



# DEEP LEARNING INFERENCE WITH INTEL® FPGAS

Class 5

# What We Discussed in the Previous Lesson

Using the Deep Learning Accelerator Suite to implement CNN networks on the FPGA with the OpenVINO™ toolkit

The different architectures and primitives that are provided as part of the DLA Suite, how they are implemented on the FPGA and how they can be customized

# Agenda

Introduction and acceleration stack overview

Developing a SW host application

Getting Acceleration Functions

Getting started with the Acceleration Stack

# Objectives

Explain How the acceleration stack enables abstraction of the FPGA in Cloud/Enterprise clusters transparently

Define the process for a host application to discover and configure an FPGA accelerator

Explain where to get accelerated functions for the FPGA

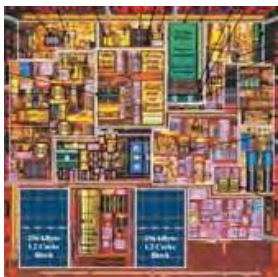
Define the process to create an accelerator function unit for the FPGA

Explain the process for adding an FPGA accelerator card into a cluster and installing the acceleration stack

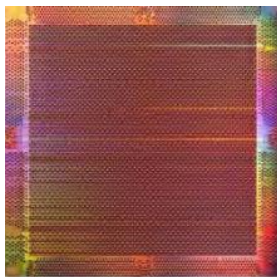
Relate sample designs to get started

# The Future is Heterogenous

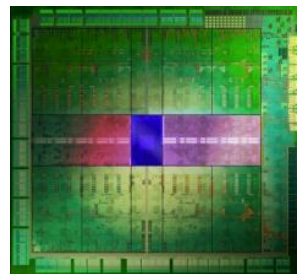
Span from CPU, GPU, FPGA to dedicated devices with consistent programming models, languages, and tools



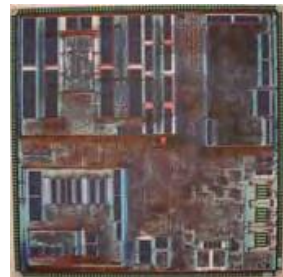
CPUs



GPUs



FPGAs



ASSP

# Separation of Concerns

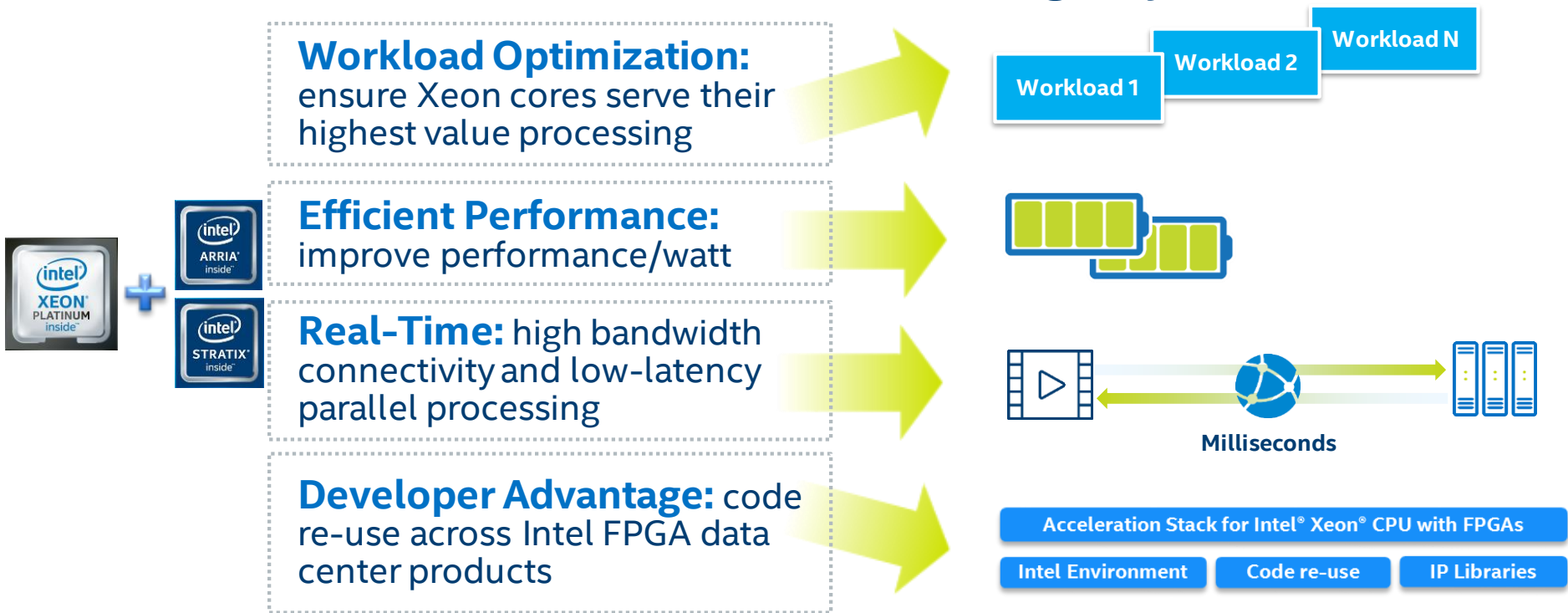
Two groups of developers:

- Domain experts concerned with getting a result
  - Host application developers leverage optimized libraries
- Tuning experts concerned with performance
  - Typical FPGA developers that create optimized libraries

Intel® Math Kernel Library is a simple example of raising the level of abstraction to the math operations

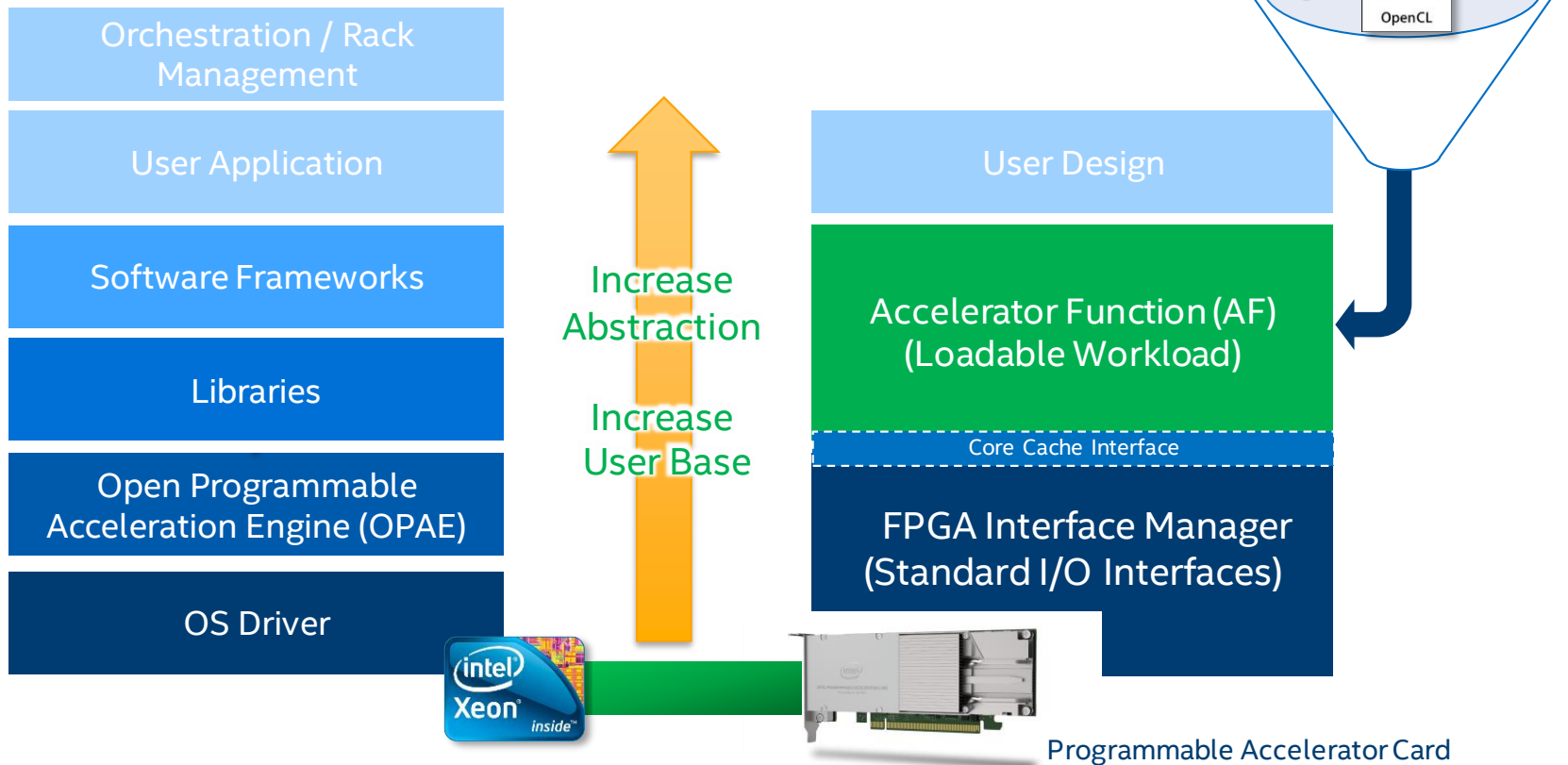
- Domain experts focus on formulating their problems
- Tuning experts focus on vectorization and parallelization

# FPGA Enabled Performance and Agility



FPGAs enhance CPU-based processing by *accelerating algorithms* and *minimizing bottlenecks*

# Using FPGAs Just Got Easier





# What Acceleration Stack enables



**ECOSYSTEM OF FPGA  
WORKLOADS**

**End Application  
User**



**APPLICATION & FPGA  
DEVELOPMENT**

**HW &  
SW Developer**



**FPGA DEPLOYMENT  
& MANAGEMENT**

**Data Center Operator  
Integrated Services  
Vendors**

**Enabled by**



# The Acceleration Stack for Intel® Xeon CPU w/ FPGA

## Applications/ Orchestration

Orchestration / Rack Level Management

User Applications



## Vertical Software Frameworks/Libs

(DL, Networking,  
Genomics, etc.)

