# FPGA Development Flow with Intel® oneAPI

## oneAPI – 가속 컴퓨팅을 개발하기 위한 스마트한 방식

2022. 04. 15.
MOASYS

# oneAPI Smart Development Series (2021)

1. Introduction to Intel oneAPI for HPC and AI-DL
   - https://www.allshowtv.com/detail.html?idx=474

2. Benchmarking the Performance of oneAPI on Heterogeneous Computing Platforms
   - https://www.allshowtv.com/detail.html?idx=660

3. Optimization and GPU Offloading Workflow with Intel oneAPI
   - https://www.allshowtv.com/detail.html?idx=826

4. Leveraging Intel® oneDNN for AI Workload
   - https://www.allshowtv.com/detail.html?idx=909

# oneAPI Smart Development Series (2022)

1. FPGA Development Flow with Intel® oneAPI Base Toolkit
   - https://www.allshowtv.com/detail.html?idx=995

2. OpenMP Offload with Intel® oneAPI HPC Toolkit
   - TBA

3. Essential DPC++ Optimization Techniques for Accelerators
   - TBA

4. Introduction to Intel® oneAPI Rendering Toolkit
   - TBA

intel software

moasys

# Contents

- ### oneAPI in 2022
  - FPGA development workflow
  - OpenMP GPU offload
  - oneAPI rendering toolkits

- ### FPGA Development Flow with one API Toolkits
  - FPGA Emulation
  - FPGA Optimization Report
  - FPGA Bitstream compilation

- ### Case Study:
  - Accelerating  Memory  Bound AI Inference Workloads

- ### Conclusion

intel software

moasys

# Overview of oneAPI for Heterogenous Computing



- Support diverse accelerator devices (XPU) such as CPU, GPU and FPGA
- Continuously evolving specifications for high performance computing and machine learning
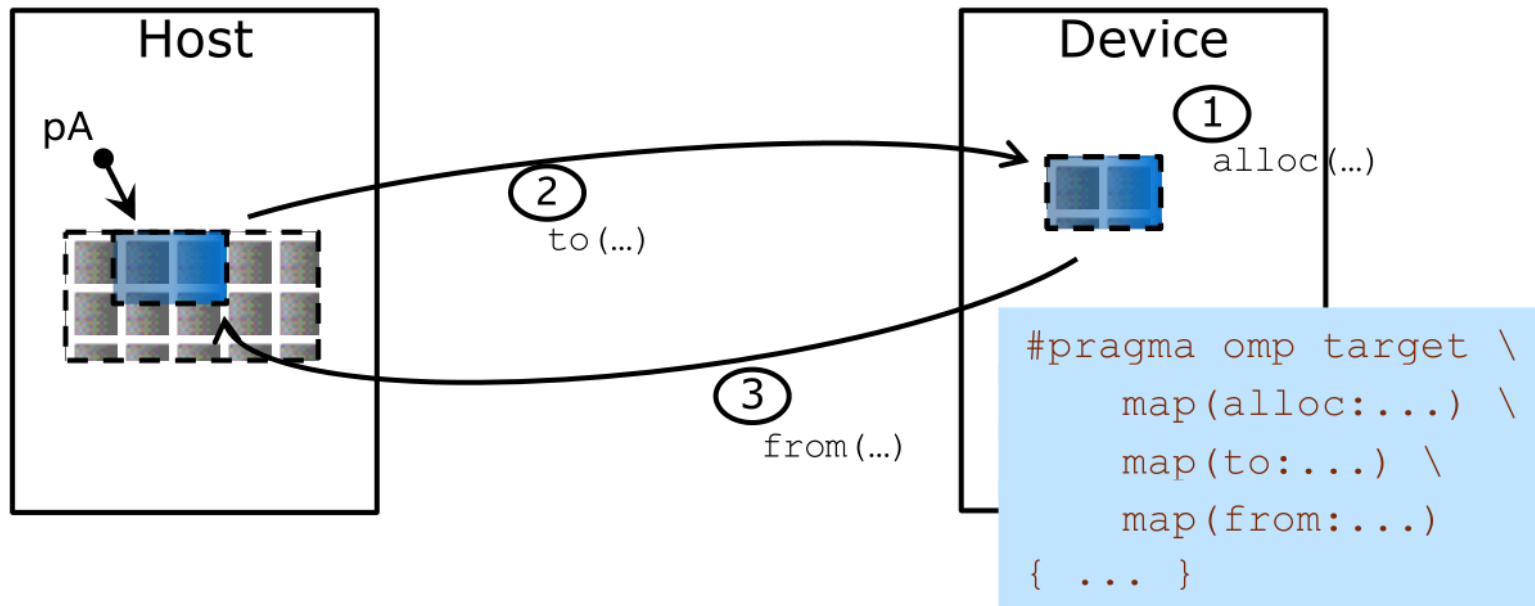
# oneAPI Industry Initiative Progress



| | 2019 | 2020 | 2021 | | | |
|---|---|---|---|---|---|---|
| | | | Q1 | Q2 | Q3 | Q4 |

**Industry Initiative**

- Industry Initiative Announced — 0.5 oneAPI Spec
- oneAPI 1.0 Spec.
- oneAPI AI Tech Advisory Board setup
- oneAPI 1.1 Provisional Spec. — Level Zero, oneDNN Graph, Advanced Ray Tracing, oneVPL
- oneAPI Prov Spec.
- oneAPI Prov Spec.
- oneAPI 1.1 Spec.

