

# Introduction to DPC++ Programming for FPGA

A part of the DPC++ Tutorial Series

Prof. Yan Luo

Acknowledgement

This work is supported by Intel Corporation 2020-2021



*Learning with Purpose*

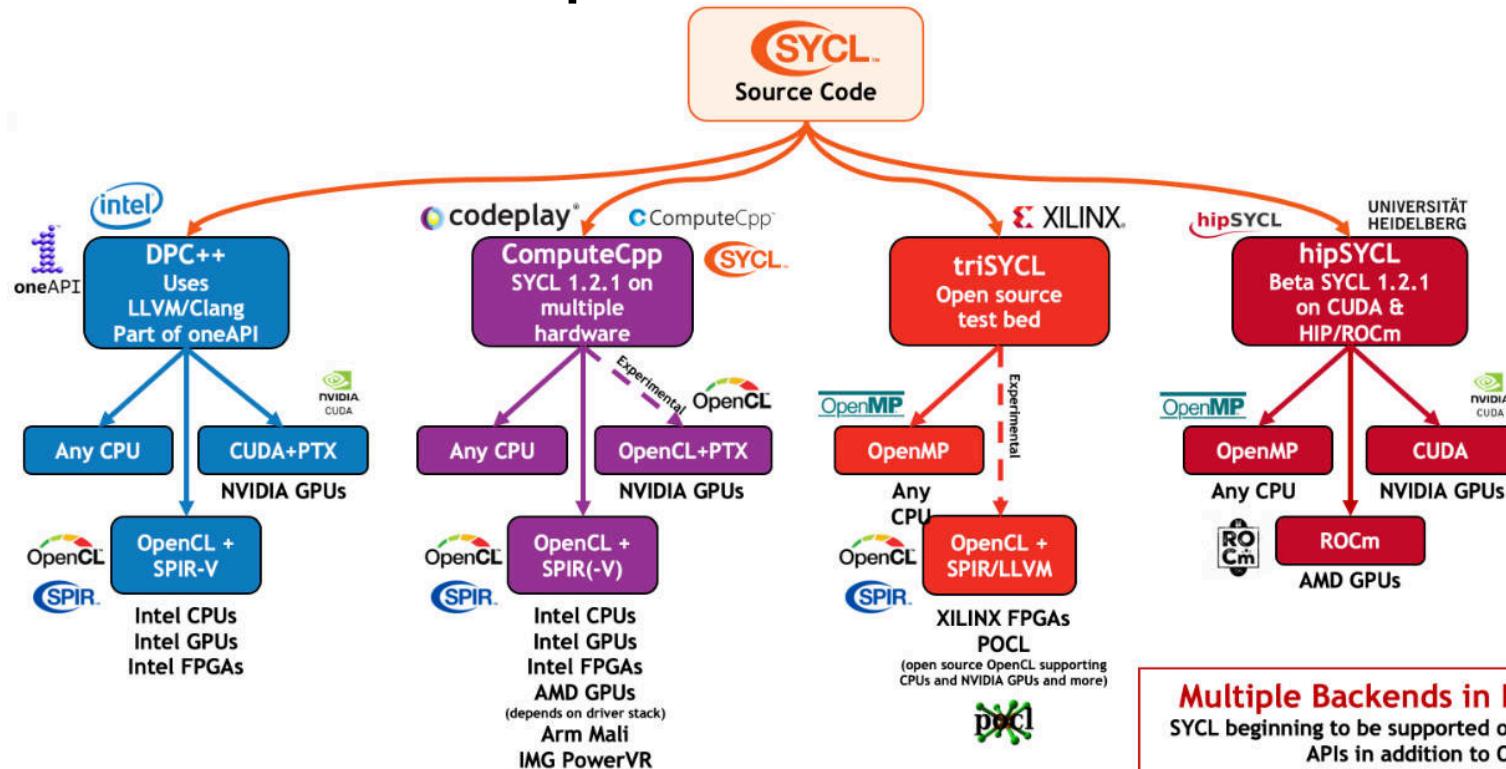
# Data Parallel C++

- A high-level language for data parallel programming
- Based on modern C++
- Single source for heterogeneous computing architectures
- Offloading computing to accelerators (e.g. FPGA and GPU)
- Speedup on data parallel workloads
  - Algorithm and parallelism analysis
  - Data/task decomposition
  - Architecture oriented performance optimization (e.g. for FPGA)



Learning with Purpose

# DPC++: an Implementation of SYCL



**Multiple Backends in Development**  
SYCL beginning to be supported on multiple low-level APIs in addition to OpenCL  
e.g. ROCm and CUDA  
For more information: <http://sycl.tech>

# A DPC++ Example

```
1 #include <CL/sycl.hpp>
2 #include <iostream>
3
4 constexpr int num=16;
5 using namespace sycl;
6
7 int main() {
8     auto r = range{num};
9     buffer<int> a{r};
10
11     queue{}.submit([&](handler& h) {
12         accessor out{a, h};
13         h.parallel_for(r, [=](item<1> idx) {
14             out[idx] = idx;
15         });
16     });
17
18     host_accessor result{a};
19     for (int i=0; i<num; ++i)
20         std::cout << result[i] << "\n";
21 }
```

device queue  
& command

kernel function  
on device

Header Files

namespace  
scope

SYCL buffer  
declaration

Lambda  
function

host access  
results in buffer

# Matrix Multiplication : How to Think in Parallel?

A part of the DPC++ Tutorial Series

Prof. Yan Luo

Acknowledgement

This work is supported by Intel Corporation 2020-2021



*Learning with Purpose*

# Matrix Multiplication

$$A \times B = C$$

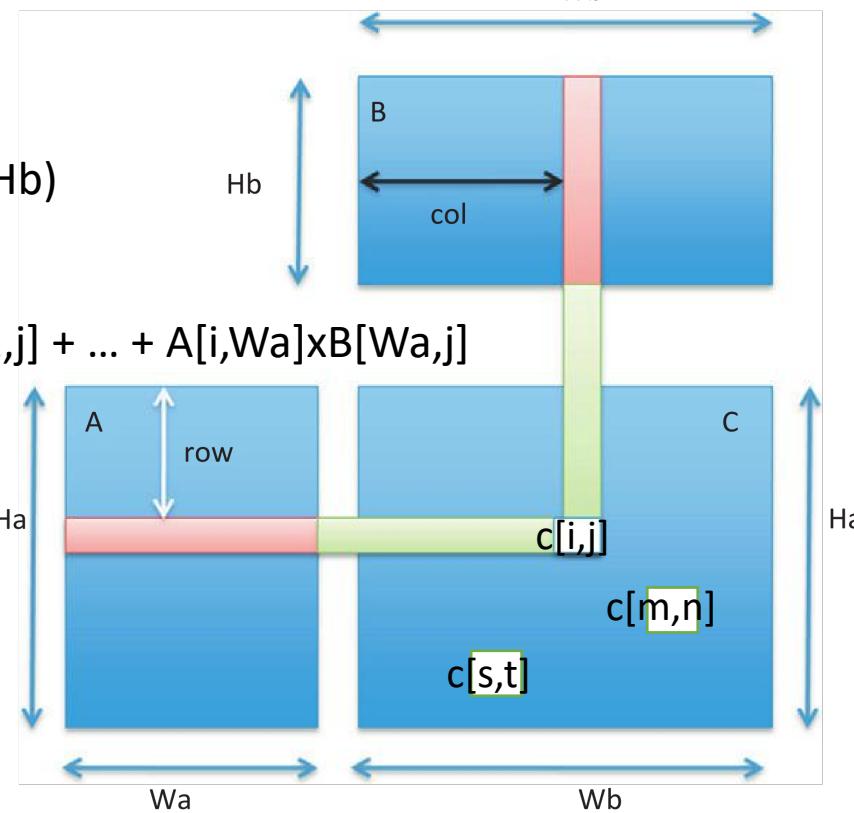
A is  $W_a \times H_a$

B is  $W_b \times H_b$

(Note:  $W_a = H_b$ )

C is  $W_b \times H_a$

$$C[i,j] = A[i,1] \times B[1,j] + A[i,2] \times B[2,j] + \dots + A[i,W_a] \times B[W_a,j]$$