Frameworks

Intel Xeon FPGA  
Acceleration Libraries



## Common Infrastructure

✓ *Simplify FPGA  
programming model*

FPGA HW & SW  
Tool Chains

Open Programmable Acceleration Engine  
(OPAE Software API)

Drivers, virtualization, API's, acceleration engine  
Intel FPGA SDK for OpenCL™, Intel Quartus® Prime

Operating Systems

OS Enablement: Linux, ESXi, Windows

## FPGA Images

FPGA Interface  
Manager (FIM)

Loadable AFU image(.gbs)

**IP Libraries:** DLA, GEMM, VirtIO, pHMM  
Compression, Encryption, etc..

## Hardware

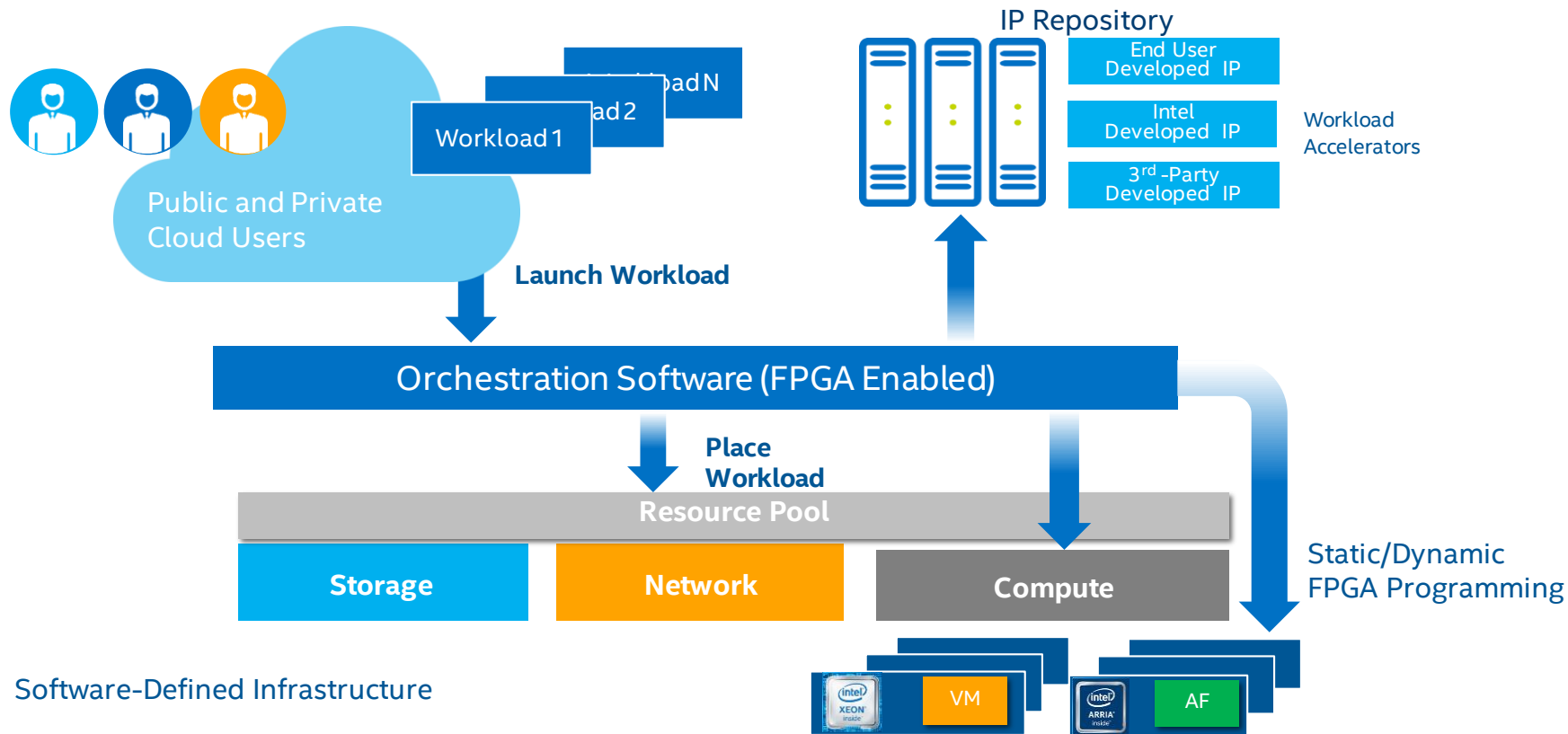
FPGA Platforms (Programmable Acceleration Cards)



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



# Orchestrating FPGA-accelerated applications

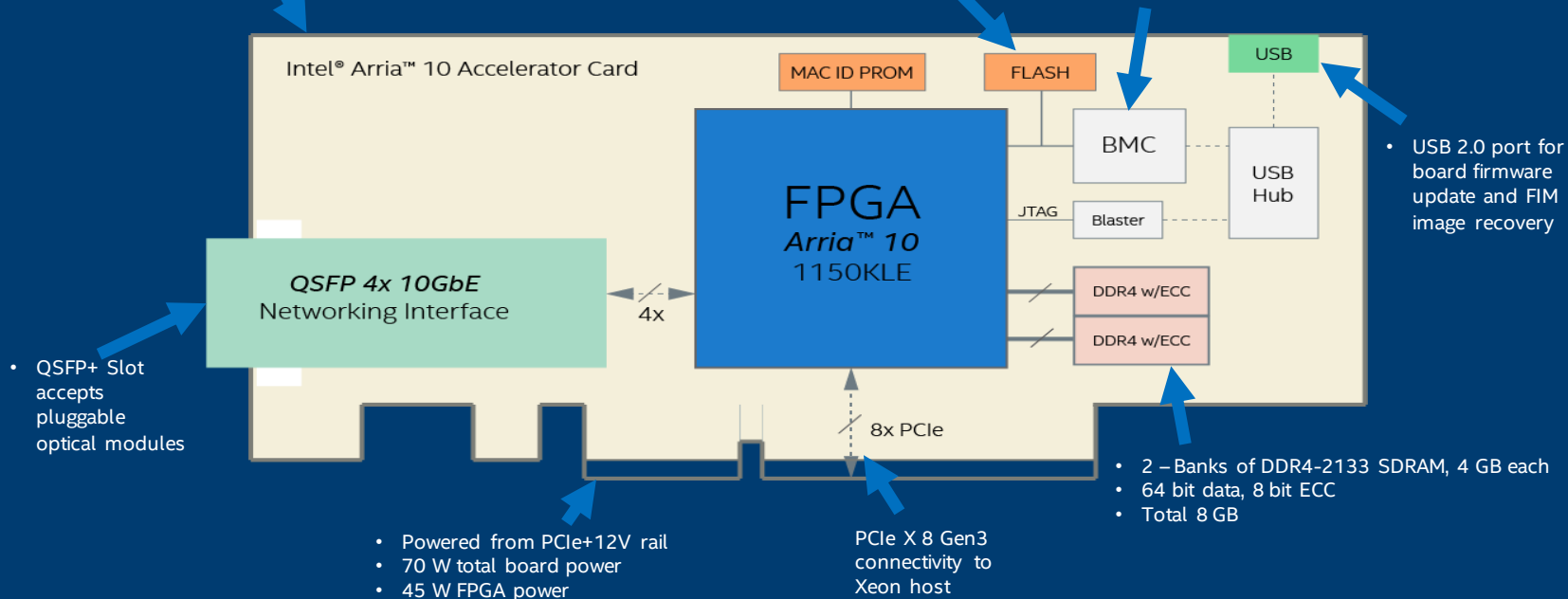


# PROGRAMMABLE ACCELERATION CARD WITH ARRIA 10 FPGA

- Low-profile (half-length, half height) PCIe\* slot card
- 168 mm × 56 mm
- Maximum component height: 14.47 mm
- PCIe × 16 mechanical

- 128 MB Flash
- For storage of FPGA configuration

- Board Management Controller (BMC)
- Server class monitor system
- Accessed via USB or PCIe



# Open Programmable Acceleration Engine (OPAE)

Simplified FPGA Programming Model for Application Developers

## Consistent API across product generations and platforms

- Abstraction for hardware specific FPGA resource details

## Designed for minimal software overhead and latency

- Lightweight user-space library (*libfpga*)

## Open ecosystem for industry and developer community

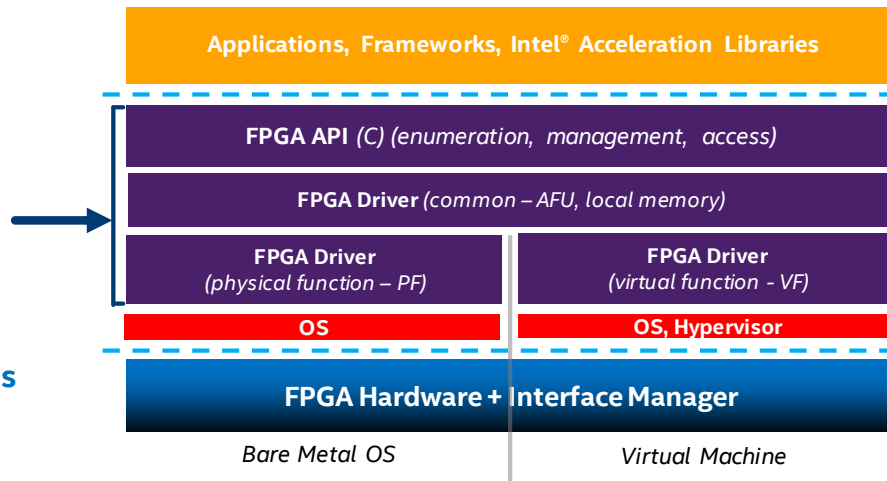
- License: FPGA API (BSD), FPGA driver (GPLv2)