**Industry Adoption**

- CodePlay delivers oneAPI support for Nvidia GPUs, + oneMKL support
- Univ. Heidelberg delivers DPC++ AMD GPU support
- Huawei extends DPC++ support for Ascent AI chipset
- Argonne, OLCF to provide SYCL support for AMD GPUs
- Fujitsu & Riken: oneDNN open source on ARM for Fugaku supercomputer
- Bittware FPGA deploys oneAPI training
- Sbercloud launches with oneAPI for AI development
- NERSC, ALCF, Codeplay to enhance LLVM SYCL GPU compiler for Nvidia GPUs
- oneAPI Developer Summits + Geo-targeted sessions — IWOCL ■ Int'l SuperComputing ■ Intel Innovation ■ SuperComputing

Aurora Blade — Building Block for the ExaScale Supercomputer — 1 oneAPI — Argonne — ENERGY — Hewlett Packard Enterprise — intel

> 2 Exaflops (2022)
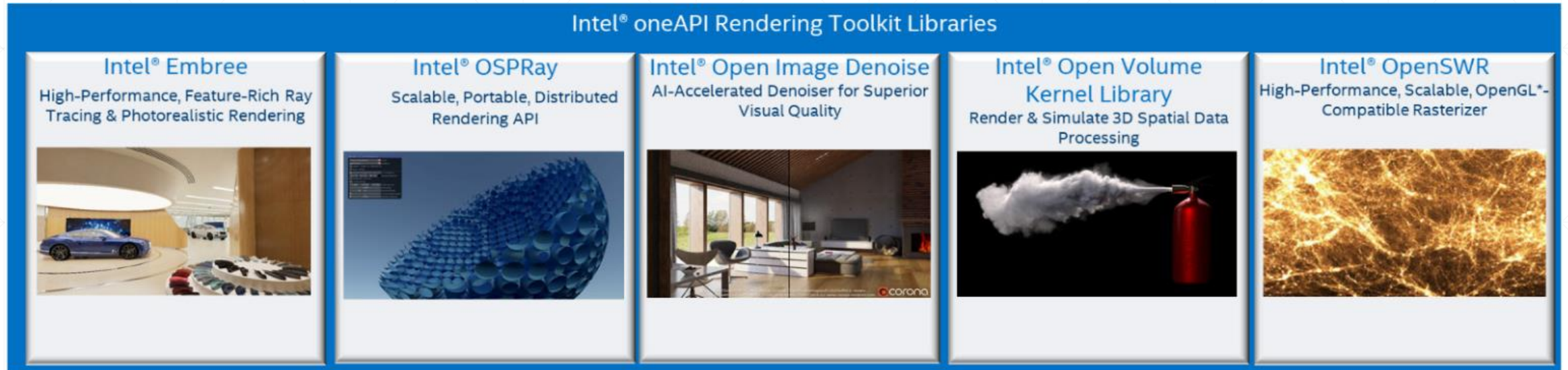
- ▪ **Cross-Vendor implementations:**
  - ▪ ARM CPU: Fujitsu & Riken (Japan)
  - ▪ NVIDIA GPU: CodePlay (USA), NERSC (USA), Argonne National Lab (USA)
  - ▪ AMD GPU: Heidelberg Computing Center (Germany), Argonne National Lab (USA), Oak Ridge National Lab(USA)