# A C++ Implementation

```
// iterate over the rows of Matrix A
for (int i = 0; i < Ha; i++)
    // iterate over the columns of Matrix B
    for (int j = 0; i < Wb; j++) {
        C[i][j] = 0;
        // element-wise multiplication and accumulation
        for (int k = 0; k < Wa; k++)
            C[i][j] += A[i][k] * B[k][j];
    }
```



Learning with Purpose

# A DPC++ Implementation

```
q.submit([&](handler &h) {
    auto A = a_buf.get_access<access::mode::read>(h);
    auto B = b_buf.get_access<access::mode::read>(h);
    auto C = sum_buf.get_access<access::mode::write>(h);
    range<2> num_items{a_rows, b_columns};
```

create  
accessors

```
h.parallel_for<class MMpara>(num_items, [=](id<2> i) {
    size_t row = i[0], col = i[1];
```

set up  
problem size

```
    C[row][col] = 0;
    for (size_t k = 0; k < Wa; k++)
        C[row][col] += A[row * Wa + k] * B[k * Wb + col];
});
```

lambda function  
for device

```
});
```



Learning with Purpose

# Demonstration

Compilation and execution of Matrix Multiplication example on Intel FPGA DevCloud



*Learning with Purpose*

# A Very Brief Introduction to FPGA Design Concepts

A part of the DPC++ Tutorial Series

Prof. Yan Luo

Acknowledgement

This work is supported by Intel Corporation 2020-2021



*Learning with Purpose*

# Agenda

- Introduction to FPGA Architecture
- Concepts of FPGA Hardware Design
- Mapping Source Code to Hardware Datapath
- Scheduling
- Parallelism Models to FPGA Hardware
- Memory Types

Some materials used in this presentation are based on Intel®  
OneAPI DPC++ FPGA Optimization Guide



Learning with Purpose

# FPGA vs CPU

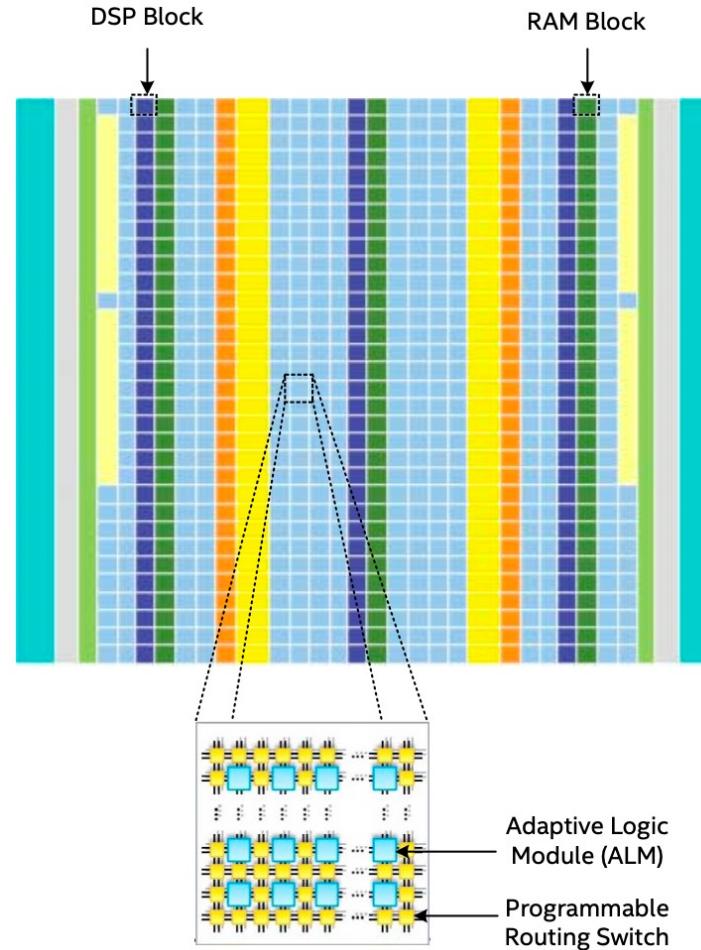
- FPGA does not have a fixed datapath
  - that is why it is “field programmable” !
- “Program” the hardware resources
  - You have a lot more control on how your design is mapped
- Function as accelerators to offload intensive computing
- Design methodologies
  - Hardware Description Language
  - High level language (like DPC++)



Learning with Purpose

# FPGA Architecture

- Adaptive Logic Module
- RAM block
- DSP block
- Programmable routing switch



# FPGA Hardware Design Concepts

- Maximum frequency  $f_{MAX}$
- Latency
- Pipelining
- Throughput
- Datapath
- Control path
- Occupancy



Learning with Purpose

$f_{MAX}$

- $f_{MAX}$  is the max clock frequency a digital circuit operates at
- The max rate at which the outputs or registers are updated
- Clock speed is limited by physical propagation delay
  - Signal propagation across combination logic between two consecutive register stages
- The propagation delay is
  - a function of the complexity of the combinational logic
- Critical path
  - The path with the most combinational logic elements
  - Limits the speed of the entire circuit



Learning with Purpose

# Latency

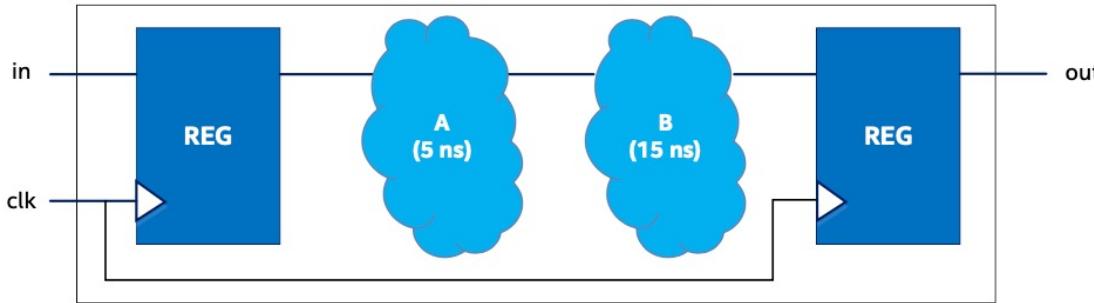
- Measures how long it takes to complete one or more operations
  - Latency of a single operation within
  - Or latency of the entire circuit
- Measure in time (ms) or cycles
  - Cycles are preferred
- Lowering latency ?
  - Yes, shorter latency is desirable in general
  - But shorter latency may decrease  $f_{MAX}$
  - So, caution on the design alternatives



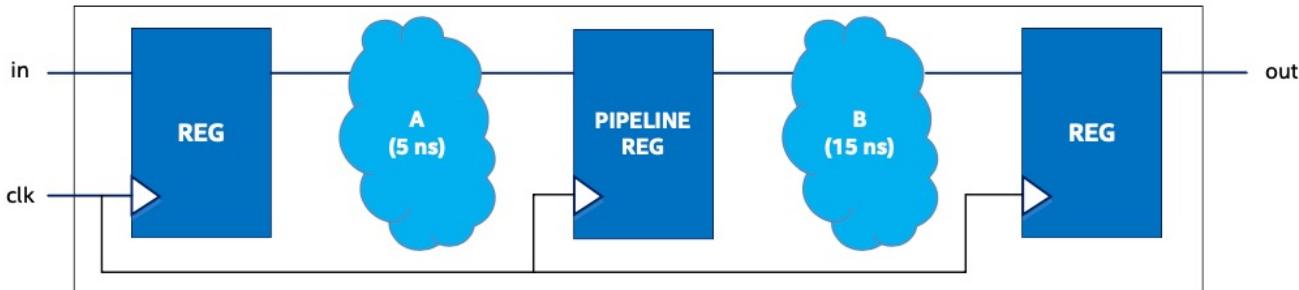
Learning with Purpose

# Pipelining to improve $f_{MAX}$

unpipelined  
 $f_{MAX} = 50\text{MHz}$   
Delay of 2 cycles



pipelined  
 $f_{MAX} = 66.7\text{MHz}$   
delay of 3 cycles



# Throughput

- The rate at which data is processed
- Higher fMAX the higher throughput
- Throughput is NOT the inverse of delay
  - Parallel processing or overlapped data processing