## FPGA driver being upstreamed into Linux kernel

## Supports both virtual machines and bare metal platforms

## Faster development and debugging of Accelerator Functions with the included AFU Simulation Environment (ASE)

## Includes guides, command-line utilities and sample code



Start developing for Intel FPGAs with OPAE today: <http://01.org/OPAE>

# FPGA Accelerator Appears as PCIe Device

## From the OS's point of view

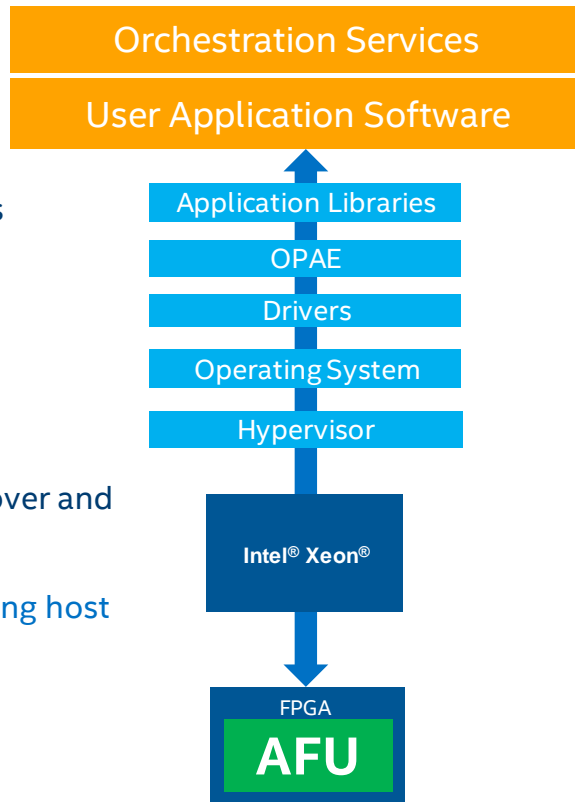
- FPGA hardware appears as a regular PCIe device
- FPGA accelerator appears as a set of features accessible by software programs running on host

## Unified C API model

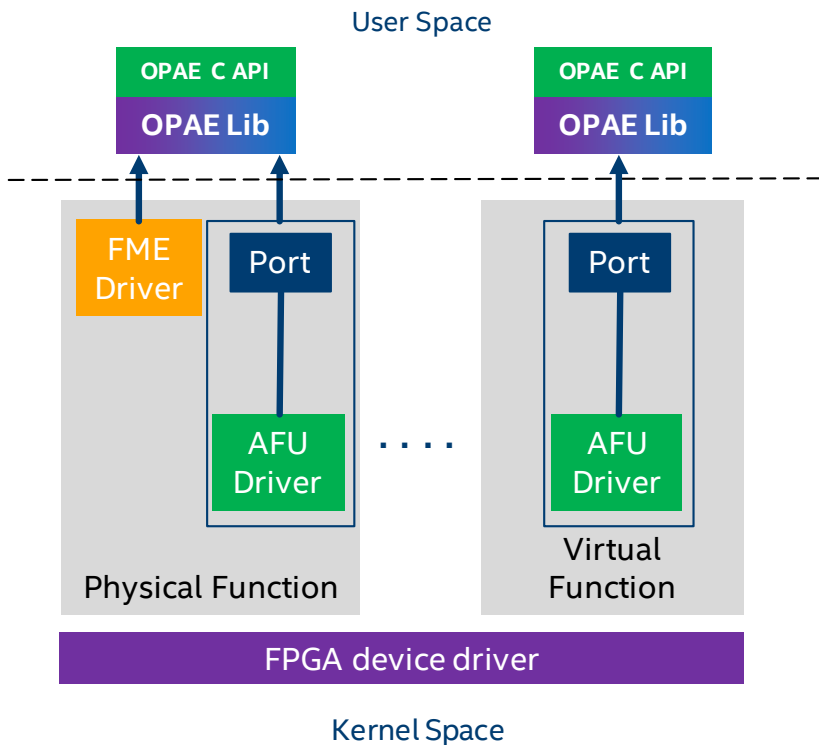
- Supports different kinds of FPGA integration and deployment. (E.g.: A single application can use the FPGA to accelerate certain algorithms)
- Resource management and orchestration services in a data center use to discover and select the FPGA resources and organize them to be used by the workloads

Architecture supports Single Root I/O Virtualization (SROIV) PCIe extension enabling host software to access the accelerator:

- Via a hypervisor/VMM (Virtual Function)
- Bypassing the VMM/Hypervisor Physical Functions



# FPGA Driver Architecture



## FME: FPGA Management Engine Driver

- Static circuits for power/thermal management, reconfiguration, debugging, error reporting, performance counters, etc.

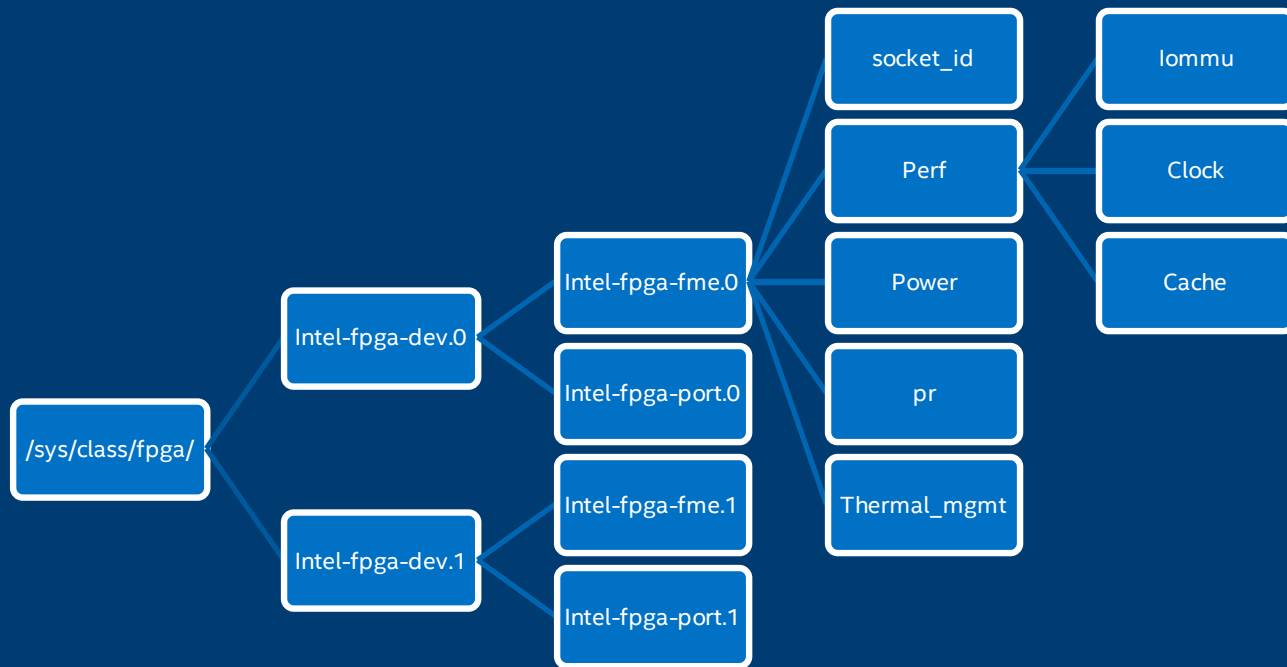
## Port:

- Interface between the static (FIM) and the reconfigurable Acceleration Function (AF) region
- Controls communication from software to the accelerator
- Expose features such as reset and debug
- There may be multiple ports exposed through a VF

## AFU: Accelerator Function Unit Driver

- Exposes a 256KB region as control registers through Port
- Reconfigurable circuits for application specific functions
- User process can share memory buffers with AFU

# HOW HOST APPLICATIONS **ENUMERATE** THE FPGA DEVICE: **SYSFS**



Ex:

2 Intel(R) FPGA devices are installed in the host

Each FPGA device has one FME and one Port (AFU)



The background is a deep blue gradient. On the left side, there is a stylized, glowing blue globe. The globe is covered in a network of white lines and dots, suggesting a global or digital theme. Several circular icons are visible on the globe's surface, including one with three people silhouettes and another with a gear-like pattern. A bright, horizontal light streak cuts across the middle of the globe. The overall aesthetic is high-tech and futuristic.

# **SW APPLICATION DEVELOPMENT FOR ACCELERATION STACK**

# The OPAE Library at a Glance

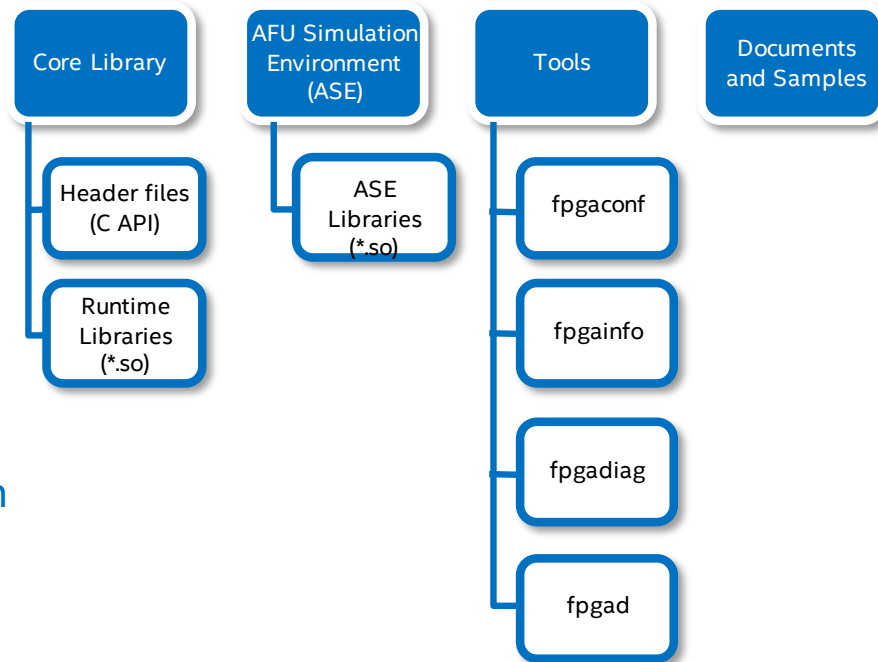
Enumerate, access, and manage FPGA resources through API objects