# Preview: OpenMP Offload Capabilities in oneAPI HPC Toolkit



```
#pragma omp target \
    map(alloc:...) \
    map(to:...) \
    map(from:...)
{ ... }
```
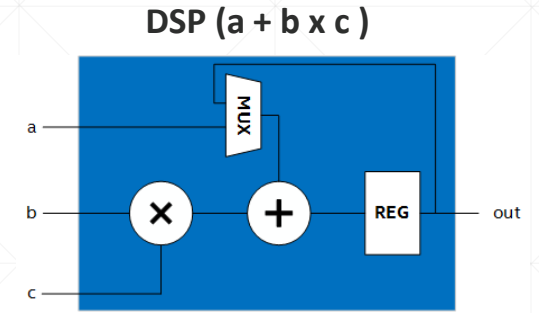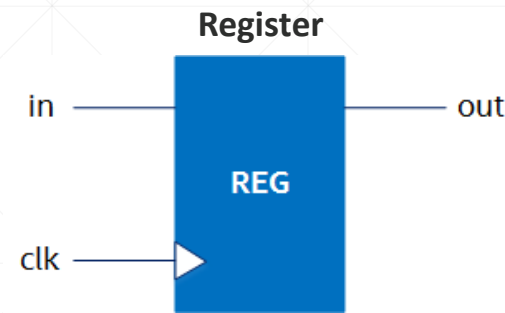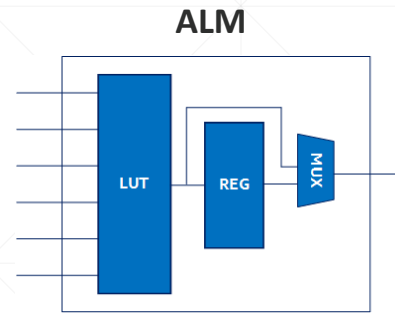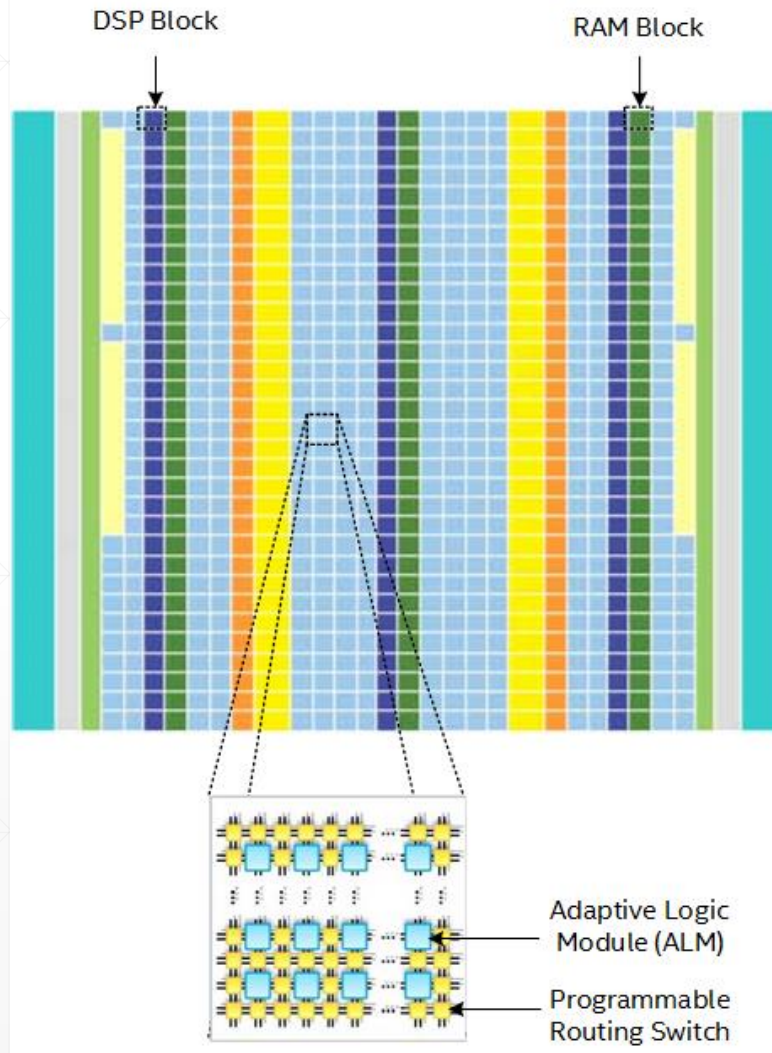
- **OpenMP: Portable, performance and productive parallel programming interface**
- **New compiler technology based on LLVM (DPCPP, ICX and IFX)**
  - Built-in support for $X^e$ GPUs
  - Mixing host and GPU parallelism for maximum efficiency
  - Support for Unified Share Memory (USM)
  - Support for offloading oneMKL routines to GPUs

# Preview: Introduction to oneAPI Rendering Toolkit

## Intel® oneAPI Rendering Toolkit Libraries

**Intel® Embree**
High-Performance, Feature-Rich Ray Tracing & Photorealistic Rendering

**Intel® OSPRay**
Scalable, Portable, Distributed Rendering API

**Intel® Open Image Denoise**
AI-Accelerated Denoiser for Superior Visual Quality

**Intel® Open Volume Kernel Library**
Render & Simulate 3D Spatial Data Processing

**Intel® OpenSWR**
High-Performance, Scalable, OpenGL*-Compatible Rasterizer

- Intel leadership in ray tracing for high performance graphics and compute
- Run at scale, from laptop to workstation and enterprise HPC cluster
  - Hybrid rendering on Intel CPUs and $X^e$ GPUs
  - Visualize huge dataset interactively
  - Persistent memory support for reducing restart times and improving I/O performance
  - Combined volume and geometric rendering
  - High-fidelity visualization including shadows, ambient occlusion, global illumination, motion blur

intel software

moasys

# FPGA Architecture Overview



**ALM**

**Register**

**DSP ( a + b x c )**

- **FPGA (Field-Programmable Gate Array)**
  - Reconfigurable semiconductor integrated circuit (IC).
  - Adaptive logic module (ALM):
    - Basic building block of FPGA
    - Look-up table (LUT) and output register to build a logic circuit
  - Register:
    - Basic storage element of FPGA: input, output and clock signal (clk)
    - Output is synchronized *every* clock cycle
  - Digital processing block (DSP):
    - Support for common fixed-point and floating-point arithmetic