Learning with Purpose

# Datapath and Control Path

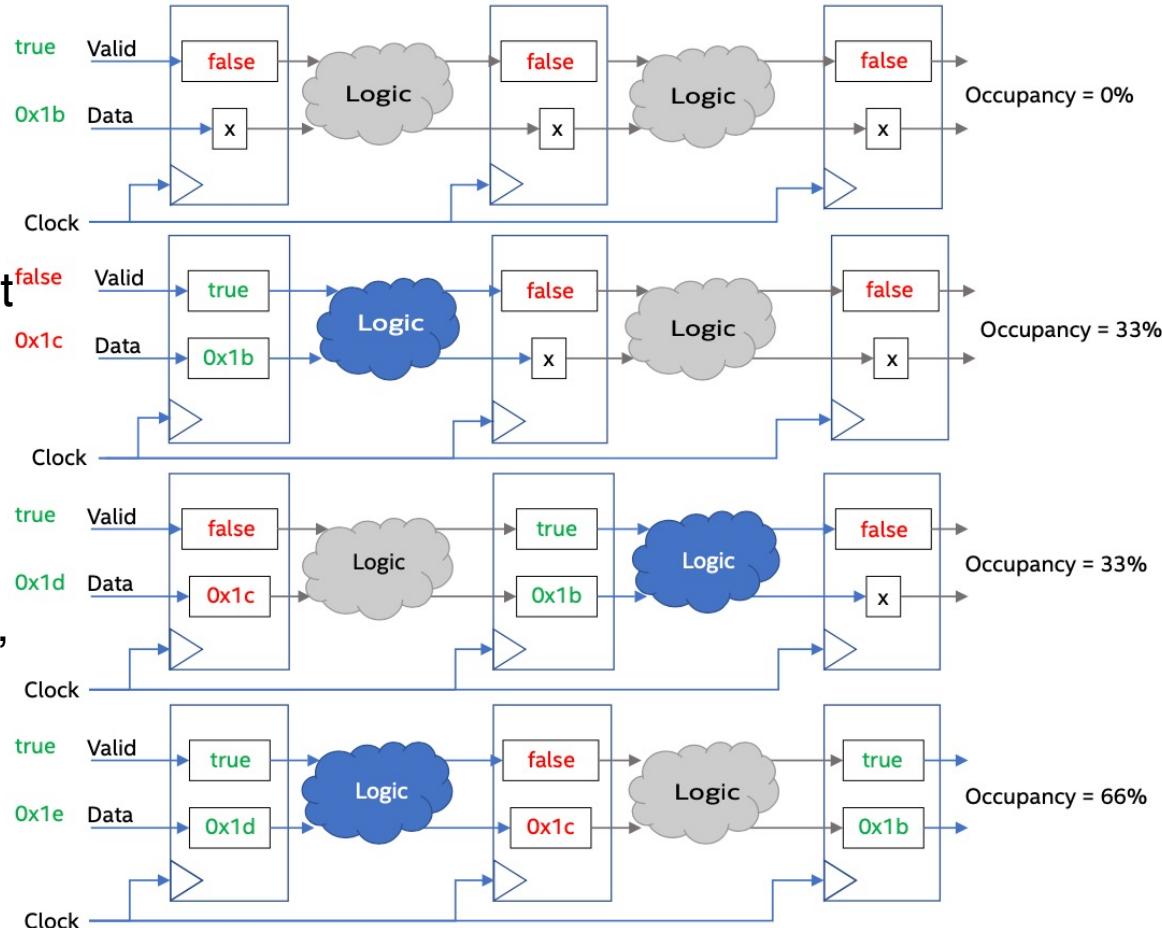
- Datapath
  - A chain of registers and combinational logic
  - Memory blocks or offchip memory units are considered as outside of the datapath
- Control Path
  - Path of the signals that control the operation of the circuit
  - Logic added for
    - Handshaking flow control
    - Loop control
    - Branch control



Learning with Purpose

# Occupancy

- Datapath Occupancy: At a point in time the proportion of the datapath that contains **valid** data
- Circuit Occupancy: average occupancy over time from program start to end
- Unoccupied portion = "bubbles"
- Decreasing bubbles increases occupancy



# DPC++ Design Analysis (I): Analyze FPGA Early Image

A part of the DPC++ Tutorial Series

Prof. Yan Luo

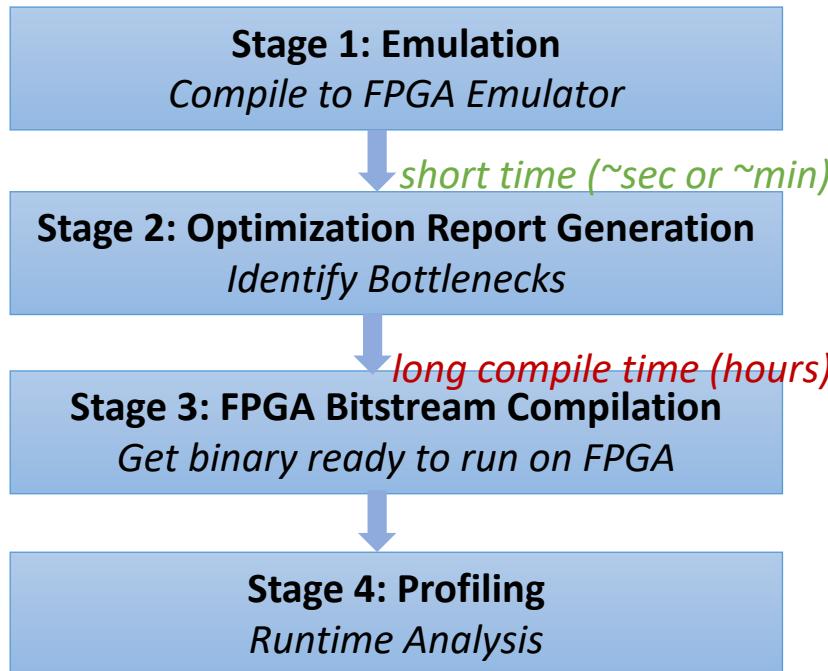
Acknowledgement

This work is supported by Intel Corporation 2020-2021



*Learning with Purpose*

# Analyze your Design before Optimization



1. Make sure the design is functionally correct
2. Look for bottlenecks through compilation reports and revise design if necessary
3. Generate FPGA binary for execution on hardware
4. Run profiling to analyze runtime performance



Learning with Purpose

# Demo: Compilation Report Generation and Analysis

Look through compilation report of Matrix Multiplication example



*Learning with Purpose*

# Static Analysis Report

The screenshot shows a static analysis report for a project named "matrix-multi-para-v1\_report". The interface includes tabs for Reports, Summary, Throughput Analysis, Area Analysis, and System Viewers. The Summary tab is active, displaying sections for Compile Info, Kernels Summary, Clock Frequency Summary, System Resource Utilization Summary, Quartus Fitter Resource Utilization Summary, Compile Estimated Kernel Resources, and Warnings Summary. The Kernels Summary table lists one kernel: MMpara\_v1, located at :0, with NDRange type, Autorun set to No, n/a workgroup size, 1 compute unit, and Not specified target frequency. The code editor on the right displays the source code for "dpc\_common.hpp".

Name	Source Location	Kernel Type	Autorun	Workgroup Size	# Compute Units	Target Frequency (MHz)
MMpara_v1	:0	NDRange	No	n/a	1	Not specified

```
1 ///////////////////////////////////////////////////////////////////////////////
2 // Copyright © 2020 Intel Corporation
3 //
4 // SPDX-License-Identifier: MIT
5 // =====
6
7 #ifndef _DP_HPP
8 #define _DP_HPP
9
10 #pragma once
11
12 #include <stdlib.h>
13 #include <exception>
14
15 #include <CL/sycl.hpp>
16
17 - namespace dpc_common {
18 // this exception handler with catch async exceptions
19 - static auto exception_handler = []()>>cl::sycl::exception_list elist) {
20 -     for (std::exception_ptr const &e : elist) {
21 -         try {
22 -             std::rethrow_exception(e);
23 -         } catch (std::exception const &e) {
24 - #if _DEBUG
25 -             std::cout << "Failure" << std::endl;
26 - #endif
27 -             std::terminate();
28 -         }
29 -     }
30 - };
31
32 // The TimeInterval is a simple RAII class.
33 // Construct the timer at the point you want to start timing.
34 // Use the Elapsed() method to return time since construction.
35
36 - class TimeInterval {
37     public:
```

# Clock Frequency and System Resource Utilization Summary

Reports Summary Throughput Analysis ▾ Area Analysis ▾ System Viewers ▾

Summary				
Clock Frequency Summary				
	Quartus Fitter: Clock Frequency (MHz)	Compile Target Frequency (MHz)	Compile Estimated Frequency (MHz)	
Kernel clock	TBD(?)	240.00	240.00	

System Resource Utilization Summary				
	Device	Static partition	Quartus Fitter: Total Used (Entire System)	Estimated: Kernel system
ALM	427200	89975	TBD	
- ALUT				22471
- REG	1708800	358572	TBD	32268
- MLAB				71
RAM	2713	492	TBD	214
DSP	1518	123	TBD	22



Learning with Purpose

# System Graph Viewer

Report: matrix-multi-par x +

file:///home/yluo/projects/dpcpp-tutorial/matrix-multi/build/matrix-multi-para-v1\_report.prj/reports/report.html# 120% ...