A common interface across different FPGA form factors

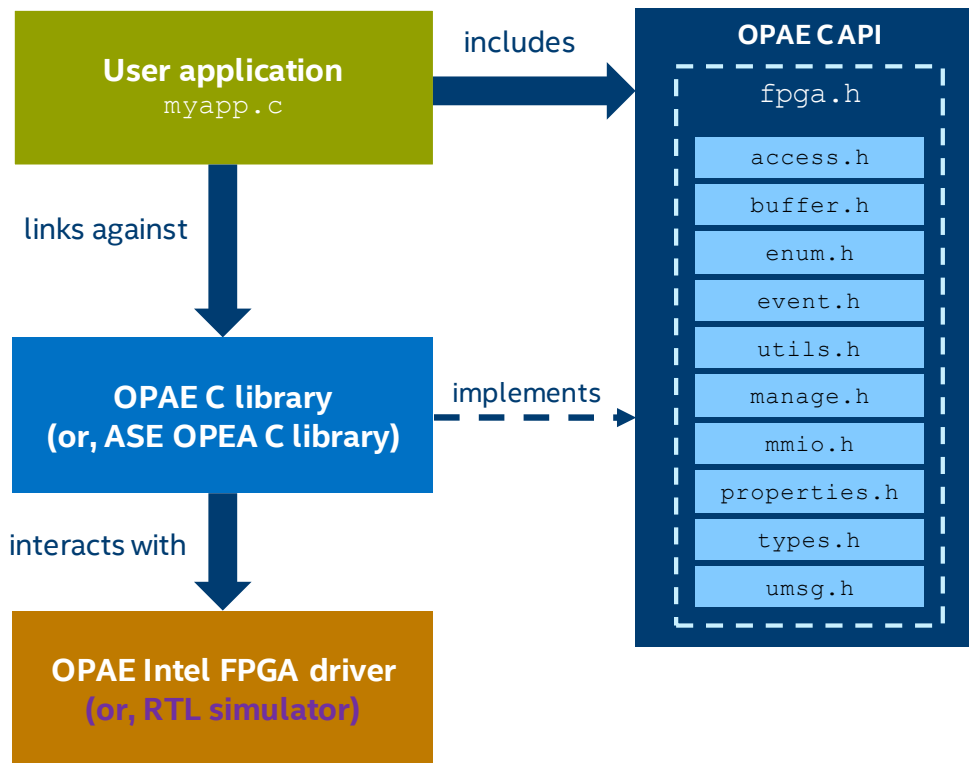
C API designed for extensibility

AFU Simulation Environment (ASE) allows developing and debugging accelerator functions and software applications without an FPGA

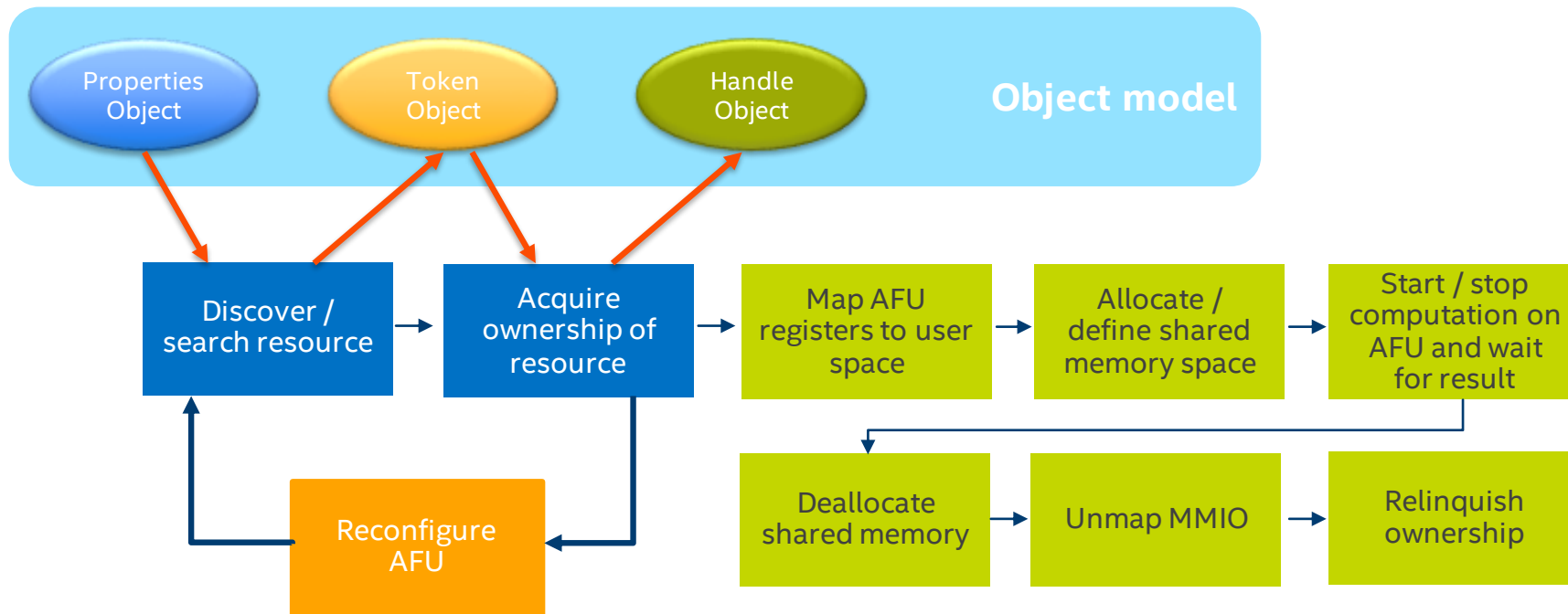
Tools for partial reconfiguration, FPGA hardware information, error reporting, etc.



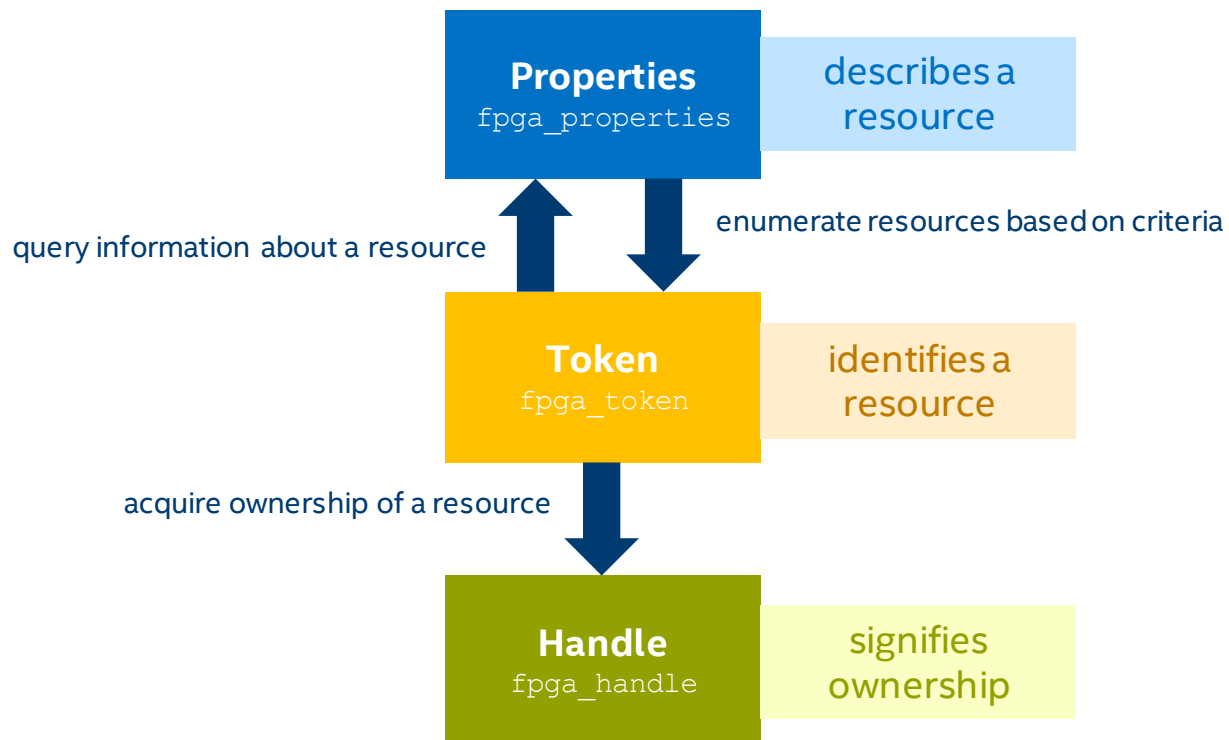
# Application Development with OPAE



# The OPAE Library Programming Model



# The OPAE Object Model



# fpga\_properties Object

An opaque type for a properties object used by application to query and search for appropriate resources

## 2 Object types for FPGA resources

- **FPGA\_DEVICE**
  - Corresponds to physical FPGA device
  - Can invoke management functions
- **FPGA\_ACCELERATOR**
  - Represents an instance of an AFU

Defined in **types.h** file

Property	FPGA*	Accelerator*	Description
Parent	No	Yes	<b>fpga_token</b> of the parent object
ObjectType	Yes	Yes	The type of the resource: either <b>FPGA_DEVICE</b> or <b>FPGA_ACCELERATOR</b>
Bus	Yes	Yes	The bus number
Device	Yes	Yes	The PCI device number
Function	Yes	Yes	The PCI function number
SocketId	Yes	Yes	The socket ID
DeviceId	Yes	Yes	The device ID
NumSlots	Yes	No	Number of AFU slots available on an <b>FPGA_DEVICE</b> resource
BBSID	Yes	No	The FPGA Interface Manager (FIM) ID of an <b>FPGA_DEVICE</b> resource
BBSVersion	Yes	No	The FIM version of an <b>FPGA_DEVICE</b> resource
VendorId	Yes	No	The vendor ID of an <b>FPGA_DEVICE</b> resource
Model	Yes	No	The model of an <b>FPGA_DEVICE</b> resource
LocalMemory Size	Yes	No	The local memory size of an <b>FPGA_DEVICE</b> resource
Capabilities	Yes	No	The capabilities of an <b>FPGA_DEVICE</b> resource
GUID	Yes	Yes	The Global Unique Identifier of an <b>FPGA_DEVICE</b> or <b>FPGA_ACCELERATOR</b> resource
NumMMIO	No	Yes	The number of MMIO space of an <b>FPGA_ACCELERATOR</b> resource
NumInterrupts	No	Yes	The number of interrupts of an <b>FPGA_ACCELERATOR</b> resource
Accelerator State	No	Yes	The state of an <b>FPGA_ACCELERATOR</b> resource: either <b>FPGA_ACCELERATOR_ASSIGNED</b> or <b>FPGA_ACCELERATOR_UNASSIGNED</b>

# Creating `fpga_properties` object

Error code

```
fpga_result fpgaGetProperties(fpga_token token, fpga_properties *prop)
```