# FPGA Design Concept: Basics



Delay: 20 ns / Latency: 2 Clks / Fmax: 50 Mhz          Delay: 15 ns / Latency: 3 Clks / Fmax: 66.6 Mhz

- **Critical Path**
  - The path between any two consecutive registers with the highest delay
- **Maximum frequency ($f_{max}$)**
  - The maximum rate of output registry update, defined as **1/crital_path_delay**
- **Latency**
  - How many clock cycles to complete operations in a digital circuit
- **Pipeline**
  - Adding registers to the critical path, which decreases the amount of logic between each register
  - Increasement in $f_{max}$ but also in latency
- **Datapath**
  - A chain of registers and combinational logic in a digital circuit that performs computations.
  - Input Reg -> A -> Pipeline Reg -> B -> Output Reg

intel software          moasys

```
q.submit([&](handler& h) {
    h.single_task([=](){
        for (int i = 0; i < kSize; ++i) {
            r[i] = a[i] + b[i]
        }
    });
});
```

```
for (int i = 1; i < n; i=>) {
    c[i] = c[i - 1) + b[i];
}
```

**pipeline register**

Naive implementation

Compiler-optimized implementation

- **Advantages of using single-work item kernel over ND-range kernel**
  - Loop iteration as basic unit of execution, identical to standard C/C++ code
  - Automatic creation of pipeline register by DPCPP compiler
  - Automatic dependencies resolution by DPCPP compiler
    - c[i - 1] stored in *pipeline register* and feedback loop is automatically created
- **Initiation Interval (II):**
  - The number of clock cycles between the launch of successive loop iterations.
  - For ideal pipeline with high throughput: II = 1

# FPGA vs. GPU: AI Inference for Autonomous Driving System



- GPU disadvantages:
  - Batch inference requires addition synchronization
  - Large batch size -> higher throughput and latency
  - Small batch size -> lower throughput and latency
- FPGA advantages:
  - Batch-less inference through pipeline parallelism (first in first out)
  - Consistent and predictable latency and high throughput
- Train on GPU and inference on FPGA for autonomous driving application

intel software

moasys

# Intel FPGA Accelerators



- CPU offloads complex computation tasks to FPGA connected via PCI express
- Targeted workload: streaming analytics, fintech, genomics, artificial intelligence

moasys

# Why is FPGA Compilation Different?



- Key difference of development flow between FPGA and CPU (or GPU)
  - Only ahead-of-time compilation is support due to *very* time consuming FPGA bitstream compilation

intel software

moasys

# Types of DPC++ FPGA Compilation

- **FPGA Emulation**
  - Generation of FPGA emulator image
  - Fastest method to verify the correctness of the code on CPU
  - Timing on emulator does not corresponds to FPGA hardware

- **FPGA Static Report**
  - Generation of FPGA early image (not executable)
  - Visualization of structure created by FPGA
  - Performance and bottleneck
  - Estimation of resource utilization

- **FPGA Hardware Compile and Profiling**
  - Generation of FPGA hardware image (bitstream)
  - Require Intel® FPGA Add-On for oneAPI Base Toolkit
  - Target Intel® Aria 10 GX, Stratix 10 SX or any custom board



**FPGA Development Flow**

- Coding
- Seconds — Emulation (Functional Valdation)
- Minutes — Static Reports
- Hours — Full Compile and Hardware Profiling
- Deploy

intel software

moasys

# FPGA Compilation Flags

```
# FPGA emulator image
dpcpp -fintelfpga -DFPGA_EMULATOR fpga_compile.cpp -o fpga_compile.fpga_emu

# FPGA early image (with optimization report): default board
dpcpp -fintelfpga -Xshardware -fsycl-link=early fpga_compile.cpp -o fpga_compile_report.a

# FPGA early image (with optimization report): explicit board
dpcpp -fintelfpga -Xshardware -fsycl-link=early -Xsboard=intel_s10sx_pac:pac_s10 fpga_compile.cpp -o fpga_compile_report.a

# FPGA hardware image: default board
dpcpp -fintelfpga -Xshardware fpga_compile.cpp -o fpga_compile.fpga

# FPGA hardware image: explicit board
dpcpp -fintelfpga -Xshardware -Xsboard=intel_s10sx_pac:pac_s10 fpga_compile.cpp -o fpga_compile.fpga
```

| Flag | Explanation |
|------|-------------|
| -fintelfpga | Performs ahead-of-time (offline) compilation for FPGA |
| -DFPGA_EMULATOR | Preprocessor for device selection |
| -Xshardware | Instructs the compiler to target FPGA hardware (Default: FPGA emulator) |
| -fsycl-link=early | Instructs the compiler to stop after creating the FPGA early image (and associated optimization report) |
| -Xsboard=<bsp:variant> | Specifies the FPGA board variant (Default: Target Intel® Aria 10 GX) |
| -Xsfast-compile | Allows faster compile time but at a cost of reduced performance of the compiled FPGA hardware image |
| -reuse-exe=<exe_name> | Instruct the compiler to attempt to reuse the existing FPGA device image for fast compilation |