Reports Summary Throughput Analysis Area Analysis System Viewers

Graph List (Beta)

- System
  - MMpara\_v1
    - MMpara\_v1.B0
    - MMpara\_v1.B1
    - Cluster 0
    - Cluster 1
  - MMpara\_v1.B2
    - Cluster 2
    - Cluster 3

Graph Viewer (Beta)

Graph Viewer (Beta)

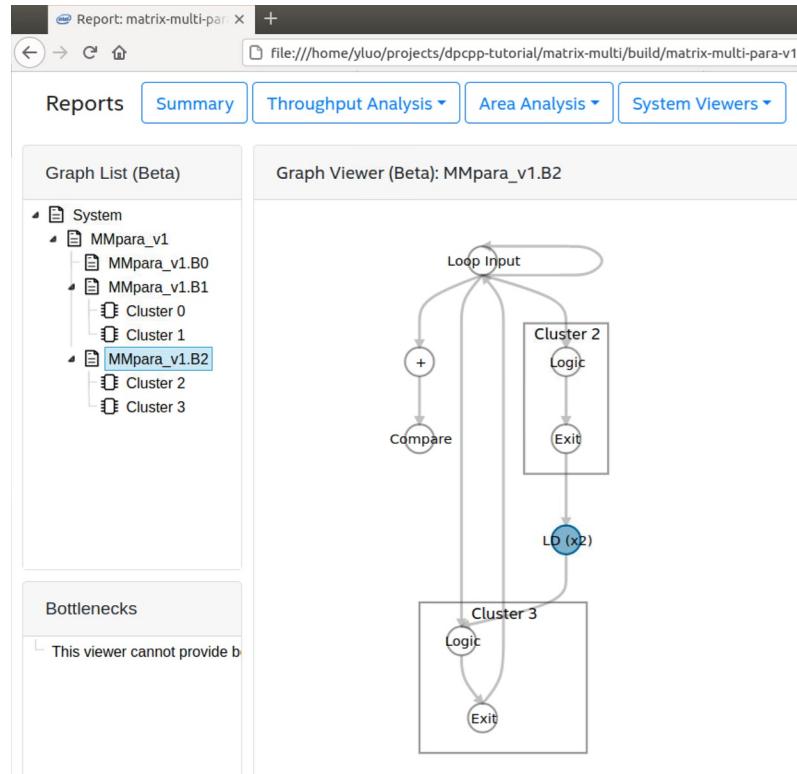
Kernel Memory Viewer

Schedule Viewer (Beta)

Graph Viewer (beta)  
Provides views of different levels of your design: system, Kernel, block, and cluster.

```
// Copyright © 2020 Intel Corporation  
// SPDX-License-Identifier: MIT  
// ======  
6  
7 ifndef _DP_HPP  
8 define _DP_HPP  
9  
10 pragma once  
11  
12 include <stdlib.h>  
13 include <exception>  
14  
15 include <CL/sycl.hpp>  
16  
17 namespace dpc_common {  
18 // this exception handler with cat  
19 static auto exception_handler = []  
20 for (std::exception_ptr const &e:  
21 try {  
22 std::rethrow_exception(e);  
23 } catch (std::exception const &  
24 #ifdef _DEBUG  
25 std::cout << "Failure" << std::endl;  
26 #endif  
27 std::terminate();  
28 }  
29 };  
30  
31 // The TimeInterval is a simple RAII  
32 // Construct the timer at the point of creation  
33 // Use the Elapsed() method to get the time  
34
```

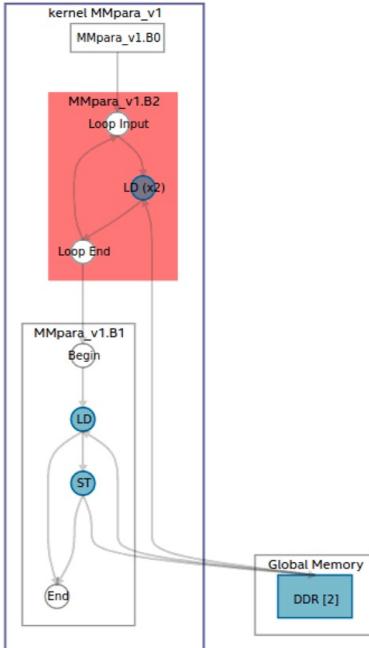
# Zoom into a Module



Learning with Purpose

# parallel for() v1

Graph Viewer (Beta)



matrix-multi-para-v1.cpp

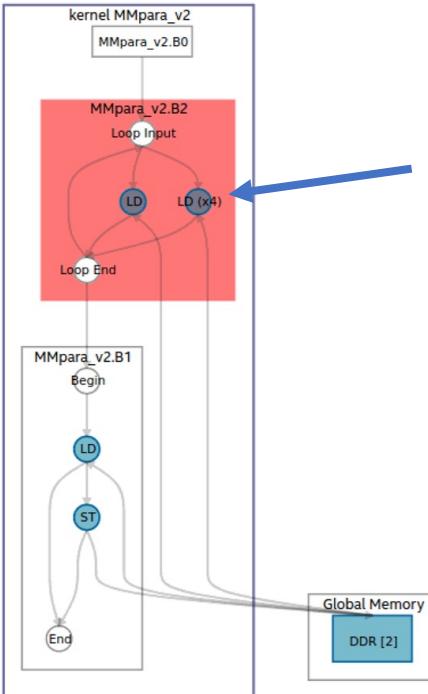
```
63 // In our case, we have a two-dimensional nd_range:  
64 // num_items: the global size, or 'all' the work in 2D, i.e. the  
size  
65 // dimensions  
66 // range<2>(1,1) : the workgroup size. (1,1) means a workgroup has  
1 work item  
67 // in each dimension  
68 // 2nd parameter is the kernel, a lambda that specifies what to do  
per  
69 // work item. The parameter of the lambda is the work item id.  
70 // DPC++ supports unnamed lambda kernel by default.  
71 auto kernel_range = nd_range<2>(num_items, range<2>(1,1));  
72 h.parallel_for<MMpara_v1>(num_items, [=](id<2> i)  
{ size_t row = i[0], col = i[1];  
74  
75 float s = 0;  
76 //##pragma unroll 2  
77 for (size_t k = 0; k < a_columns; k++)  
| s += a[row][k] * b[k][col];  
78  
79 sum[row][col] = c[row][col] + s;  
80 })  
81 );  
82 );  
83  
84 #if FPGA || FPGA_PROFILE  
85 // Query event e for kernel profiling information  
86 // (blocks until command groups associated with e complete)  
87 double kernel_time_ns =  
88 e.get_profiling_info<info::event_profiling::command_end()-  
e.get_profiling_info<info::event_profiling::command_start>();  
89  
90 // Report profiling info  
91 std::cout << "Kernel compute time: " << kernel_time_ns * 1e-6 <<  
92 ns  
93  
94 //*****  
95 // Demonstrate matrix multiplication both in sequential on CPU and in  
parallel on device.  
96 //*****  
97 //*****  
98 int main()  
99 {  
// Create device selector for the device of your interest
```