Initializes memory pointed at by *prop* to represent properties object

Populates with properties of the resource referred to by *token*

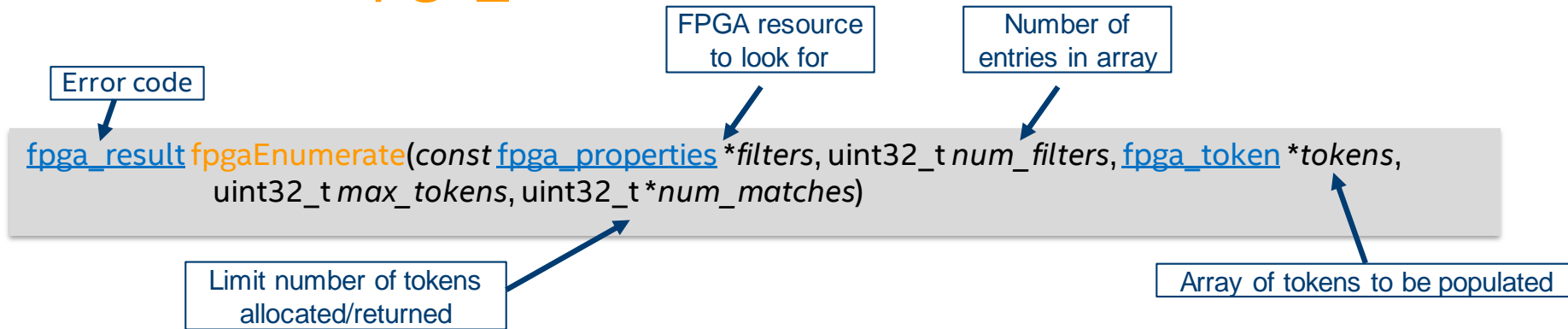
- Passing NULL *token* creates empty properties object which would match all FPGA resources in the enumeration query
- Refine query criteria using `fpgaPropertiesSet*()` functions

Individual properties can be queried using `fpgaPropertiesGet*()` accessor functions

Destroy `fpga_properties` object using `fpgaDestroyProperties()` function

Located in `properties.h` file

# Create an `fpga_token`



Function that searches for FPGA resources in system that match criteria (may be more than one)

- All accelerators assigned to a host interface, all FPGAs of a specific type, etc.

Creates `fpga_token` objects and populates the array with these tokens

- Number of tokens in the returned tokens array, either `max_tokens` or `num_matches` whichever is smaller

Free the memory with tokens no longer needed using the `fpgaDestroyToken()` function

Located in `enum.h` file



# Receive `fpga_handle`

Error code

Allows resource to be opened multiple times  
(`FPGA_OPEN_SHARED`)

```
fpga_result fpgaOpen(fpga_token token, fpga_handle *handle, int flags)
```

Acquires ownership of FPGA resource referenced by *token*

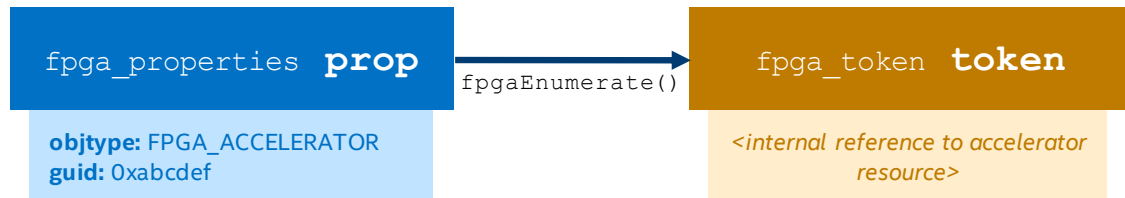
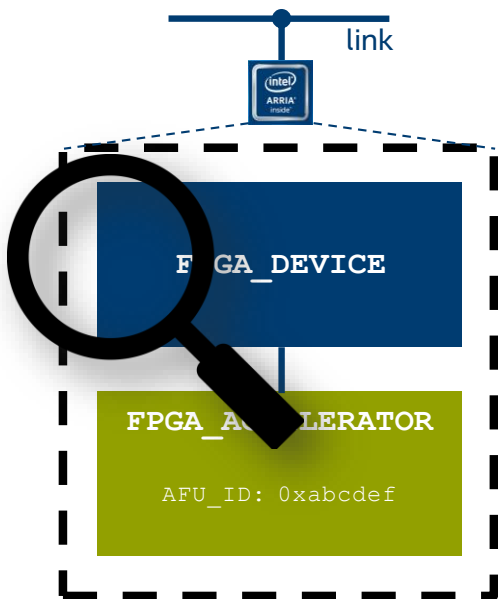
Ownership required to interact with accelerator function

Remains open until `fpga_Close()` function called or process terminates

- Can also reset accelerator using `fpga_Reset()` function

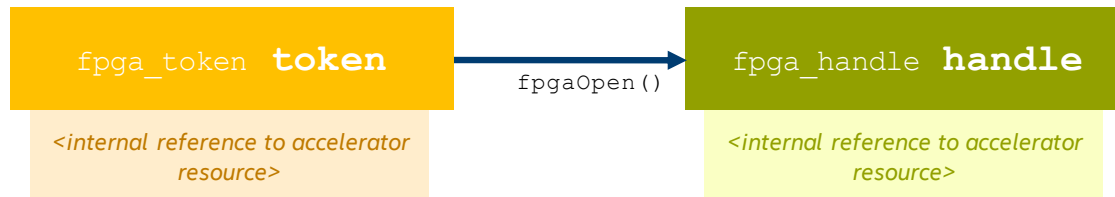
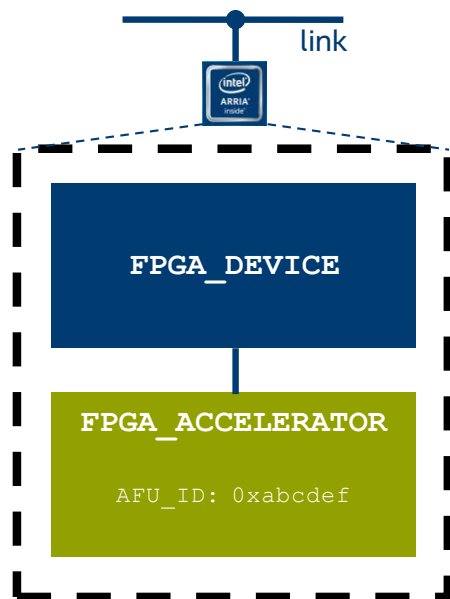
*Located in `access.h` file*

# Enumeration and Discovery



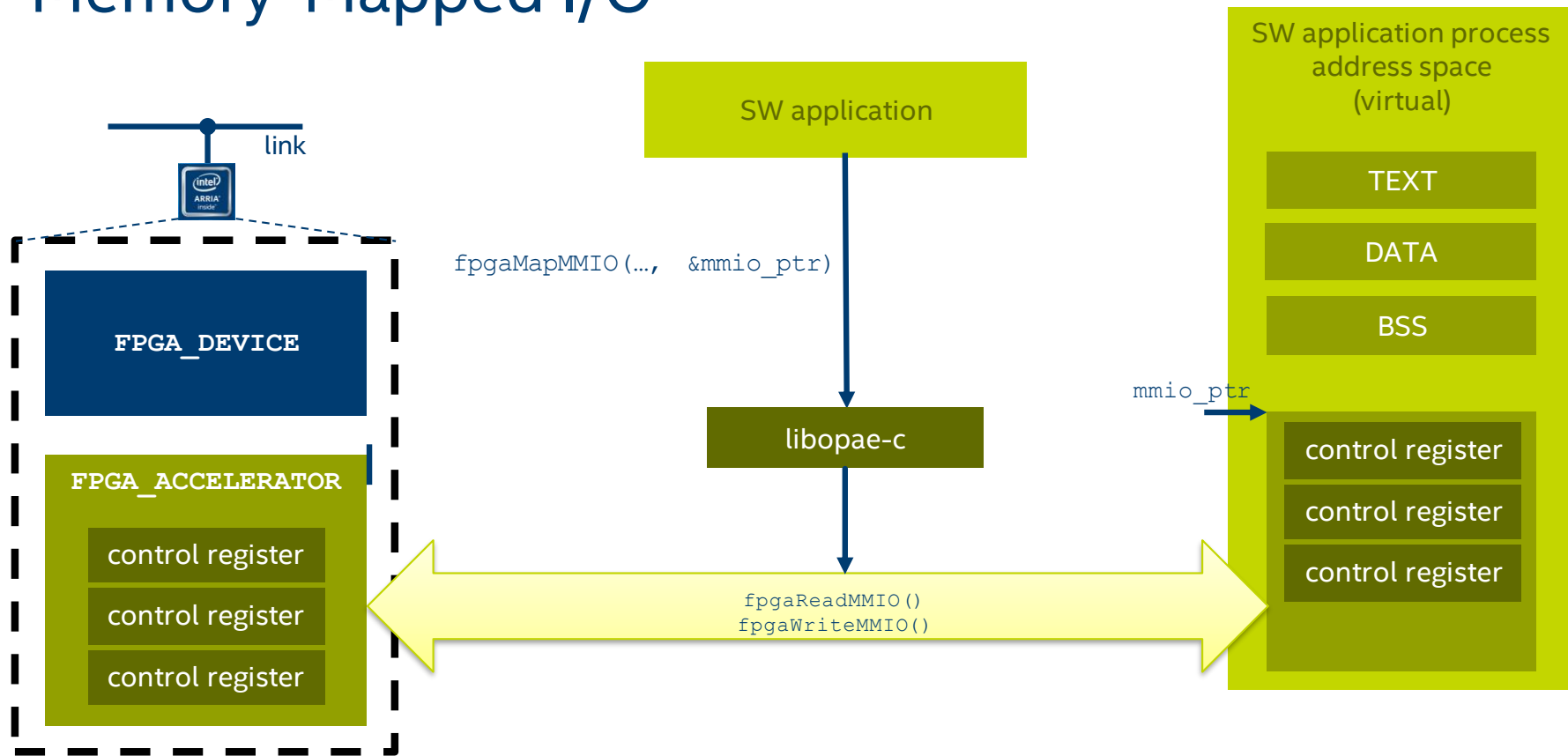
```
fpga_properties prop;  
fpga_token      token;  
fpga_guid       myguid; /* 0xabcdef */  
  
fpgaGetProperties(NULL, &prop);  
  
fpgaPropertiesSetObjectType(prop, FPGA_ACCELERATOR);  
fpgaPropertiesSetGUID(prop, myguid);  
  
fpgaEnumerate(&prop, 1, &token, 1, &n);  
  
fpgaDestroyProperties(&prop);
```