# FPGA Device Selector

```cpp
#include <CL/sycl.hpp>

// FPGA device selectors are defined in this utility header, along with
// all FPGA extensions such as pipes and fpga_reg
#include <sycl/ext/intel/fpga_extensions.hpp>

int main() {
// Select either:
//  - the FPGA emulator device (CPU emulation of the FPGA)
//  - the FPGA device (a real FPGA, can be used for simulation too)
#if defined(FPGA_EMULATOR)
    ext::intel::fpga_emulator_selector device_selector;
#else
    ext::intel::fpga_selector device_selector;
#endif

queue q(device_selector);
...

// Print out the device information.
std::cout << "Running on device: "
          << q.get_device().get_info<info::device::name>() << "\n";
...
}
```

- FPGA emulator and FPGA are distintive devices
- Only a head of time compilation supported, making the *default_selector* less useful

# FPGA Sample Code: VecAdd

```cpp
std::vector<int> vec_a(kSize), vec_b(kSize), vec_r(kSize);

for (int i = 0; i < kSize; i++) {
    vec_a[i] = rand();
    vec_b[i] = rand();
}

// Create buffers to share data between host and device.
// The runtime will copy the necessary data to the FPGA device memory when the kernel is launched.
buffer buf_a(vec_a);
buffer buf_b(vec_b);
buffer buf_r(vec_r);

// Submit a command group to the device queue.
q.submit([&](handler& h) {
    // The SYCL runtime uses the accessors to infer data dependencies.
    // A "read" accessor must wait for data to be copied to the device before the kernel can start.
    accessor a(buf_a, h, read_only);
    accessor b(buf_b, h, read_only);
    accessor r(buf_r, h, write_only, no_init);

    // The task's for loop is executed in pipeline parallel on the FPGA
    h.single_task<VectorAdd>([=]() [[intel::kernel_args_restrict]] {
        for (int i = 0; i < kSize; ++i) {
            r[i] = a[i] + b[i];
        }
    });
});
```

# FPGA Fast Compilation: Reuse FPGA Hardware Image

```
# Initial compilation
dpcpp -fintelfpga -Xshardware <files.cpp> -o out.fpga

# Subsequent recompilation
dpcpp <files.cpp> -o out.fpga -reuse-exe=out.fpga -Xshardware -fintelfpga
```

- **AOT compilation takes several hours to generate FPGA hardware image**
  - **-reuse-exe=<exe_name>** to instruct the compiler to reuse the existing FPGA device image.
  - Separating the host and device code into separate files, i.e device link method

- **Reuse FPGA image**
  - If no change in device code is detected by compiler, only host code is recompiled (few minutes)
  - If changes in device is detected by compiler, FPGA device is fully recompiled (few hours)

- **Caveats:**
  - Strong coupling between host and devices code within single source file
  - Very limited in practice, probably only suitable for simple design
  - Recompilation triggering when compiler falsely detect change in device code

# FPGA Fast Compilation: Device Link Method



```
# Compile device code
dpcpp -fintelfpga -fsycl-link=image kernel.cpp -o dev_image.a -Xshardware

# Compile host code (redo after host code changed)
dpcpp -fintelfpga host.cpp -c -o host.o

# Create device link (redo after host code changed)
dpcpp -fintelfpga host.o dev_image.a -o fast_recompile.fpga
```

# Hough Transformation for Boundary Detection



Image Space

Parameter Space

Input (edge detected)

Output (line detected)

```
for all x
    for all y
        if edge point is at (x,y)
            for all thetas:
                rho = x * cos(theta) + y * sin(theta)
                accumulator(rho, theta) += 1
            end
        end
    end
end
```

- Consider a set of $(x_i, y_i)$ corresponds to detected edges of image object:
  - $y_i = mx_i + b$ (point/image space) $\longleftrightarrow$ $b = -mx_i + y_i$ (line/parameter space)
  - Detection of vertical lines: parameterized representation $(\rho, \theta)$

intel software

moasys

# FPGA Performance Optimization: Hough Transformation

```cpp
queue_event = device_queue.submit([&](sycl::handler &cgh) {
//Create accessors
    auto _pixels       = pixels_buf.get_access<sycl::access::mode::read>(cgh);
    auto _sin_table    = sin_table_buf.get_access<sycl::access::mode::read>(cgh);
    auto _cos_table    = cos_table_buf.get_access<sycl::access::mode::read>(cgh);
    auto _accumulators = accumulators_buf.get_access<sycl::access::mode::read_write>(cgh);

    //Call the kernel
    cgh.single_task<class Hough_transform_kernel>([=]() {
        for (uint y=0; y<HEIGHT; y++) {
            for (uint x=0; x<WIDTH; x++){
                unsigned short int increment = 0;
                if (_pixels[(WIDTH*y)+x] != 0) {
                    increment = 1;
                } else {
                    increment = 0;
                }

                for (int theta=0; theta<THETAS; theta++){
                    int rho = x*_cos_table[theta] + y*_sin_table[theta];
                    _accumulators[(THETAS*(rho+RHOS))+theta] += increment;
                }
            }
        }
    });
});
```