Learning with Purpose

# parallel for() v2

Graph Viewer (Beta)



matrix-multi-para-v2.cpp

```

51 // Create an accessor for each buffer with access permission: read, write or
52 // read/write. The accessor is a mean to access the memory in the buffer.
53 auto a = a_buf.get_access<access::mode::read, access::target::global_buffer>(h);
54 auto b = b_buf.get_access<access::mode::read, access::target::global_buffer>(h);
55 auto c = c_buf.get_access<access::mode::read, access::target::global_buffer>(h);
56
57 // The sum_accessor is used to store (with write permission) the sum data.
58 auto sum = sum_buf.get_access<access::mode::write>(h);
59
60 // Use parallel_for to run vector addition in parallel on device. This
61 // executes the kernel.
62 // 1st parameter is the number of work items in total and in a workgroup
63 // In our case, we have a two-dimensional nd_range:
64 // num_items: the global size, or 'all' the work in 2D, i.e. the size
65 // of the matrix Sum in 'row' and 'column' dimensions
66 // range<2>(1,1) : the workgroup size. (1,1) means a workgroup has 1 work item
67 // in each dimension
68 // 2nd parameter is the kernel, a lambda that specifies what to do per
69 // work item. The parameter of the lambda is the work item id.
70 // DPC++ supports unnamed lambda kernel by default.
71 auto kernel_range = nd_range<2>(num_items, range<2>(1,1));
72 h.parallel_for<MMpara_v2>(num_items, [=](id<2> i)
73 {
    size_t row = i[0], col = i[1];
74
75     float s = 0;
76     #pragma unroll 4
77     for (size_t k = 0; k < a_columns; k++)
78         s += a[row][k] * b[k][col];
79
80     sum[row][col] = c[row][col] + s;
81 });
82
83 #if FPGA || FPGA_PROFILE
84 // Query event e for kernel profiling information
85 // (blocks until command groups associated with e complete)
86 double kernel_time_ns =
87     e.get_profiling_info<info::event_profiling::command_end() - info::event_profiling::command_start>();
88
89 // Report profiling info
90 std::cout << "Kernel compute time: " << kernel_time_ns * 1e-6 << " ms\n";
91 #endif
92 }
93
94 //*****
95 // Demonstrate matrix multiplication both in sequential on CPU and in parallel on device.
96 //*****
97
98 int main()
99 {
100     // Create device selector for the device of your interest.
101 #if FPGA_EMULATOR
102     // DPC++ extension: FPGA emulator selector on systems without FPGA card.

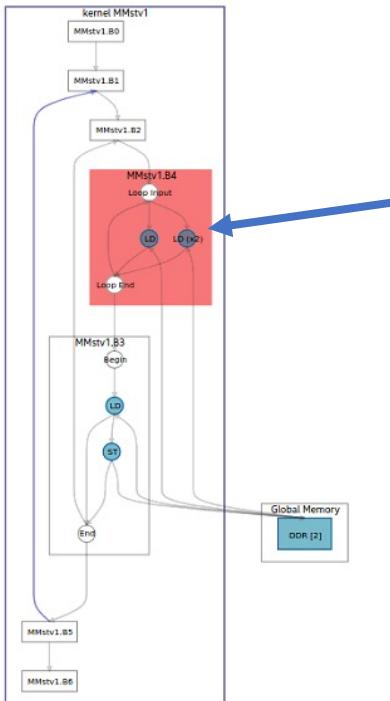
```



Learning with Purpose

# single\_task() v1

Graph Viewer (Beta)



matrix-multi-st-v1.cpp

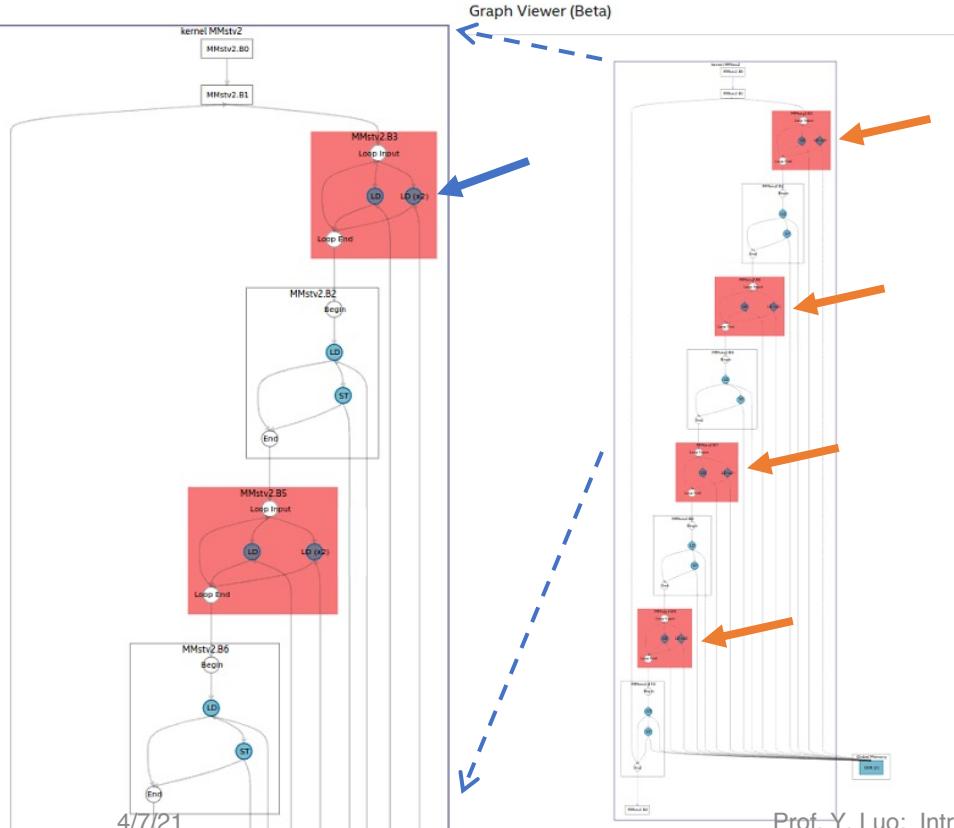
```

50 v event e = q.submit([&](handler &h) {
51   // Create an accessor for each buffer with access permission: read, write or
52   // read/write. The accessor is a mean to access the memory in the buffer.
53   auto a = a_buf.get_access<access::mode::read, access::target::global_buffer>(h);
54   auto b = b_buf.get_access<access::mode::read, access::target::global_buffer>(h);
55   auto c = c_buf.get_access<access::mode::read, access::target::global_buffer>(h);
56
57   // The sum_accessor is used to store (with write permission) the sum data.
58   auto sum = sum_buf.get_access<access::mode::write>(h);
59
60   // Use single_task to run vector addition in parallel on device. This
61   // executes the kernel.
62   // A kernel that is executed on one thread using NDRange(1,1,1) is enqueued
63   // using the cl::sycl::single_task API:
64   // single_task<typename kernel_lambda_name>(<=>() [[intel::kernel_args_restrict]]
65   h.single_task<MMstv1>(<=>() [[intel::kernel_args_restrict]])
66   {
67     size_t row, col;
68     float s = 0;
69     for (row = 0; row < a_rows; row++)
70     {
71       for (col = 0; col < b_columns; col++)
72       {
73         #pragma unroll 2
74         for (size_t k = 0; k < a_columns; k++)
75           s += a[row][k] * b[k][col];
76         sum[row][col] = c[row][col] + s;
77       }
78     });
79 #if FPGA || FPGA_PROFILE
80   // Query event e for kernel profiling information
81   // (blocks until command groups associated with e complete)
82   double kernel_time_ns =
83     e.get_profiling_info<info::event_profiling::command_end>() -
84     e.get_profiling_info<info::event_profiling::command_start>();
85
86   // Report profiling info
87   std::cout << "Kernel compute time: " << kernel_time_ns * 1e-6 << " ms\n";
88 #endif
89
90 }
91
92
93 //*****
94 // Demonstrate matrix multiplication both in sequential on CPU and in parallel on device.
95 //*****
96 v int main() {
97   // Create device selector for the device of your interest.
98 #if FPGA_EMULATOR
99   // DPC++ extension: FPGA emulator selector on systems without FPGA card.
100   INTEL::fpga_emulator_selector d_selector;
101 #elif defined(FPGA) || defined(FPGA_PROFILE)
102   // DPC++ extension: FPGA selector on systems with FPGA card.
  
