

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



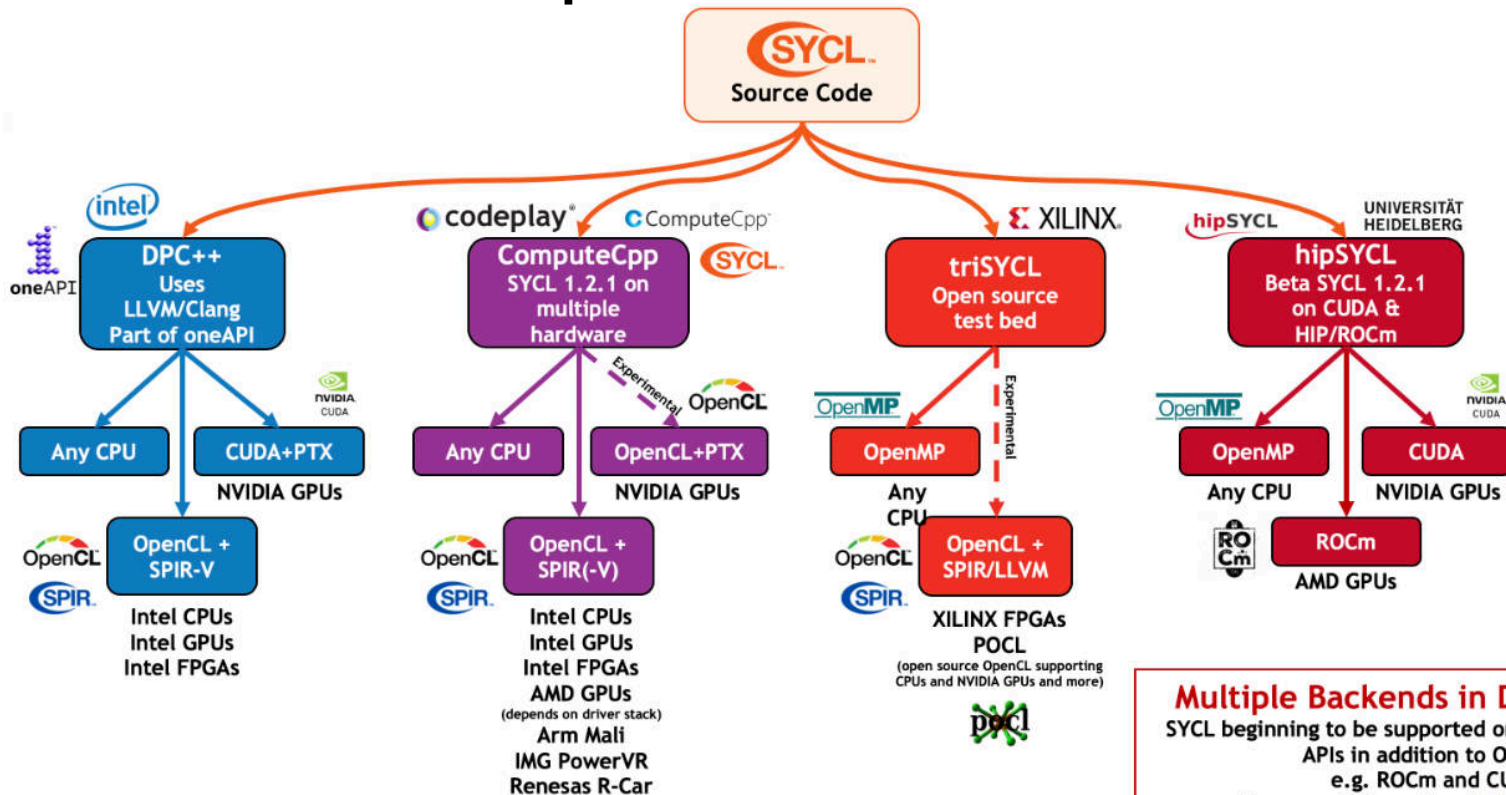
Learning with Purpose

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)