- _accumulators matrix is updated by scanning over possible value of thetas

# FPGA Performance Optimization: Early Image Report



- Generate early image (not executable) with architecture optimization report:
  - `dpcpp -fintelfpga -Xshardware` **-fsycl-link=early** `main.cpp hough_transform.cpp -o fpga_compile_report.a`
- HTML report: *fpga_compile_report.prj/reports/report.html*

# FPGA Performance Optimization: Avoid Aliasing of Kernel Arguments



- Pointer aliasing occurs when the same memory location is accessed using different names
- **kernel_args_restrict** for more aggressive compiler optimizations and improved FPGA performance

# FPGA Performance Optimization: Local Memory



- Slow data retrieval from global memory, increasing initiation interval (II)
- Local memory is referred to as on-chip memory created from FPGA's RAM blocks

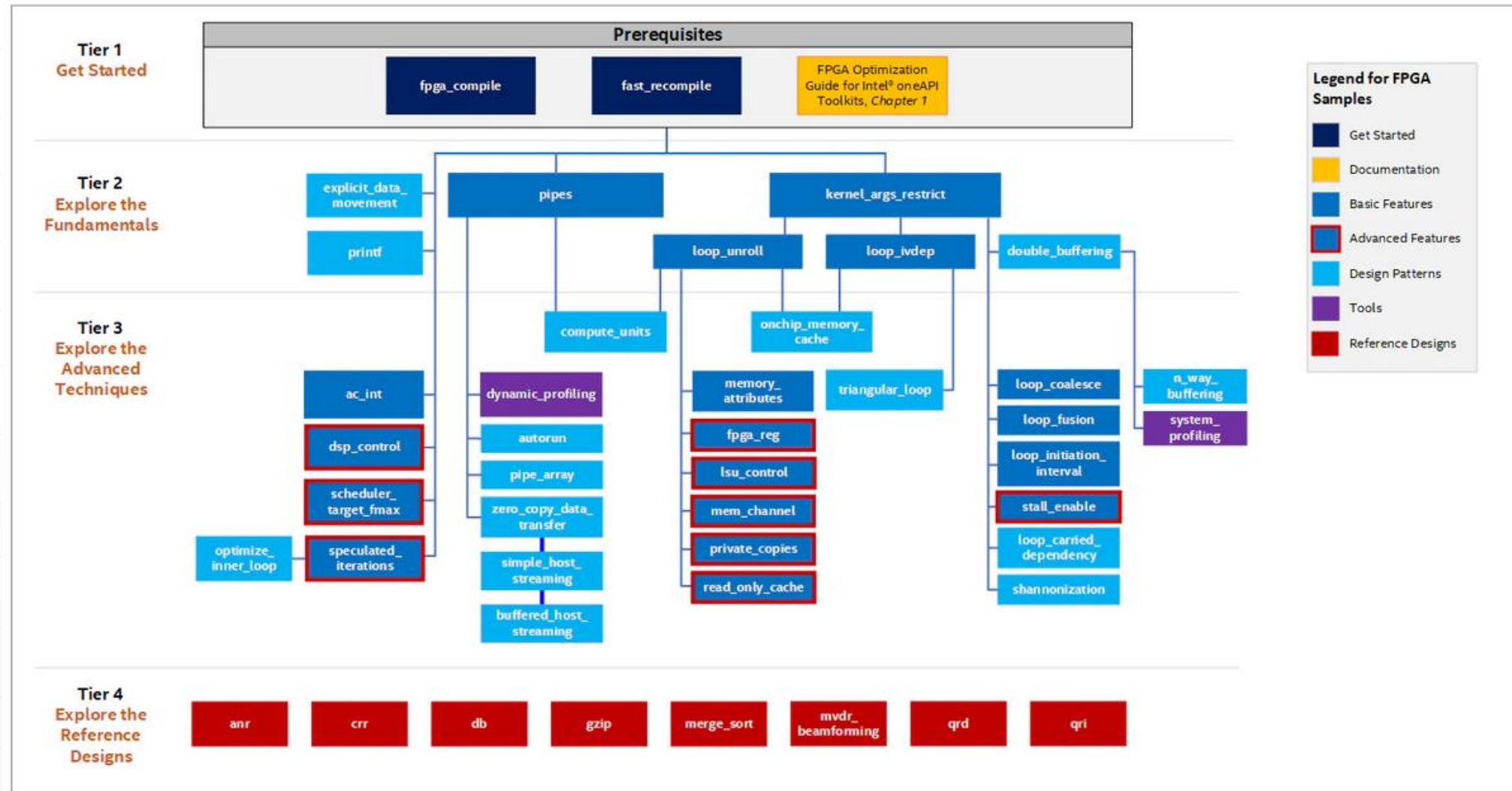# FPGA Performance Optimization: Loop Unroll and IVDEP