```



Learning with Purpose

# single\_task() v2



```
matrix-multi-st-v2.cpp
```

```

50 // Submit a command group to the queue by a lambda function that contains the
51 // data access permission and device computation (kernel).
52 event = q.submit([&]{
53     // Create an array to read buffer with access permission: read, write or
54     // read/write. The accessor is a mean to access the memory in the buffer
55     auto a = a_buf.get_access<access::mode::read, access::target::global_buffer>(h);
56     auto b = b_buf.get_access<access::mode::read, access::target::global_buffer>(h);
57     auto c = c_buf.get_access<access::mode::read, access::target::global_buffer>(h);
58
59     // The sum_accessor is used to store (with write permission) the sum data.
60     auto sum = sum_buf.get_access<access::mode::write>(h);
61
62     // Use single_task to run vector addition in parallel on device. This
63     // executes the kernel.
64     // A kernel that is executed on one thread using NRange(1,1,1) is enqueued
65     // using the cl::sycl::single_task API:
66     // single_task<typename kernel_lambda_name>(<=)(());
67     h.single_task<MMstv2>(<=) [[intel::kernel_args_restrict]];
68
69 {
70     float s = 0;
71     #pragma unroll 4
72     for (size_t i = 0; i < a_rows * b_columns; i++) {
73         size_t row, col;
74         row = i / widthC;
75         col = i % widthC;
76         #pragma unroll 2
77         for (size_t k = 0; k < a_columns; k++) {
78             s += a[row][k] * b[k][col];
79             sum[row][col] = c[row][col] + s;
80         }
81     };
82 }
83
84 #if FPGA || FPGA_PROFILE
85 // Query event e for kernel profiling information
86 // (blocks until command groups associated with e complete)
87 double kernel_time_ns =
88     e.get_profiling_info<event_profiling::command_end>() -
89     e.get_profiling_info<event_profiling::command_start>();
90
91 // Report profiling info
92 std::cout << "Kernel compute time: " << kernel_time_ns * 1e-6 << " ms\n";
93 #endif
94 }
95
96
97 //*****
98 // Demonstrate matrix multiplication both in sequential on CPU and in parallel on device.
99 //*****
100
101 int main() {
102     // Create device selector for the device of your interest.
103     #if DPC++ extension: FPGA emulator selector on systems without FPGA card.
104     INTEL::fpga_emulator_selector d_selector;
105     #elif defined(FPGA) || defined(FPGA_PROFILE)
106     // DPC++ extension: FPGA selector on systems with FPGA card.
107     INTEL::fpga_selector d_selector;
108     #else
109     // The default device selector will select the most performant device.
110     default_selector d_selector;
111     #endif

```

# DPC++ Design Analysis (II): Analyze Runtime Profiling Data

A part of the DPC++ Tutorial Series

Prof. Yan Luo

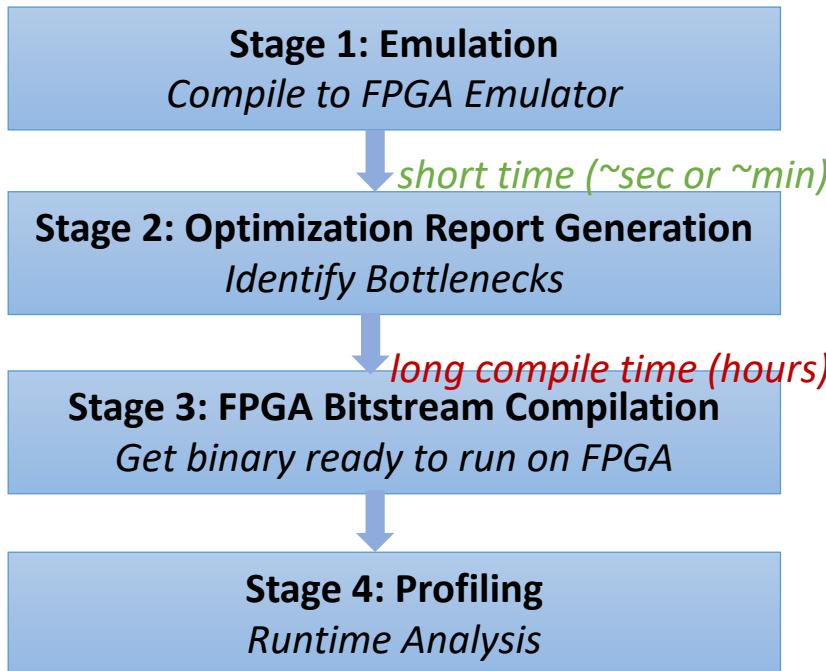
Acknowledgement

This work is supported by Intel Corporation 2020-2021



*Learning with Purpose*

# Analyze your Design before Optimization



1. Make sure the design is functionally correct
2. Look for bottlenecks through compilation reports and revise design if necessary
3. Generate FPGA binary for execution on hardware
4. Run profiling to analyze runtime performance

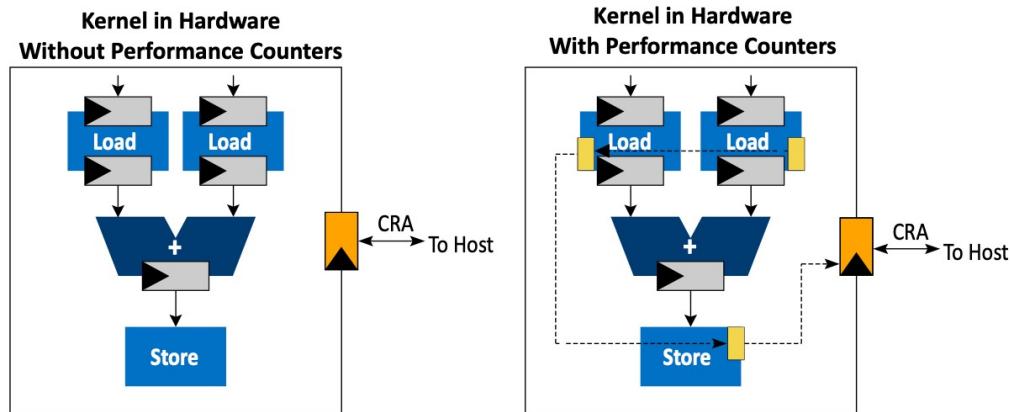


Learning with Purpose

# Compile for profiling

- Intel FPGA Dynamic Profiler
  - User performance counters to collect kernel performance data
    - Stall, occupancy, idle, bandwidth usage, etc.
  - The data can be viewed in Intel VTune Profiler

To enable profiling,  
add `-xsprofile`  
flag to compiler  
command



Learning with Purpose

# Profiling data collection

- Execute FPGA binary with profiling tool. Do not generate json file

```
aocl profile -no-json ./matrix-multi-para-v1.fpga_profile
```

- Generate json results file using profile.mon

```
aocl profile -no-run profile.mon ./matrix-multi-para-v1.fpga_profile
```

- Load profile results folder to Intel VTune toolkit
  - Need to put the files (profile.mon and profile.json) in a folder
  - Check “import multiple trace files from a directory” option when importing profiling data in VTune.