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         {
14             h.parallel_for(r, [=](item<1> idx) {
15                 out[idx] = idx;
16             });
17         }
18     });
19
20     host_accessor result{a};
21     for (int i=0; i<num; ++i)
22         std::cout << result[i] << "\n";
23 }
```

Header Files

namespace
scope

device queue
& command

SYCL buffer
declaration

kernel function
on device

Lambda
function

host access
results in buffer

Prof. Y. Luo: Intro to DPC++

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



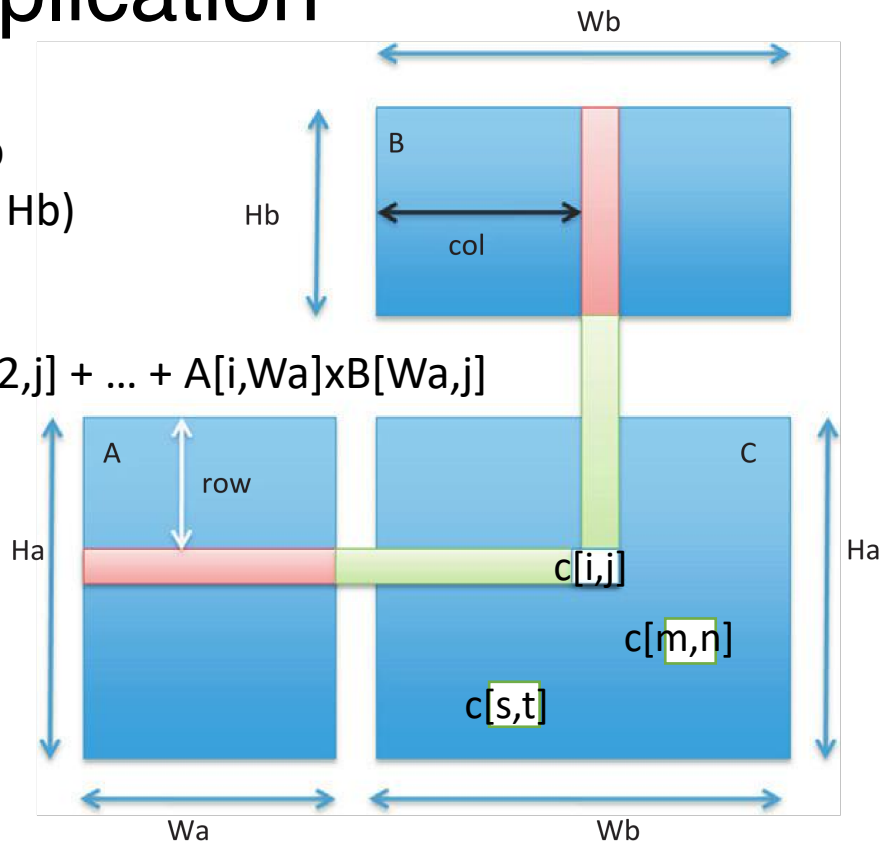
Learning with Purpose

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; j < 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];
    }
```

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};  
  
    h.parallel_for<class MMpara>(num_items, [=](id<2> i) {  
        size_t row = i[0], col = i[1];  
  