# Acquire and Release Accelerator Resource



```
fpga_token token;  
// ... enumeration ...  
fpga_handle handle;  
  
fpgaOpen(token, &handle, 0);  
.  
.  
.  
fpgaClose(handle);
```

# Memory-Mapped I/O



# MMIO Read/Write

Diagram illustrating the parameters for the `fpgaWriteMMIO64` function:

- Error code (points to `fpga_result`)
- Handle of previously opened accelerator resource (points to `fpga_handle`)
- Number of MMIO space to access (points to `mmio_num`)

```
fpga_result fpgaWriteMMIO64(fpga_handle handle, uint32_t mmio_num, uint64_t offset, uint64_t value)
```

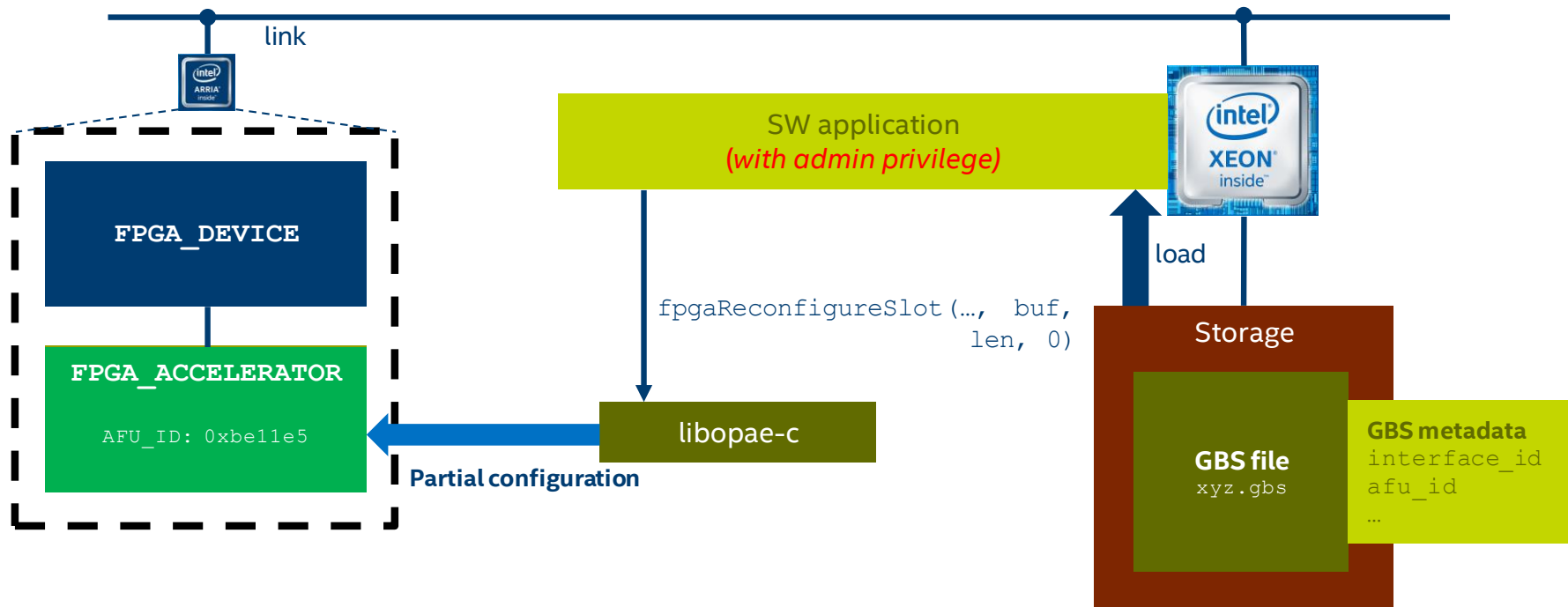
Performs 64bit MMIO space write of value to specified byte offset

- Also supports 32 bit MMIO writes using `fpgaWriteMMIO32()` function

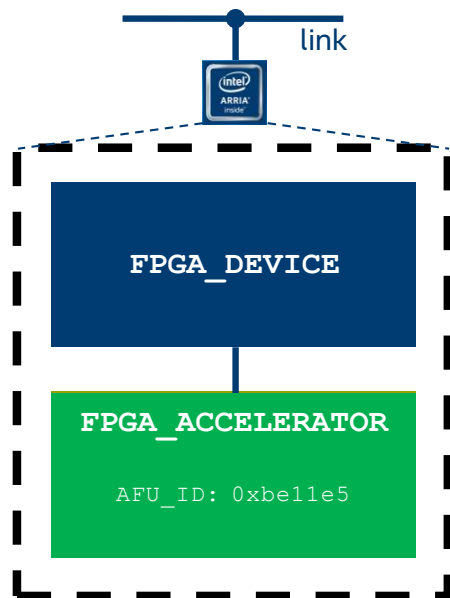
Reads use `fpgaReadMMIO64()` or `fpgaReadMMIO32()` functions

Located in `mmio.h` file

# Management and Reconfiguration



# Management and Reconfiguration

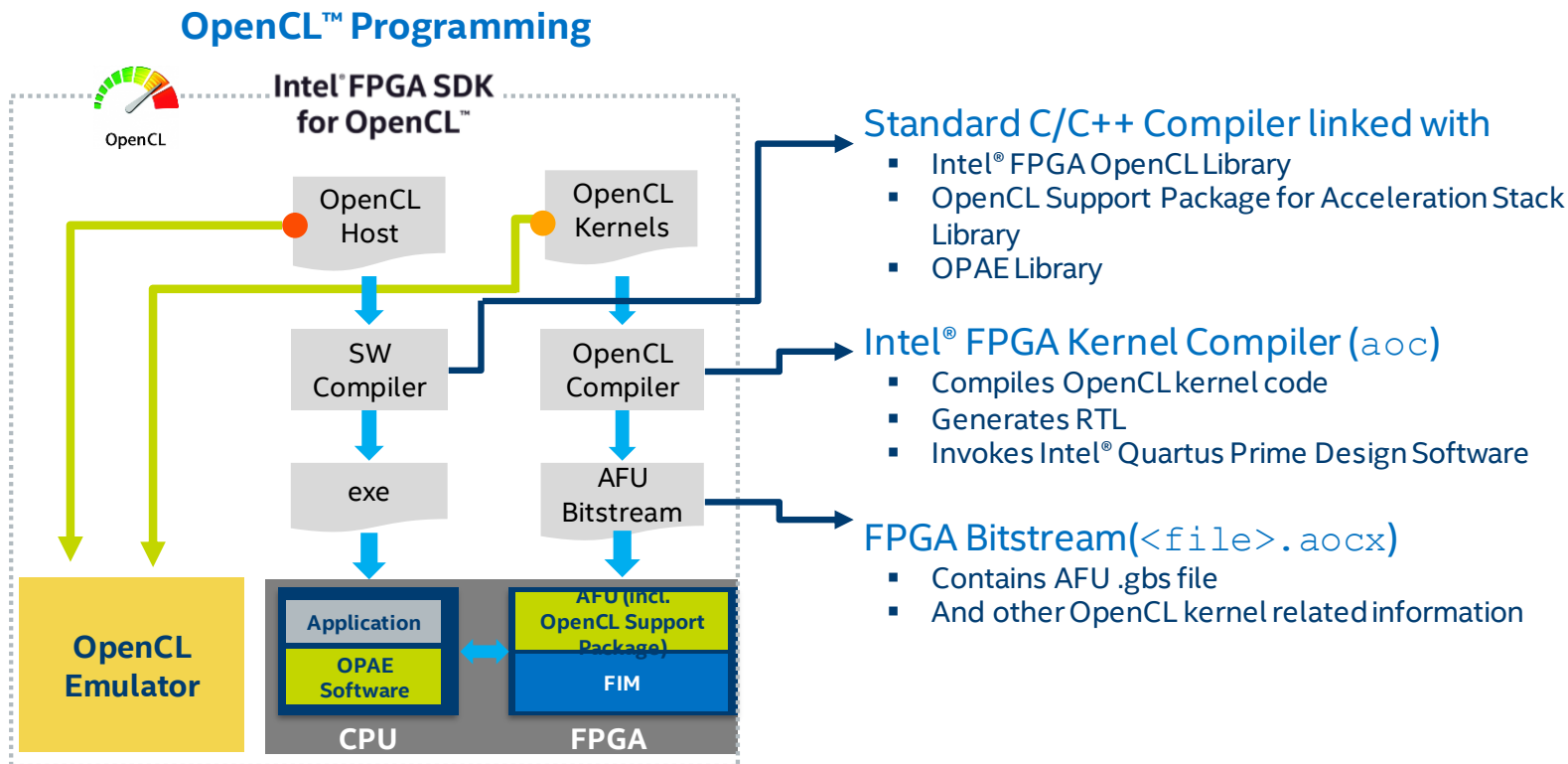


```
fpga_handle handle;           /* handle to device */
FILE          *gbs_file;
void          *gbs_ptr;
size_t        gbs_size;

/* Read bitstream file */
gbs_ptr = malloc(gbs_size);
fread(gbs_ptr, 1, gbs_len, gbs_file);

/* Program GBS to FPGA */
fpgaReconfigureSlot(handle, 0, gbs_ptr, gbs_size, 0);
/* ... */
```

# OpenCL Development Approach





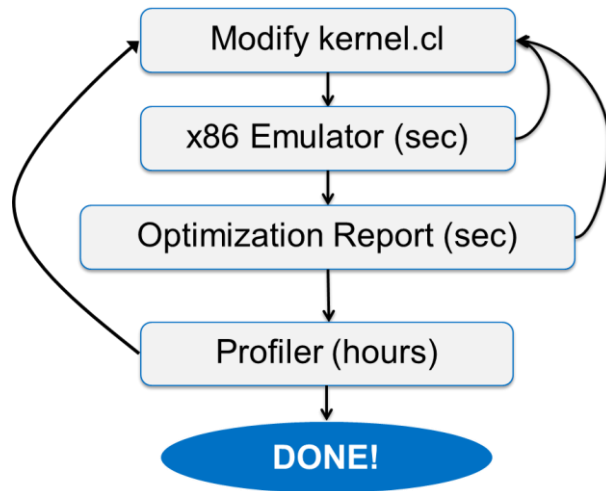
# Software Application Flow using OpenCL™