- Improved parallelism by duplicating the compute logic with loop unrolling
- Ignore loop carries dependency based on iteration distance across multiple loops

# FPGA Performance Optimization: Hough Transformation Summary

```
cgh.single_task([=]() [[intel::kernel_args_restrict]] {
    //Load from global to local memory
    short accum_local[RHOS*2*THETAS];
    for (int i = 0; i < RHOS*2*THETAS; i++) {
        accum_local[i] = 0;
    }
    for (uint y=0; y<HEIGHT; y++) {
        for (uint x=0; x<WIDTH; x++){
        unsigned short int increment = 0;
            if (_pixels[(WIDTH*y)+x] != 0) {
                increment = 1;
            } else {
                increment = 0;
            }

            #pragma unroll 32
            [[intel::ivdep]]
            for (int theta=0; theta<THETAS; theta++){
                int rho = x*_cos_table[theta] + y*_sin_table[theta];
                accum_local[(THETAS*(rho+RHOS))+theta] += increment;
            }
        }
    }
    //Store from local to global memory
    for (int i = 0; i < RHOS*2*THETAS; i++) {
        _accumulators[i] = accum_local[i];
    }
}):
```
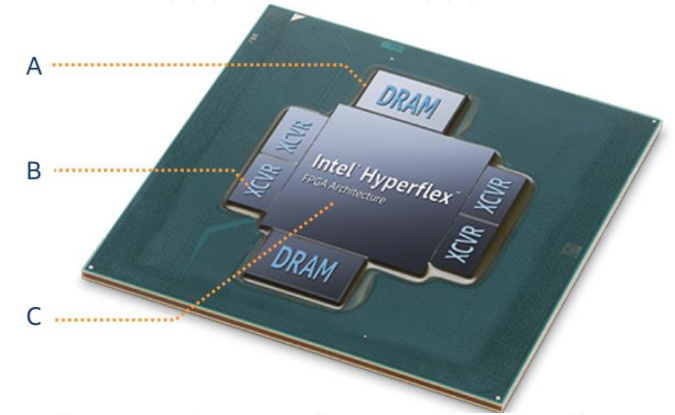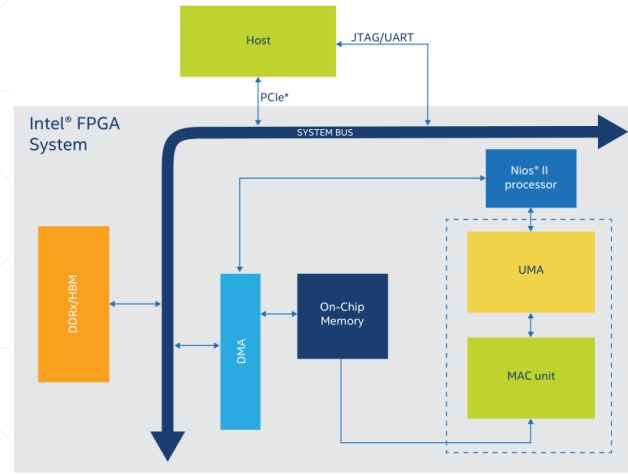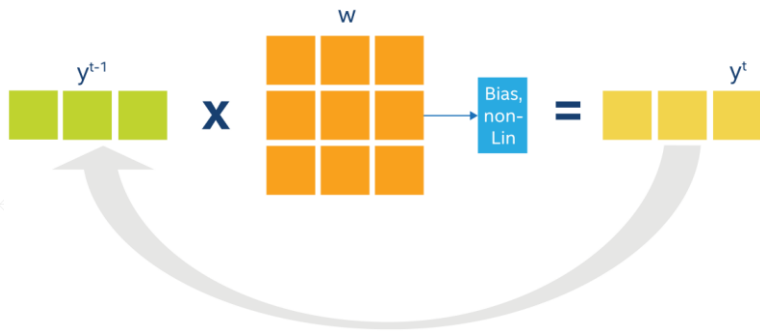
# Explore DPC++ Through Intel® FPGA  Code Samples

