernel_mmult, 1, offset,
    global, local, 0, nullptr,
    &kernel_events[2]);
    
```

## Single Out-of-Order Command Queue

The following figure shows an example with a single out-of-order command queue. The scheduler can dispatch commands from the queue in any order. You must set up event dependencies and synchronizations explicitly, if required.

Figure 41: Example with Single Out-of-Order Command Queue



X22783-042519

The following code snippet from the *Concurrent Kernel Execution Example* from [SDAccel Getting Started Examples](#) on GitHub sets up a single out-of-order command queue and enqueues commands:

```

cl_command_queue ooo_queue = clCreateCommandQueue(
    world.context, world.device_id,
    CL_QUEUE_PROFILING_ENABLE | CL_QUEUE_OUT_OF_ORDER_EXEC_MODE_ENABLE,
    &err);

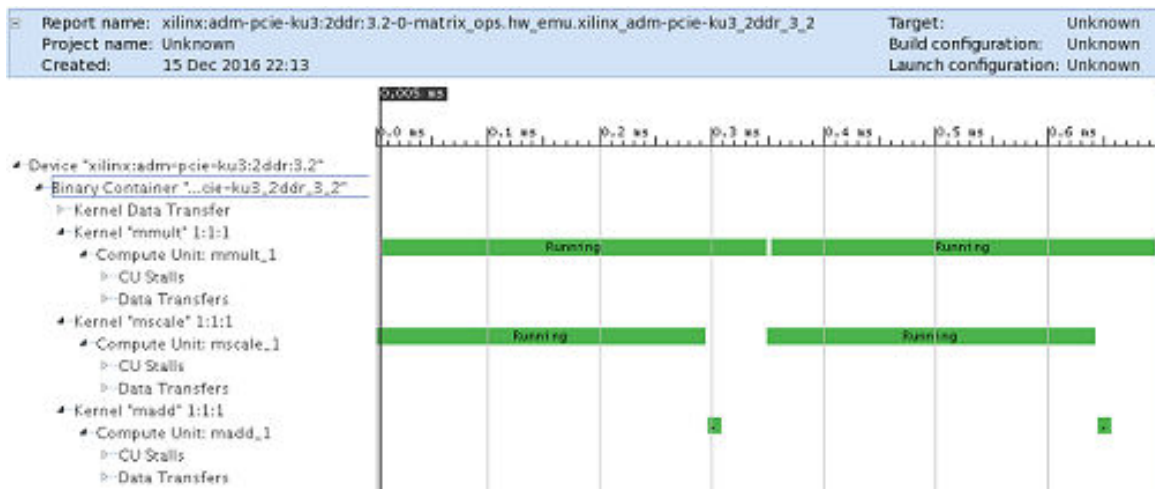
clEnqueueNDRangeKernel(ooo_queue, kernel_mscale, 1, offset, global,
    local, 0, nullptr, &ooo_events[0]);

clEnqueueNDRangeKernel(ooo_queue, kernel_madd, 1, offset, global,
    local, 1,
    &ooo_events[0], // Event from previous call
    &ooo_events[1]);

clEnqueueNDRangeKernel(ooo_queue, kernel_mmult, 1, offset, global,
    local, 0,
    nullptr, // Does not depend on previous call
    &ooo_events[2])
    
```

The following figure shows the Application Timeline view where that the compute unit `mmult_1` is running in parallel with the compute units `mscale_1` and `madd_1`, using both multiple in-order queues and single out-of-order queue methods.

**Figure 42: Application Timeline View Showing `mult_1` Running with `mscale_1` and `madd_1`**



## Using the `clEnqueueMigrateMemObjects` API to Transfer Data

The OpenCL framework provides a number of APIs for transferring data between the host and the device. Typically, data movement APIs, such as `clEnqueueWriteBuffer` and `clEnqueueReadBuffer`, implicitly migrate memory objects to the device after they are enqueued. They do not guarantee when the data is transferred. This makes it difficult for the host application to overlap the placements of the memory objects onto the device with the computation carried out by kernels.

The OpenCL 1.2 framework introduced a new API, `clEnqueueMigrateMemObjects`. Using this API, memory migration can be explicitly performed ahead of the dependent commands. This allows the application to preemptively change the association of a memory object, through regular command queue scheduling, to prepare for another upcoming command. This also permits an application to overlap the placement of memory objects with other unrelated operations before these memory objects are needed, potentially hiding transfer latencies. After the event associated by the `clEnqueueMigrateMemObjects` API are marked `CL_COMPLETE`, the memory objects specified in `mem_objects` are successfully migrated to the device associated with `command_queue`.

The `clEnqueueMigrateMemObjects` API can also be used to direct the initial placement of a memory object after creation, possibly avoiding the initial overhead of instantiating the object on the first enqueued command to use it.

Another advantage of using the `clEnqueueMigrateMemObjects` API is that it can migrate multiple memory objects in a single API call. This reduces the overhead of scheduling and calling functions for transferring data for more than one memory object.

The following code snippet shows the usage of the `clEnqueueMigrateMemObjects` API from the *Vector Multiplication for XPR Device* example in the host category from [SDAccel Getting Started Examples](#) on GitHub.