No different from traditional OpenCL™ flow

- C based development and optimization flow to create AFUs and Host Application
- Standard OpenCL™ FPGA application using the Intel® FPGA SDK for OpenCL
  - FPGA OpenCL™ debug and profiling tools supported
- More information on using [OpenCL with FPGAs](#)

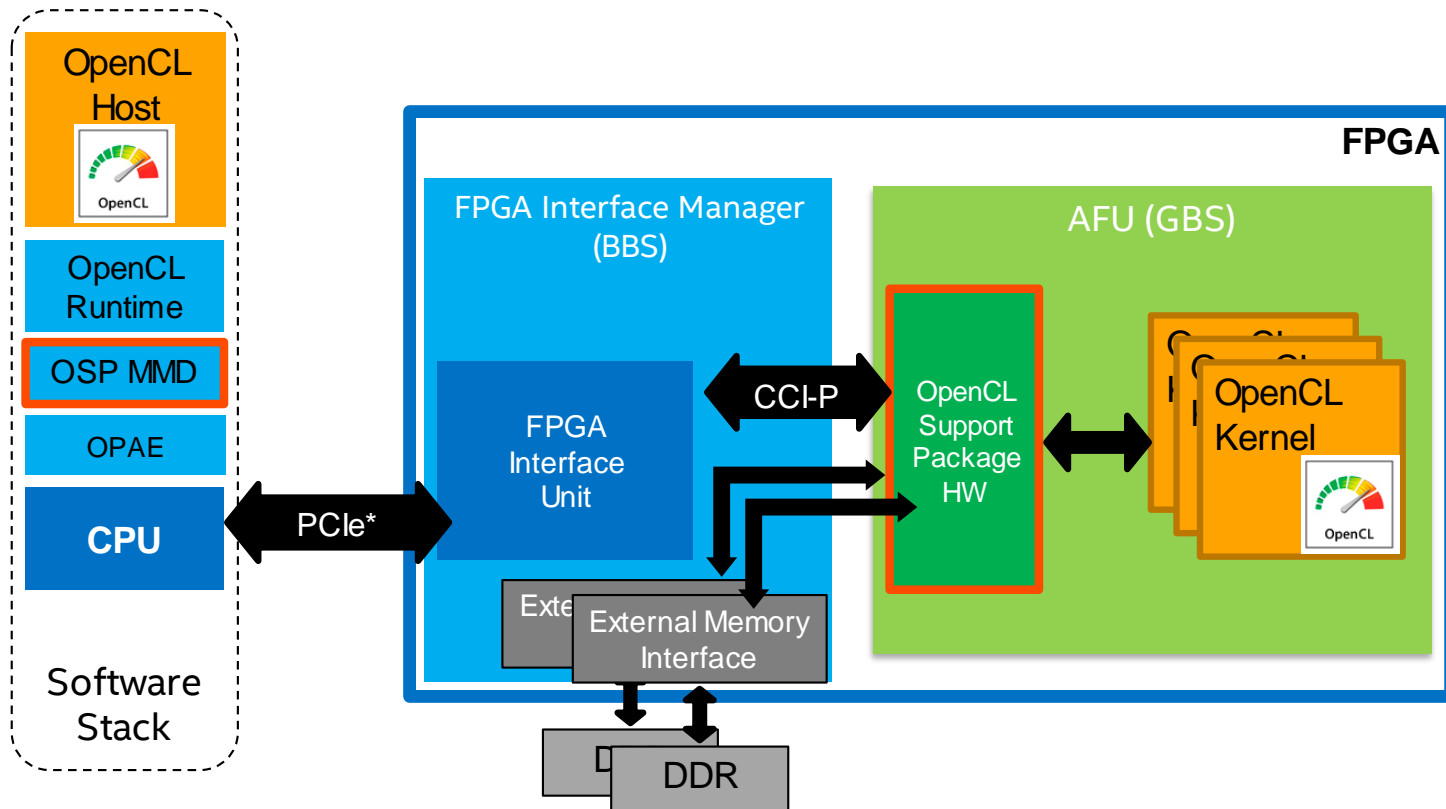
The Acceleration Stack abstracted away from user

- OPAE part of the Host Run-Time
  - Host does not need to interact with OPAE SW directly
- OpenCL™ Support Package(OSP) part of the FPGA Interface Manager
  - Kernel Avalon interface translated to CCI-P by the OSP



To learn more about using OpenCL with FPGAs, visit [Intel FPGA Customer Training page](#)

# OpenCL™ Adds HW and SW Abstraction

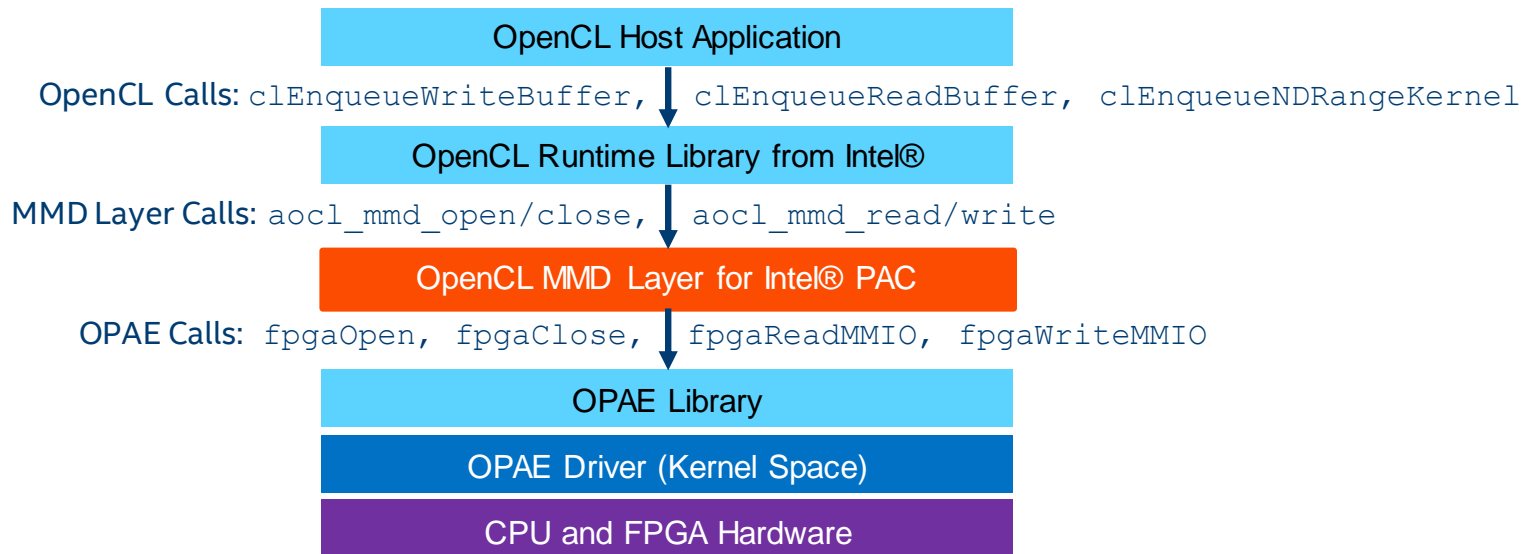


# Software Stack and the OpenCL™ MMD Layer

Memory-Mapped Device (MMD) SW connects board driver to OpenCL runtime

- For the Acceleration Stack it links the OpenCL API to OPAE

Uses DMA BBB software driver unmodified and runs it in a separate thread.



# OpenCL™ with Acceleration Stack Features

Standard OpenCL FPGA application can be used as is for the Acceleration Stack

- No code changes necessary, just switch the OSP in the compilation environment
  - Acceleration Stack functionality built into the OpenCL Support Package

Easily leverage the capabilities of the Acceleration Stack

- Allow OpenCL to be used with virtualization
- No need develop RTL conforming to CCI-P interface
- No need to code at the lower OPAE level

Switch between OpenCL and Non-OpenCL Acceleration Stack applications without rebooting

- FPGA Interface Manager for both are the same

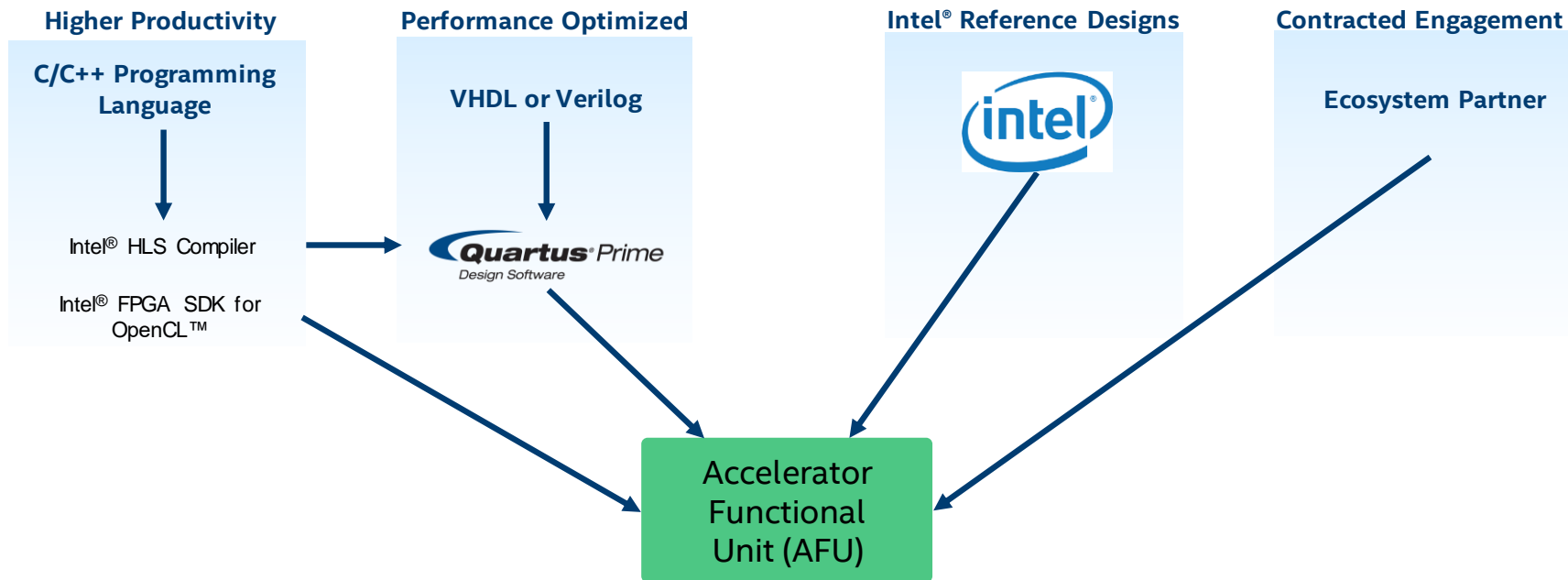
The background is a deep blue gradient. On the left side, there is a stylized, glowing blue globe. The globe is covered with various icons, including a cloud with three dots, a person silhouette, and a circular arrow. A bright, horizontal light streak cuts across the globe. Radiating from the globe are numerous thin, white lines that extend across the blue background, creating a sense of motion and data flow.