https://www.intel.com/content/www/us/en/developer/articles/code-sample/explore-dpcpp-through-intel-fpga-code-samples.html
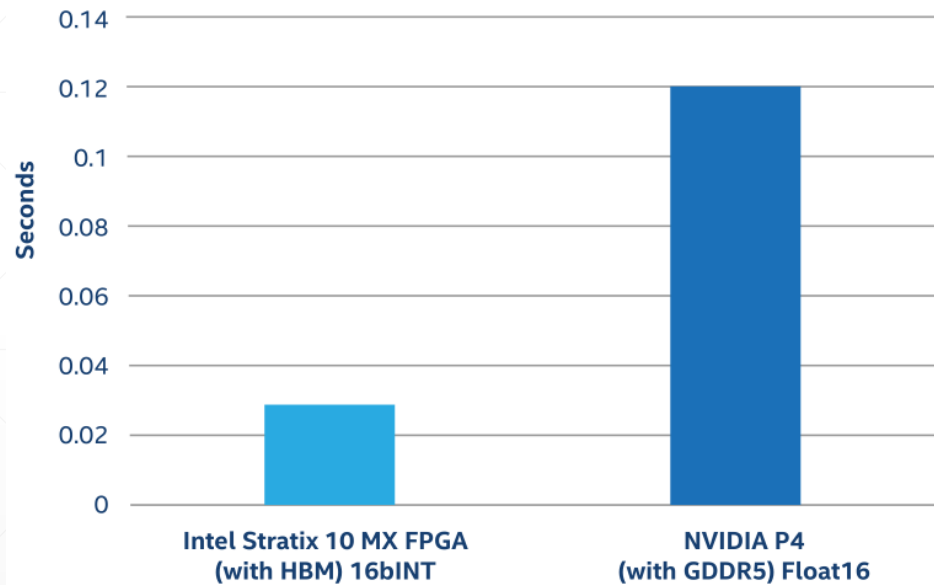
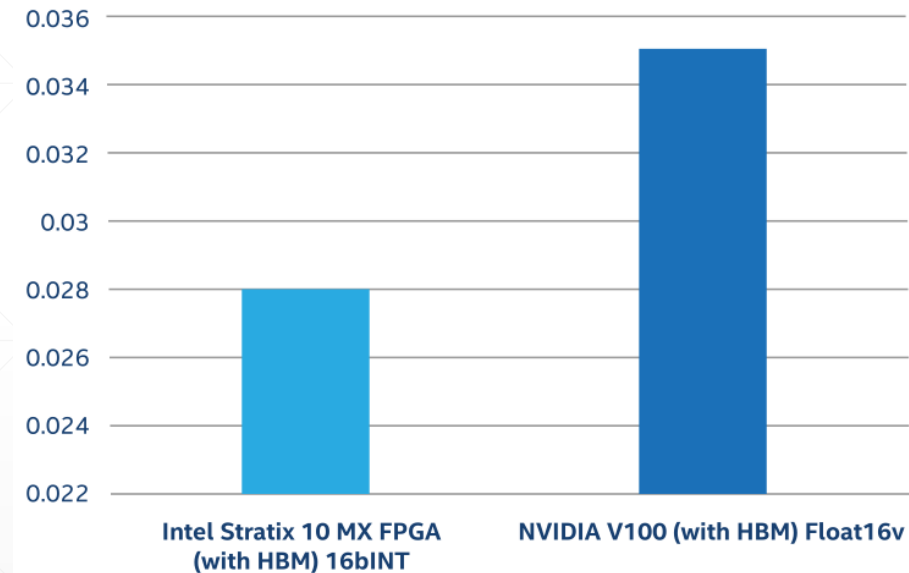# Case Study: Accelerating  Memory  Bound AI Inference Workloads



- **Recurrent Neural Network (RNN):**
  - Reflect the influence of past data on current data
  - Finance, genome mapping, speech AI, automatic speech recognition (ASR), etc
  - Repeated application of weight matrix each time step, i.e low latency required for real time application
  - Suitable for FPGA accelerator thanks to pipeline parallelism
- **Intel Stratix 10 MX FPGA:**
  - Single package of a state-of-the-art Intel® FPGA with Samsung High Bandwidth Memory 2 (HBM2)
  - 10x bandwidth with highest performance per Watt in comparison to SDRAM
  - Programmable high performance AI inference accelerator for Intel FPGAs

# Case Study: Accelerating  Memory  Bound AI Inference Workloads



**Inference Latency**
Intel® Stratix® 10 MX FPGA(PIE) vs.
NVIDIA P4 GPU with GDDR5 Memory

(Mozilla DeepSpeech, 1s of audio, single batch )

**Inference Latency**
Intel® Stratix® 10 MX FPGA(PIE) vs.
NVIDIA V100 GPU with HBM

(Mozilla DeepSpeech, 1s of audio, single batch )

- Mozilla DeepSpeech algorithm was accelerated using Intel FPGA
  - Stratix 10MX has ~ 4x lower latency than NVIDIA P4
  - Stratix 10MX has ~ 20% lower latency than NVIDIA V100 (HBM)

# Conclusion

- **FPGA offers advantages over traditional accelerators such as CPU and GPU**
  - Highly efficiency via pipeline parallelism
  - Data dependency across parallel work automatically resolved by DPCPP compiler

- **FPGA development workflow with oneAPI**
  - Emulation -> optimization report -> bitstream compilation
  - Single-work item kernel instead ND-range kernel
  - Optimization techniques

- **FPGA case study**
  - Stratix 10 MX  offer lower latency than GPUs for RNN workload

- **Contact**: sales@moasys.com
  - GPU, FPGA code migration
  - Optimization and parallelization consultant
  - Customized HPC education

intel software

moasys