```
int err = clEnqueueMigrateMemObjects(
    world.command_queue,
    1,
    &d_mul_c,
    CL_MIGRATE_MEM_OBJECT_HOST,
    0,
    NULL,
    NULL);
```

# Topological Optimization

This chapter focuses on the topological optimization, and looks at the attributes related to the rough layout and implementation of multiple compute units and their impact on performance.

---

## Multiple Compute Units

Depending on available resources on the FPGA, multiple compute units of the same kernel (or different kernels) can be created to run in parallel, which improves the system processing time and throughput.

Different kernels are provided as separate .xo files on the `xocc` link line. Multiple kernel compute units can be added by using the `--nk` option:

```
xocc -l --nk <kernel_name:number(:compute_unit_name1.compute_unit_name2...)>
```

**Note:** Each of the individual kernels must be individually driven by the host code.

---

## Using Multiple DDR Banks

In the SDAccel™ environment, supported acceleration cards provide one, two, or four DDR banks and up to 80 GB/s raw DDR bandwidth.



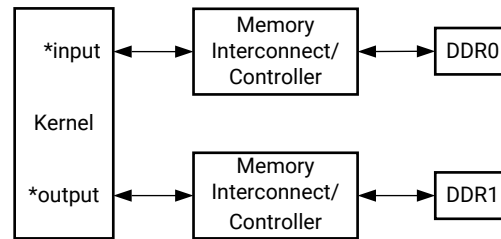
**RECOMMENDED:** For kernels moving large amount of data between the FPGA and the DDR, Xilinx recommends that you direct the SDAccel compiler and runtime library to use multiple DDR banks.

In addition to DDR banks, the host application can access PLRAM to transfer data directly to a kernel. This feature is enabled using the `xocc --sp` option with compatible platforms.

To take advantage of multiple DDR banks or PLRAMs, use the `--sp` option to map the individual arguments of the accelerator to the desired DDR banks or PLRAM in the xclbin. This mapping will be automatically picked up by the host executable.

The following block diagram shows the *Global Memory Two Banks Example* from the “[kernel\\_to\\_gmem](#)” category on [SDAccel Getting Started Examples](#) on GitHub that connects the input pointer to DDR bank 0 and output pointer to DDR bank 1.

Figure 43: Global Memory Two Banks



X22633-050619

## Connecting Kernel Ports to Memory Banks

### Creating Multiple AXI Interfaces

OpenCL™ kernels, C/C++ kernels, and RTL kernels have different methods for assigning function parameters to AXI interfaces.

- For OpenCL kernels, the `--max_memory_ports` option is required to generate one AXI4 interface for each global pointer on the kernel argument. The AXI4 interface name is based on the order of the global pointers on the argument list.

The following code is taken from the example `gmem_2banks_ocl` in the `kernel_to_gmem` category from the [SDAccel Getting Started Examples](#) on GitHub:

```

__kernel __attribute__((reqd_work_group_size(1, 1, 1)))
void apply_watermark(__global const TYPE * __restrict input,
__global TYPE * __restrict output, int width, int height) {
    ...
}
    
```

In this example, the first global pointer `input` is assigned an AXI4 name `xi_gmem0`, and the second global pointer `output` is assigned a name `axi_gmem1`.

- For C/C++ kernels, multiple AXI4 interfaces are generated by specifying different “bundle” names in the HLS INTERFACE pragma for different global pointers. For more information, refer to the *SDAccel Environment Programmers Guide* ([UG1277](#)).

The following is a code snippet from the `gmem_2banks_c` example that assigns the `input` pointer to the bundle `gmem0` and the `output` pointer to the bundle `gmem1`. The bundle name can be any valid C string, and the AXI4 interface name generated will be `m_axi_<bundle_name>`. For this example, the input pointer will have AXI4 interface name as `axi_gmem0`, and the output pointer will have `m_axi_gmem1`.