        C[row][col] = 0;  
        for (size_t k = 0; k < Wa; k++)  
            C[row][col] += A[row * Wa + k] * B[k * Wb + col];  
    });  
});
```

create
accessors

set up
problem size

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
- Trivia

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



Learning with Purpose

FPGA vs CPU

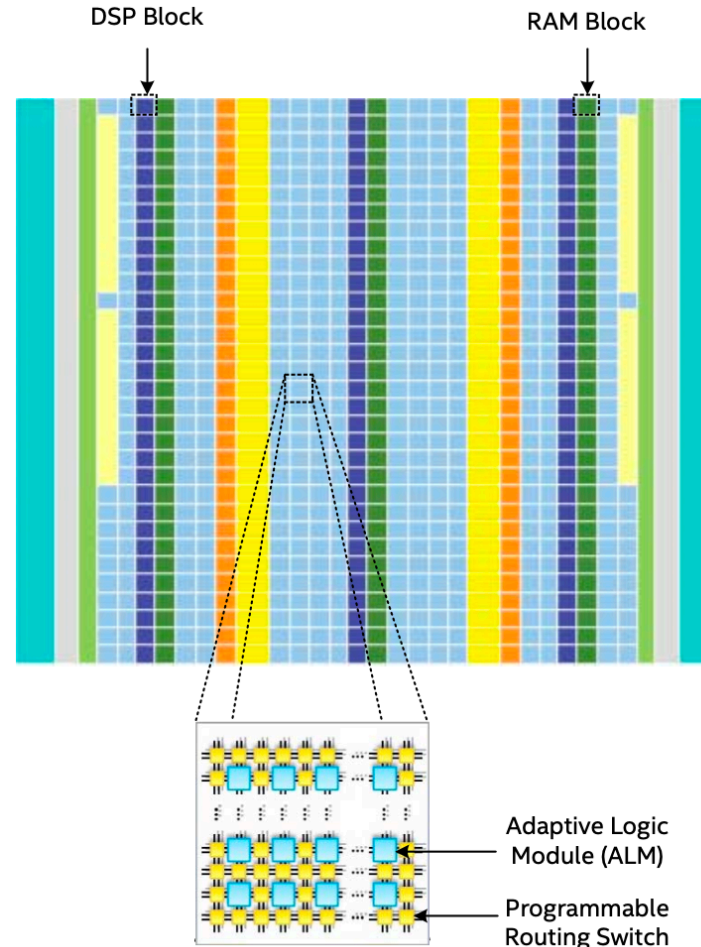
- FPGA does not have a fixed datapath
 - that is why it is “field programmable” !
- “Program” the hardware resources
 - You have a lot more control on how your design is mapped
- Function as accelerators to offload intensive computing
- Design methodologies
 - Hardware Description Language
 - High level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



Learning with Purpose

FPGA Architecture

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