Learning with Purpose

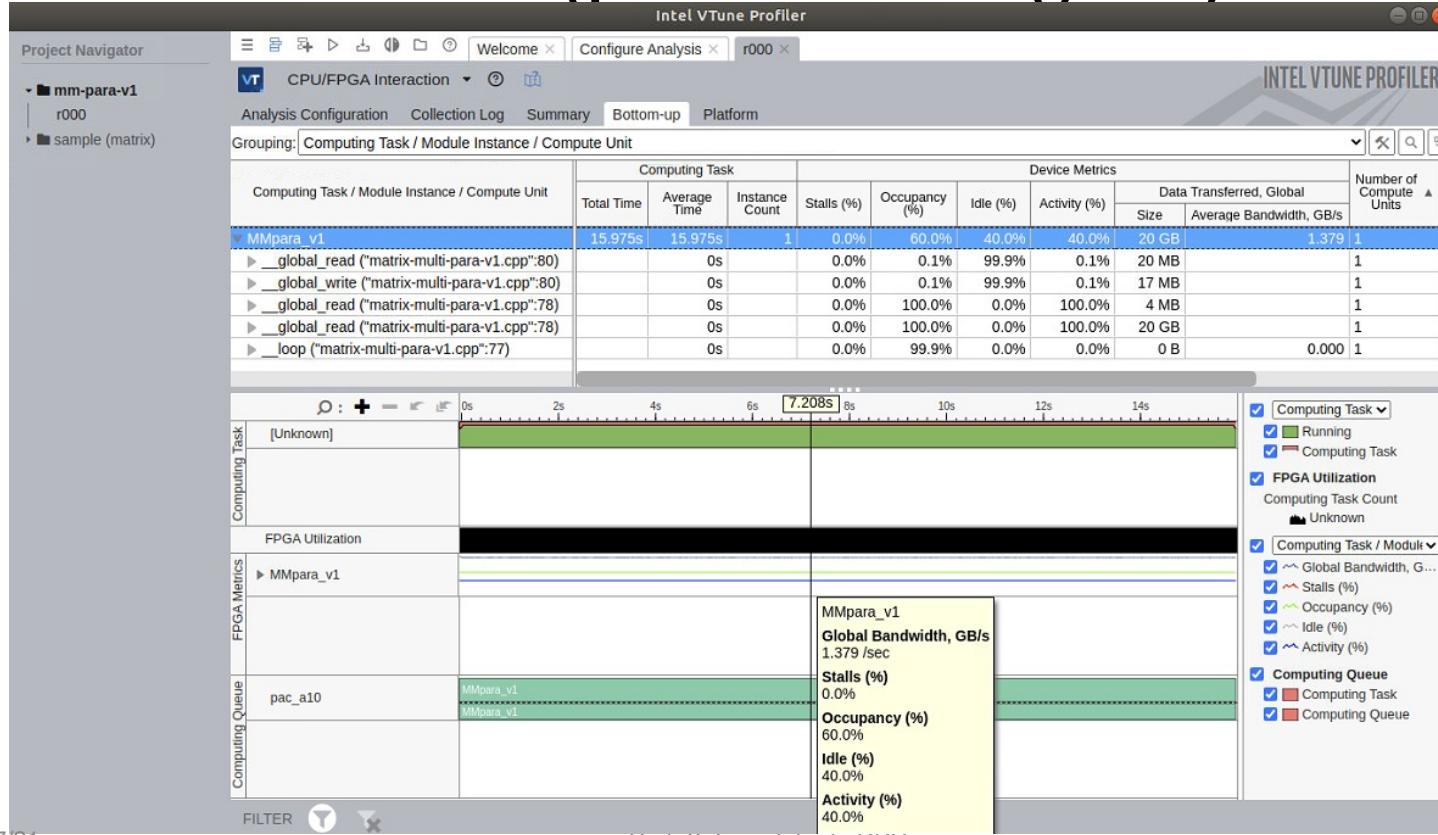
# Demo: Profiling Analysis

Run-time profiling analysis on Matrix Multiplication example



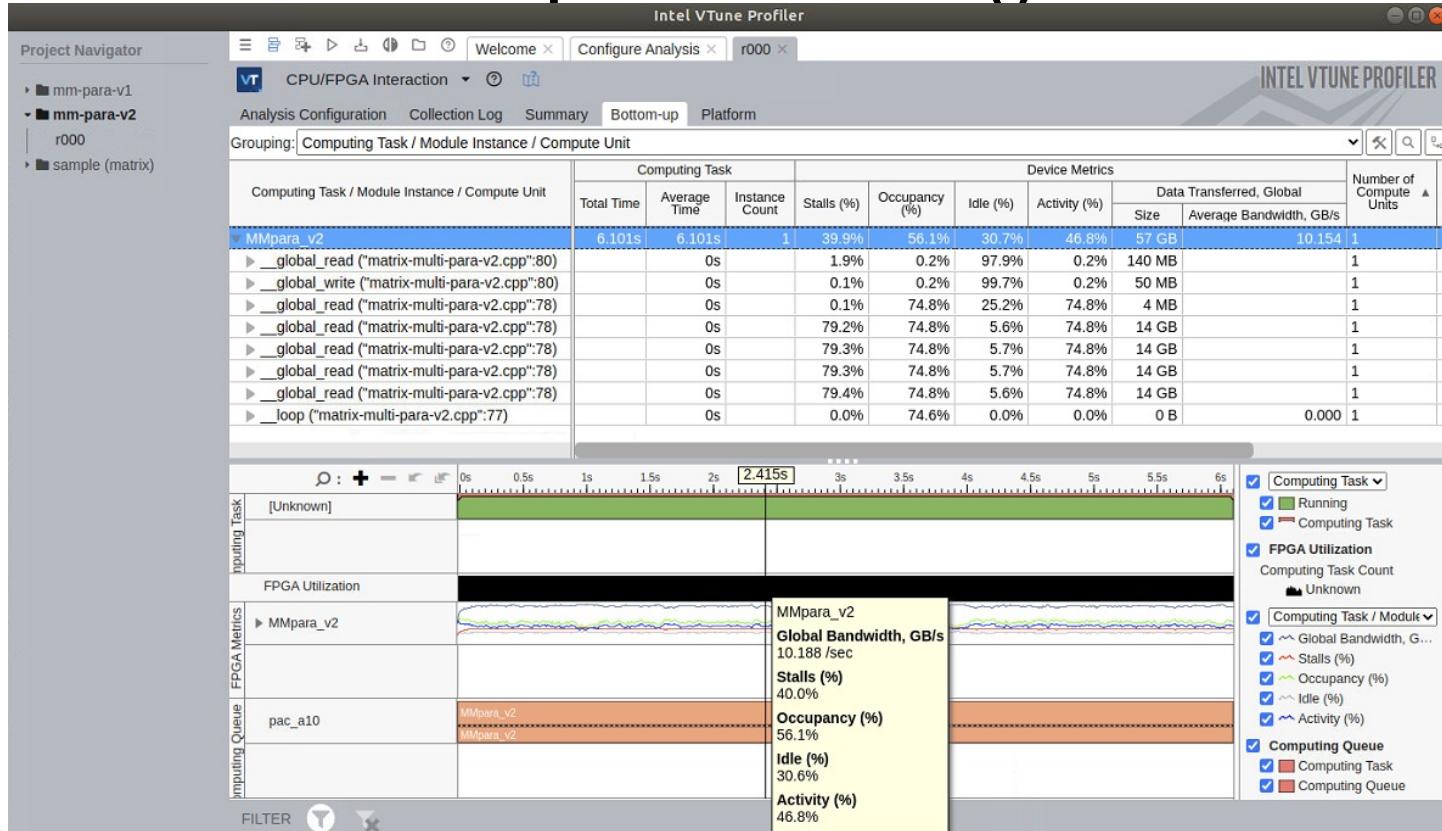
*Learning with Purpose*

# VTune Profiler (parallel for() v1)



Learning with Purpose

# VTune Profiler parallel for() v2



# DPC++ Design Optimization: Reducing Global Memory Access with Local Memory

A part of the DPC++ Tutorial Series

Prof. Yan Luo

Acknowledgement

This work is supported by Intel Corporation 2020-2021



*Learning with Purpose*

# Motivation

- We observe a large amount of global memory access
  - Matrix-multi in those implementations are a memory-bound
  - The total sizes of matrices A, B = 25.6MB  
800x1600 and 1600x3200 floating point numbers in A, B respectively
  - We read 20GB from global memory!
- The matrix-multi calculation uses the same element from source matrices multiple times
  - Reads an entire row of A and an entire column of B to calculate **each** value of C
  - Adjacent elements in C require reading much of the same data (row from matrix A or a column from matrix B)



Learning with Purpose

# Tiling

- Tiling is buffering data onto fast on-chip storage where it will be repeatedly accessed (caching)
  - Very common technique
  - Used when multiple instances need to access overlapping parts of data set
- Data must be partitioned into blocks to fit into local memory
  - Only work-items within a workgroup can share data
  - Local memory size and geometry set at compile time
  - Workgroup sizes (block sizes) must be known at compile time