```

#pragma HLS INTERFACE m_axi port=input offset=slave bundle=gmem0
#pragma HLS INTERFACE m_axi port=output offset=slave bundle=gmem1
    
```

- For RTL kernels, the port names are generated during the import process by the RTL kernel wizard. The default names proposed by the RTL kernel wizard are `m00_axi` and `m01_axi`. If not changed, these names have to be used when assigning a DDR bank through the `--sp` option.

## Assigning AXI Interfaces to Global Memory



**IMPORTANT!** When using more than one DDR interface, Xilinx requires you to specify the DDR memory bank for each kernel/CU using the `--sp` option, and specify in which SLR the kernel is placed. Refer to the XOCC command in the *SDx Command and Utility Reference Guide (UG1279)* for details of the `--sp` command option, and the *SDAccel Environment User Guide (UG1023)* for details on SLR placement.

AXI4 interfaces are connected to DDR banks using the `--sp` option. The `--sp` option value is in the format of `<kernel_instance_name>.<interface_name>:<DDR_bank_name>`.

The complete list of DDR memories or alternative communication memories, such as PLRAM and HBM, can be found through the `platforminfo` command.

**Note:** Not all global memory options are supported by all platforms.

The following is the command line example that connects the input pointer (`M_AXI_GMEM0`) to DDR bank 0 and the output pointer (`M_AXI_GMEM1`) to DDR bank 1:

```
xocc --max_memory_ports apply_watermark
--sp apply_watermark_1.m_axi_gmem0:DDR[0]
--sp apply_watermark_1.m_axi_gmem1:DDR[1]
```

You can use the Device Hardware Transaction view to observe the actual DDR Bank communication, and to analyze DDR usage.

Figure 44: Device Hardware Transaction View Transactions on DDR Bank



## Assigning Kernels to SLR regions

Assigning ports to DDR banks requires that the kernel will have to be physically routed on the FPGA to connect to the assigned DDR. Currently, large FPGAs use stacked silicon devices with several Super Logic Regions (SLRs). By default, the SDAccel environment will place the compute units in the same SLR as the shell. This might not always be desirable, especially when specific DDR banks are used that might be in another SLR region. As a result, Xilinx recommends to use the `--slr` option to map kernels to be close to the used DDR memory. For example, the `apply_watermark_1` kernel can be mapped to SLR 1 by applying the following link option:

```
xocc -l --slr apply_watermark_1:SLR1
```

To better understand the platform attributes, such as the number of DDRs and SLR regions, use `platforminfo`. For more information, refer to the *SDx Command and Utility Reference Guide* ([UG1279](#)).

# Examples

To help you quickly get started with the SDAccel™ environment, the [SDAccel Getting Started Examples](#) on GitHub hosts many examples to demonstrate good design practices, coding guidelines, design pattern for common applications, and most importantly optimization techniques to maximize application performance. The on-boarding examples are divided into several main categories. Each category has various key concepts that are illustrated by individual examples in both OpenCL™ C and C/C++, when applicable. All examples include a Makefile for running software emulation, hardware emulation, and running on hardware, and a `README.md` file that explains the example in detail.



# Additional Resources and Legal Notices

---

## Xilinx Resources

For support resources such as Answers, Documentation, Downloads, and Forums, see [Xilinx Support](#).

---

## Documentation Navigator and Design Hubs

Xilinx<sup>®</sup> Documentation Navigator (DocNav) provides access to Xilinx documents, videos, and support resources, which you can filter and search to find information. DocNav is installed with the SDSoc™ and SDAccel™ development environments. To open it:

- On Windows, select **Start** → **All Programs** → **Xilinx Design Tools** → **DocNav**.
- At the Linux command prompt, enter `docnav`.

Xilinx Design Hubs provide links to documentation organized by design tasks and other topics, which you can use to learn key concepts and address frequently asked questions. To access the Design Hubs:

- In DocNav, click the **Design Hubs View** tab.
- On the Xilinx website, see the [Design Hubs](#) page.

**Note:** For more information on DocNav, see the [Documentation Navigator](#) page on the Xilinx website.

---

## References

1. *SDAccel Environment Release Notes, Installation, and Licensing Guide* ([UG1238](#))

2. *SDAccel Environment Profiling and Optimization Guide* ([UG1207](#))
3. *SDAccel Environment Getting Started Tutorial* ([UG1021](#))
4. [SDAccel™ Development Environment web page](#)
5. [Vivado® Design Suite Documentation](#)
6. *Vivado Design Suite User Guide: Designing IP Subsystems using IP Integrator* ([UG994](#))
7. *Vivado Design Suite User Guide: Creating and Packaging Custom IP* ([UG1118](#))
8. *Vivado Design Suite User Guide: Partial Reconfiguration* ([UG909](#))
9. *Vivado Design Suite User Guide: High-Level Synthesis* ([UG902](#))
10. *UltraFast Design Methodology Guide for the Vivado Design Suite* ([UG949](#))
11. *Vivado Design Suite Properties Reference Guide* ([UG912](#))
12. [Khronos Group web page](#): Documentation for the OpenCL standard
13. [Xilinx® Virtex® UltraScale+™ FPGA VCU1525 Acceleration Development Kit](#)
14. [Xilinx® Kintex® UltraScale™ FPGA KCU1500 Acceleration Development Kit](#)
15. [Xilinx® Alveo™ web page](#)

---

## Please Read: Important Legal Notices

The information disclosed to you hereunder (the "Materials") is provided solely for the selection and use of Xilinx products. To the maximum extent permitted by applicable law: (1) Materials are made available "AS IS" and with all faults, Xilinx hereby DISCLAIMS ALL WARRANTIES AND CONDITIONS, EXPRESS, IMPLIED, OR STATUTORY, INCLUDING BUT NOT LIMITED TO WARRANTIES OF MERCHANTABILITY, NON-INFRINGEMENT, OR FITNESS FOR ANY PARTICULAR PURPOSE; and (2) Xilinx shall not be liable (whether in contract or tort, including negligence, or under any other theory of liability) for any loss or damage of any kind or nature related to, arising under, or in connection with, the Materials (including your use of the Materials), including for any direct, indirect, special, incidental, or consequential loss or damage (including loss of data, profits, goodwill, or any type of loss or damage suffered as a result of any action brought by a third party) even if such damage or loss was reasonably foreseeable or Xilinx had been advised of the possibility of the same. Xilinx assumes no obligation to correct any errors contained in the Materials or to notify you of updates to the Materials or to product specifications. You may not reproduce, modify, distribute, or publicly display the Materials without prior written consent. Certain products are subject to the terms and conditions of Xilinx's limited warranty, please refer to Xilinx's Terms of Sale which can be viewed at <https://>

[www.xilinx.com/legal.htm#tos](http://www.xilinx.com/legal.htm#tos); IP cores may be subject to warranty and support terms contained in a license issued to you by Xilinx. Xilinx products are not designed or intended to be fail-safe or for use in any application requiring fail-safe performance; you assume sole risk and liability for use of Xilinx products in such critical applications, please refer to Xilinx's Terms of Sale which can be viewed at <https://www.xilinx.com/legal.htm#tos>.

### **AUTOMOTIVE APPLICATIONS DISCLAIMER**

AUTOMOTIVE PRODUCTS (IDENTIFIED AS "XA" IN THE PART NUMBER) ARE NOT WARRANTED FOR USE IN THE DEPLOYMENT OF AIRBAGS OR FOR USE IN APPLICATIONS THAT AFFECT CONTROL OF A VEHICLE ("SAFETY APPLICATION") UNLESS THERE IS A SAFETY CONCEPT OR REDUNDANCY FEATURE CONSISTENT WITH THE ISO 26262 AUTOMOTIVE SAFETY STANDARD ("SAFETY DESIGN"). CUSTOMER SHALL, PRIOR TO USING OR DISTRIBUTING ANY SYSTEMS THAT INCORPORATE PRODUCTS, THOROUGHLY TEST SUCH SYSTEMS FOR SAFETY PURPOSES. USE OF PRODUCTS IN A SAFETY APPLICATION WITHOUT A SAFETY DESIGN IS FULLY AT THE RISK OF CUSTOMER, SUBJECT ONLY TO APPLICABLE LAWS AND REGULATIONS GOVERNING LIMITATIONS ON PRODUCT LIABILITY.

### **Copyright**

© Copyright 2016-2019 Xilinx, Inc. Xilinx, the Xilinx logo, Alveo, Artix, Kintex, Spartan, Versal, Virtex, Vivado, Zynq, and other designated brands included herein are trademarks of Xilinx in the United States and other countries. OpenCL and the OpenCL logo are trademarks of Apple Inc. used by permission by Khronos. HDMI, HDMI logo, and High-Definition Multimedia Interface are trademarks of HDMI Licensing LLC. AMBA, AMBA Designer, Arm, ARM1176JZ-S, CoreSight, Cortex, PrimeCell, Mali, and MPCore are trademarks of Arm Limited in the EU and other countries. All other trademarks are the property of their respective owners.