FPGA Hardware Design Concepts

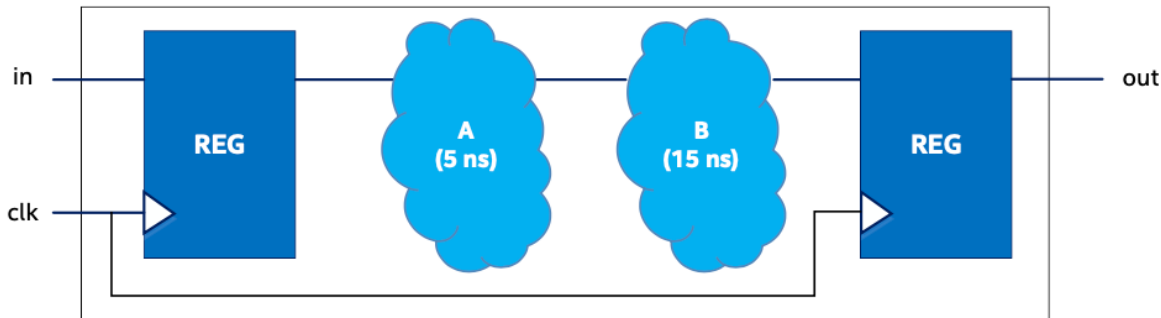
- Maximum frequency f_{MAX}
- Latency
- Pipelining
- Throughput
- Datapath
- Control path
- Occupancy



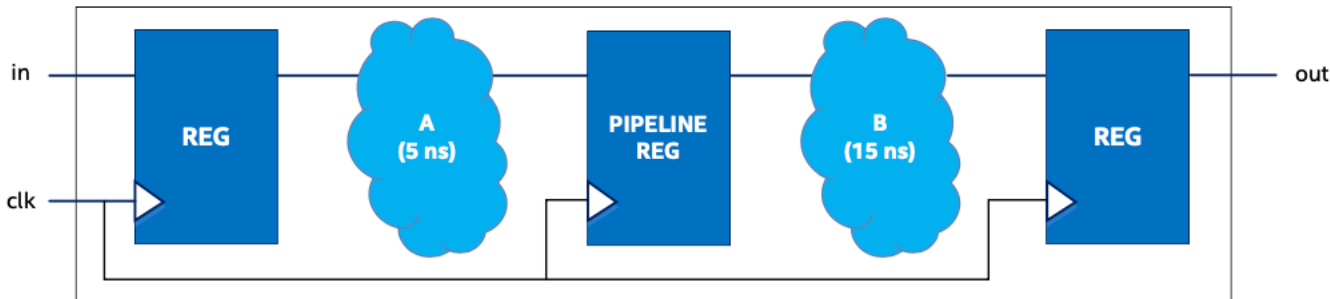
Learning with Purpose

Pipelining to improve f_{MAX}

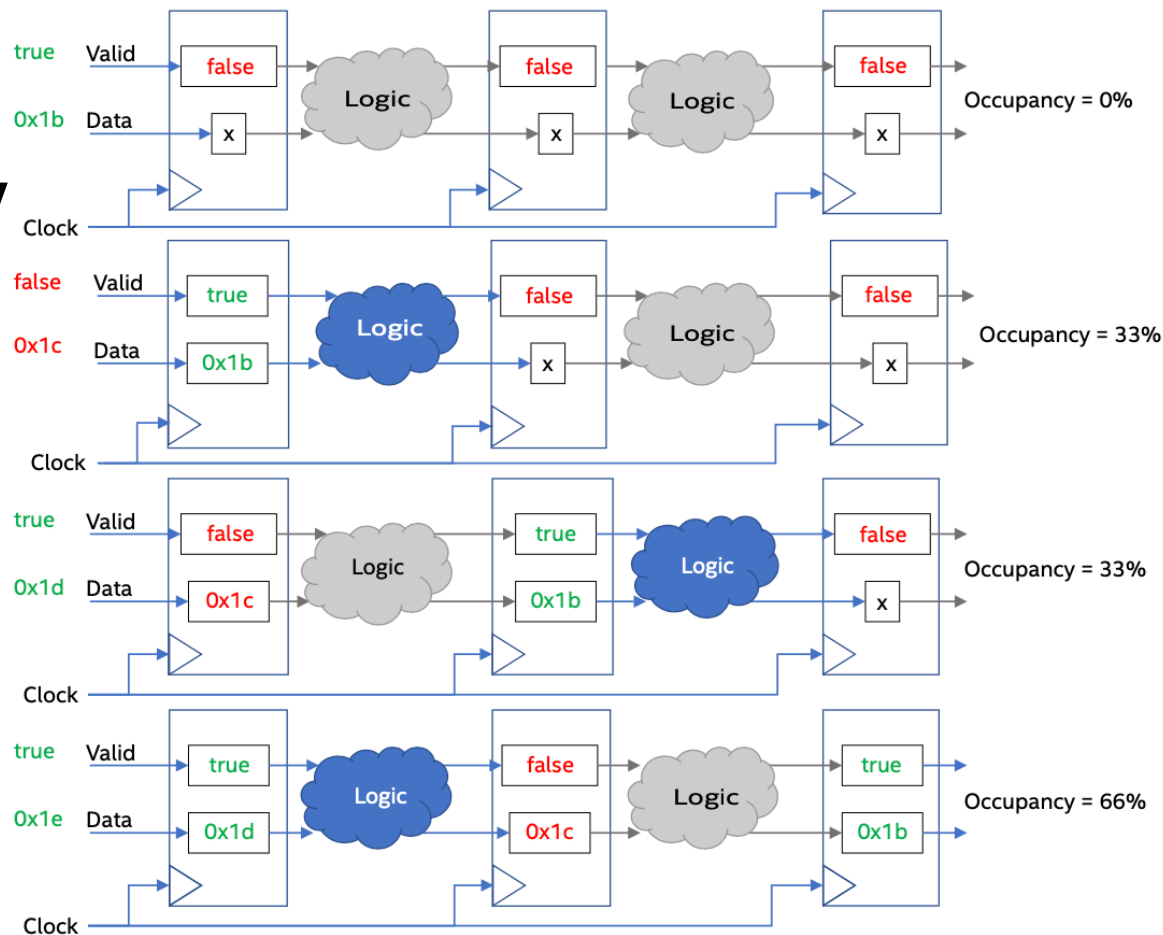
unpipelined
 $f_{MAX} = 50\text{MHz}$



pipelined
 $f_{MAX} = 66.7\text{MHz}$



Occupancy



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

A part of the DPC++ Tutorial Series

Prof. Yan Luo

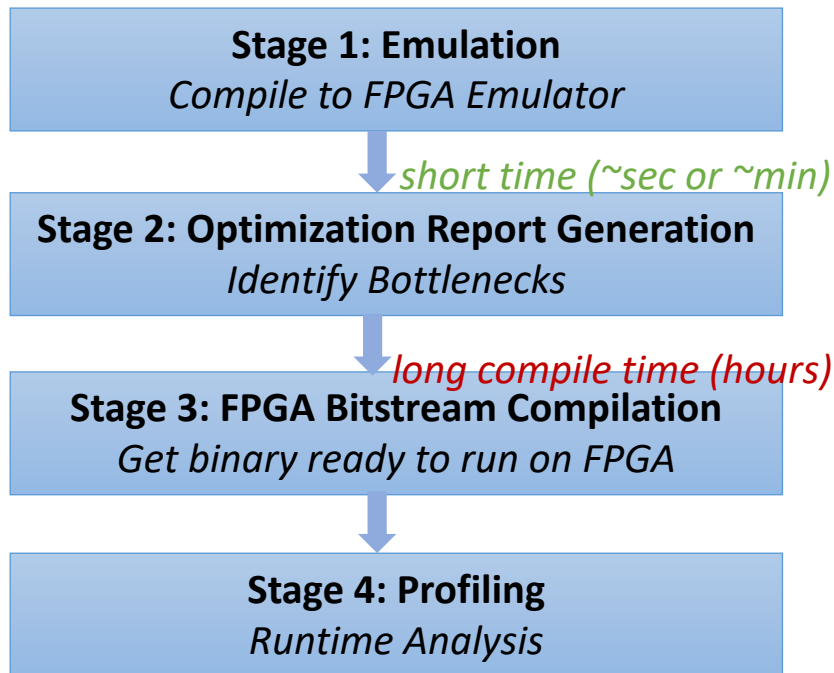
Acknowledgement

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
3. revise the design to eliminate bottlenecks
4. repeat 1,2,3
5. profiling to analyze runtime performance



Learning with Purpose

Demonstration

Look through compilation report of Matrix Multiplication example



Learning with Purpose