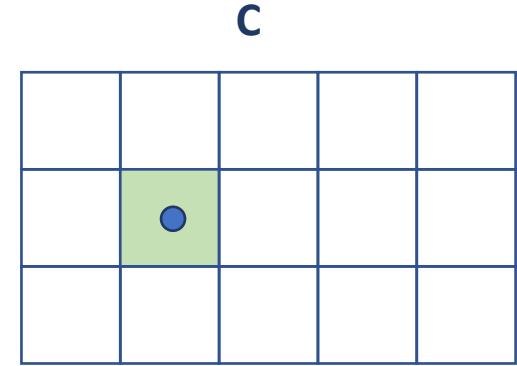


Learning with Purpose

# Tiling Illustration

		A			
		0	1	2	3
0	0	Light Blue	Dark Purple	Orange	Yellow
	1	Light Blue	Purple	Brown	Yellow
2	Light Blue	Dark Purple	Orange	Yellow	

		B				
		0	1	2	3	4
0	0	Light Blue	Light Blue	Light Blue		
	1	Dark Purple	Purple	Brown	Orange	Yellow
2	Light Blue	Dark Purple	Orange	Yellow		
3	Light Blue	Dark Purple	Orange	Yellow		



- Perform element-wise multiplication on tiles (at different times or in parallel)

$$A_{10} \otimes B_{01} \quad A_{11} \otimes B_{11} \quad A_{12} \otimes B_{21} \quad A_{13} \otimes B_{31}$$

- The accumulated sum is from multiple tile-wise multiplications

$$C_{11} = A_{10} \otimes B_{01} + A_{11} \otimes B_{11} + A_{12} \otimes B_{21} + A_{13} \otimes B_{31}$$

# Overview of the Design

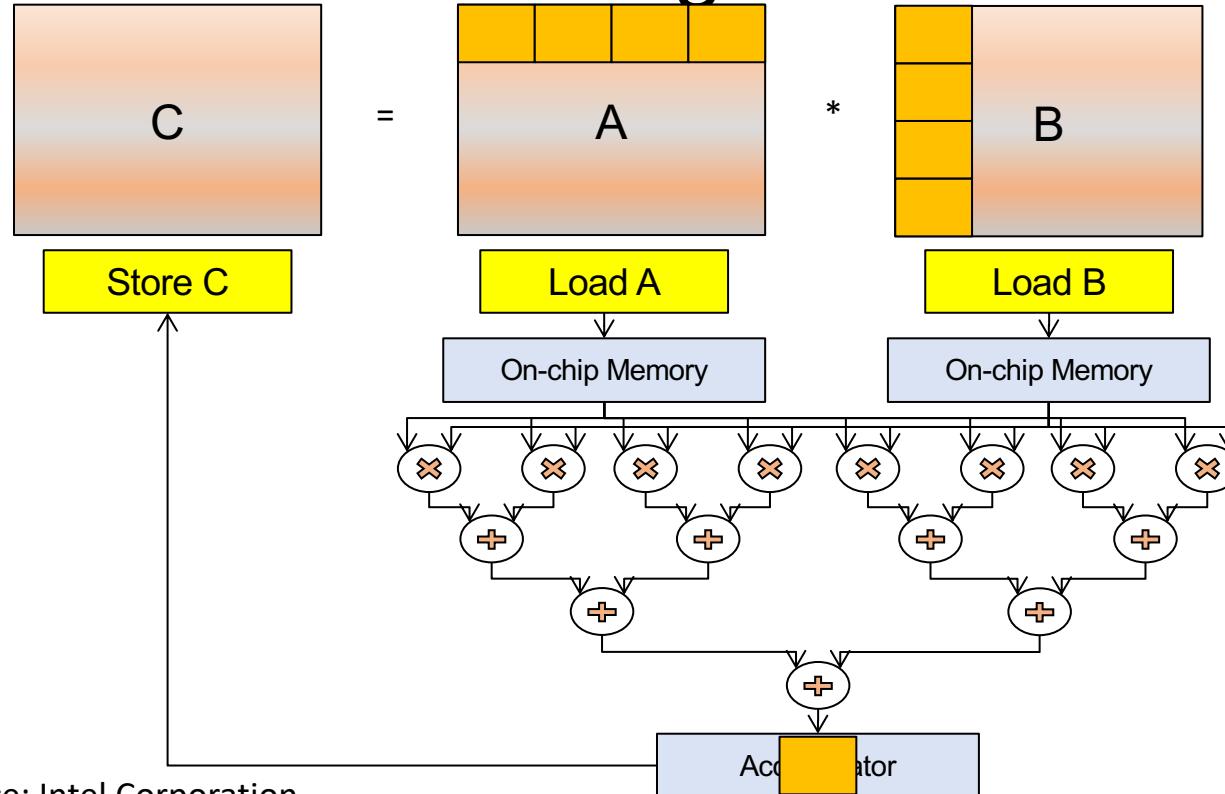


Diagram source: Intel Corporation

4/7/21

Prof. Y. Luo: Intro to DPC++



Learning with Purpose

43

# Implementation Details

- Calculate the number of rounds for tile-wise multiplication
- Create a single-task kernel
- Allocate local memory for tiles
- Loop unrolling optimization



Learning with Purpose

# Tile-wise multiplication kernel (step 1)

- Calculate the number of tiles, tile row and column ID

```
for(int i=0; i < a_rows*a_columns/(BLOCK_SIZE*BLOCK_SIZE); i++) {  
    // the block indices of the block in A  
    auto block_row_a = i / (a_columns/BLOCK_SIZE);  
    auto block_col_a = i % (a_columns/BLOCK_SIZE);  
    // we need to calculate dot-product with all the blocks in B where  
    // the row number is equal to block_col_a  
    auto block_row_b = block_col_a;  
    for(int j=0; j<b_columns/BLOCK_SIZE; j++)  
    {  
        .....  
    }  
}
```



Learning with Purpose

# Tile-wise multiplication kernel (step 2)

- Allocation of local memory for tile
  - $\text{BLOCK\_SIZE} \times \text{BLOCK\_SIZE}$  floats
  - Need one local memory block to hold a tile from A, B (and D)

```
// allocate local memory to hold a block of data from A, B
[[intel::numbanks(NUM_BANKS), intel::bankwidth(BANK_WIDTH)]] float local_mem_a[BLOCK_SIZE][BLOCK_SIZE];
[[intel::numbanks(NUM_BANKS), intel::bankwidth(BANK_WIDTH)]] float local_mem_b[BLOCK_SIZE][BLOCK_SIZE];
[[intel::numbanks(NUM_BANKS), intel::bankwidth(BANK_WIDTH)]] float local_mem_d[BLOCK_SIZE][BLOCK_SIZE];
```

- Note **NUMBANKS** and **BANK\_WIDTH**
  - Large numbers help but are limited by the FPGA's MLAB resources.



Learning with Purpose

# Tile-wise multiplication kernel (step 3)

- Then load the tile data from global memory to local memory for computation
  - DPC++ compiler optimize with pipelined Load-Store units

```
// load blocks of data to local memory from global memory
for (m=0; m < BLOCK_SIZE; m++)
    for ( n=0; n < BLOCK_SIZE; n++)
    {
        local_mem_a[m][n] = a[block_row_a*BLOCK_SIZE + m][block_col_a*BLOCK_SIZE + n];
        local_mem_b[m][n] = b[block_row_b*BLOCK_SIZE + m][j*BLOCK_SIZE + n];
        local_mem_d[m][n] = d[block_row_a*BLOCK_SIZE + m][j*BLOCK_SIZE + n];
    }
```



Learning with Purpose

# Tile-wise multiplication kernel (step 4)

- Element wise multiplication and accumulation within a Tile

```
// element-wise multiplication and accumulation
for (m=0; m < BLOCK_SIZE; m++) {
    for ( n=0; n < BLOCK_SIZE; n++) {
        s = 0;
        #pragma unroll
        for (k=0; k < BLOCK_SIZE; k++)
            s += local_mem_a[m][k] * local_mem_b[k][n];
        // add to Matrix D
        // the corresponding row and col in D
        row = block_row_a * BLOCK_SIZE + m;
        col = j * BLOCK_SIZE + n;
        //d[row][col] += s;
        local_mem_d[m][n] += s;
    }
}
```

- Note that we are doing  $D = A \times B + C$ . so the final accumulation is D.
- We will add C to D at a later step.



Learning with Purpose