# ACCELERATOR FUNCTIONAL UNITS (AFU)

# How Can FPGA Accelerators Be Created?

## Self-Developed

## Externally-Sourced



# Growing List of Accelerator Solution Partners

## Easing Development and Data Center Deployment of Intel FPGAs For Workload Optimization

Accelize

ALGO-LOGIC

ADAPTIVE MICRO-WARE

bigstream  
ACCELERATING INTELLIGENCE

b com

CAST

CTACCEL  
Accelerated Computing

Falcon  
COMPUTING

iAbra

Levyx

MEGH  
COMPUTING

Myrtle

NAGASE  
NAGASE & CO., LTD.

napatech

rENIAC

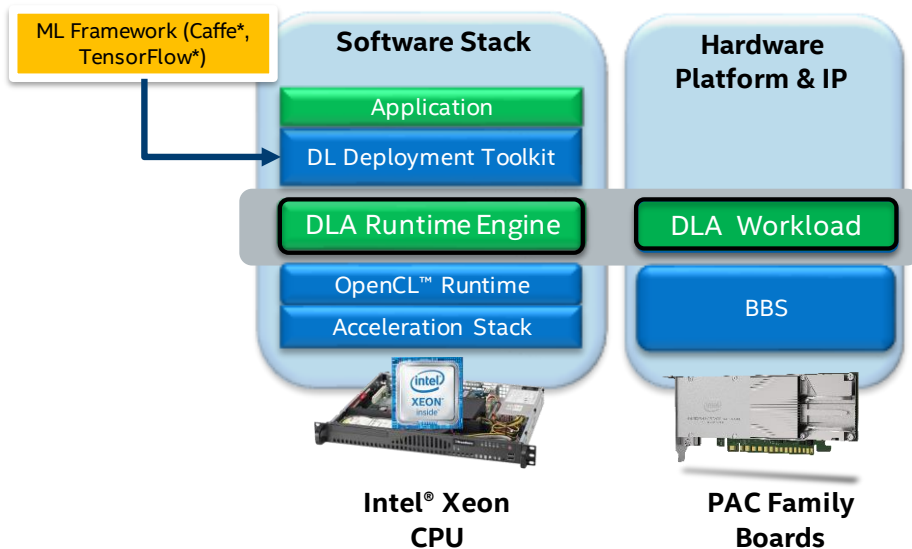
SOC  
System-On-Chip Technologies

enyx

swarm64

# Machine Learning on Intel® FPGA Platform

## Acceleration Stack Platform Solution

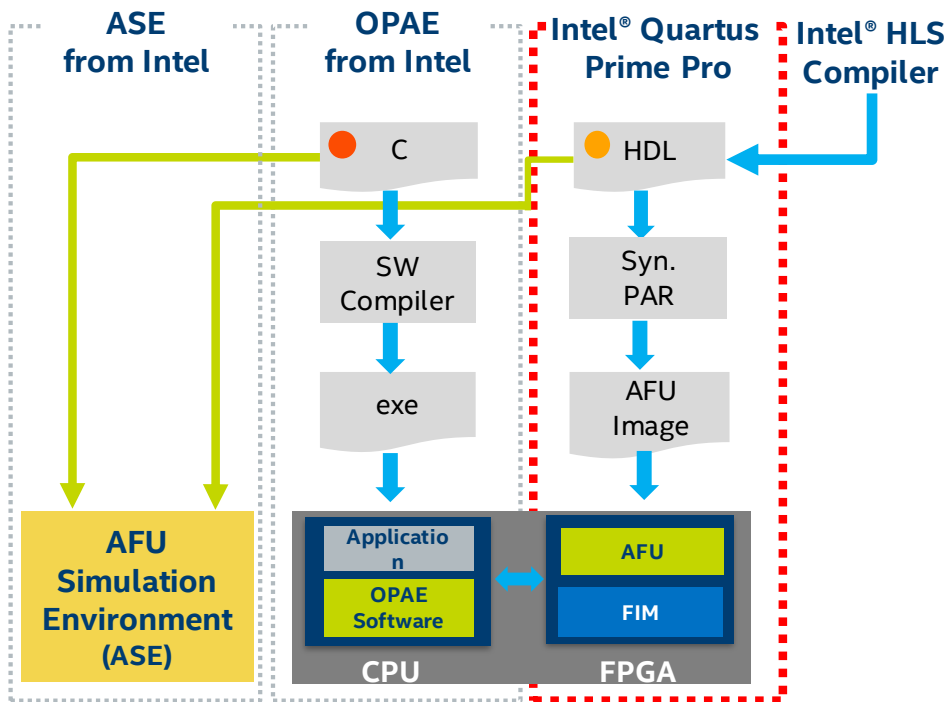


Ships with bitstreams for FPGA as part of OpenVINO™ toolkit

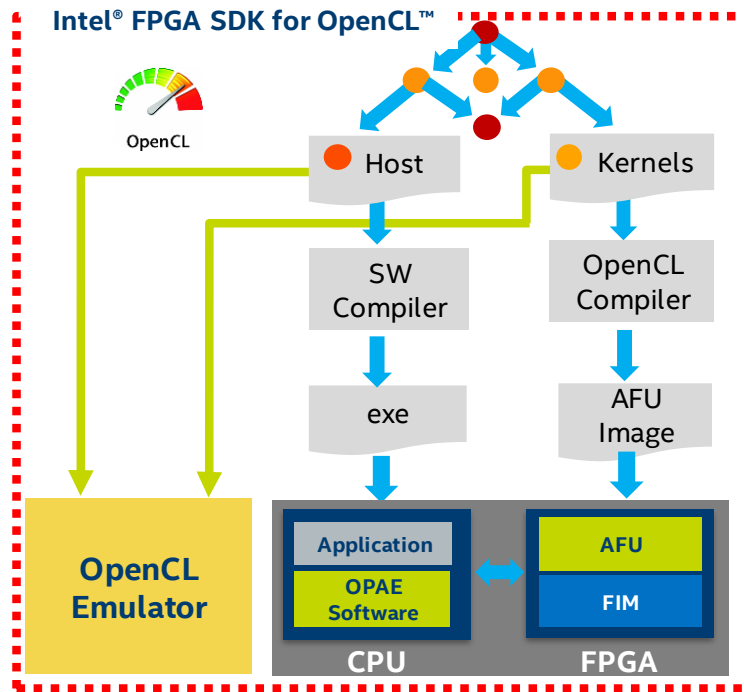


# Accelerator Functional Unit Development

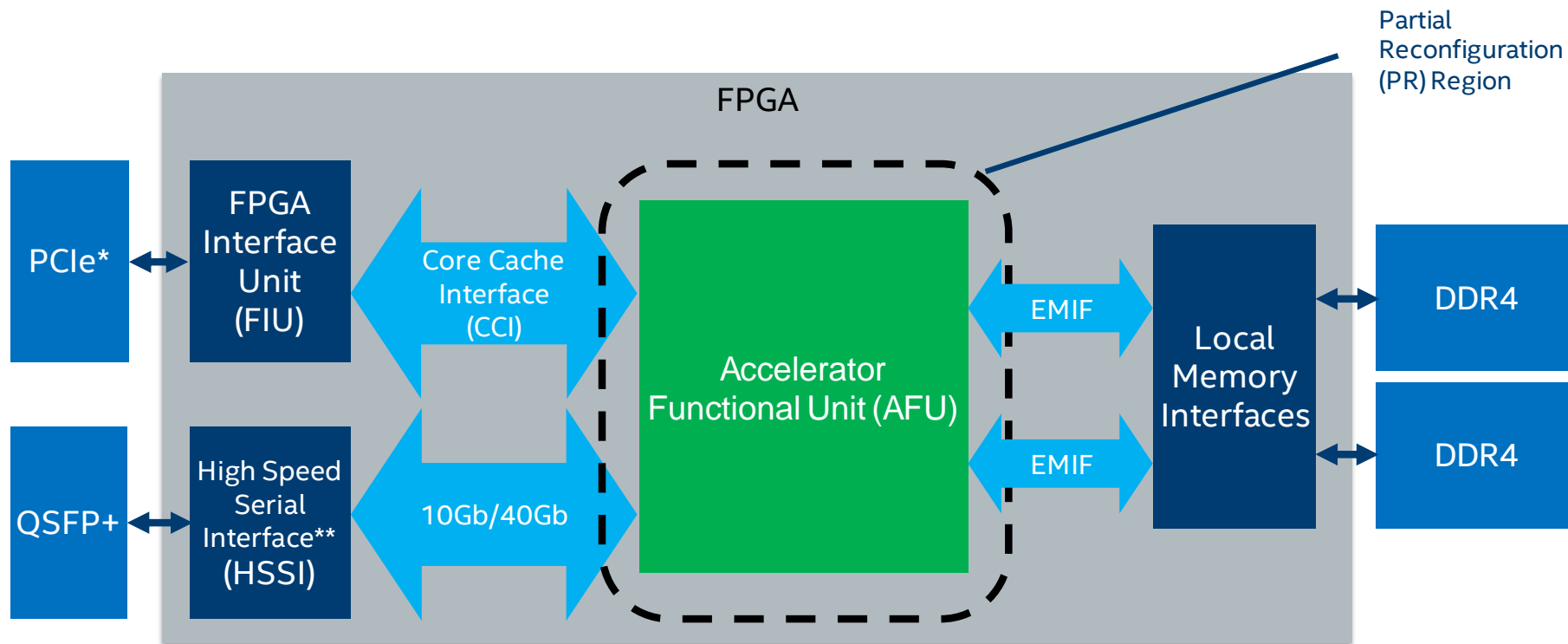
## HDL Programming



## OpenCL Programming

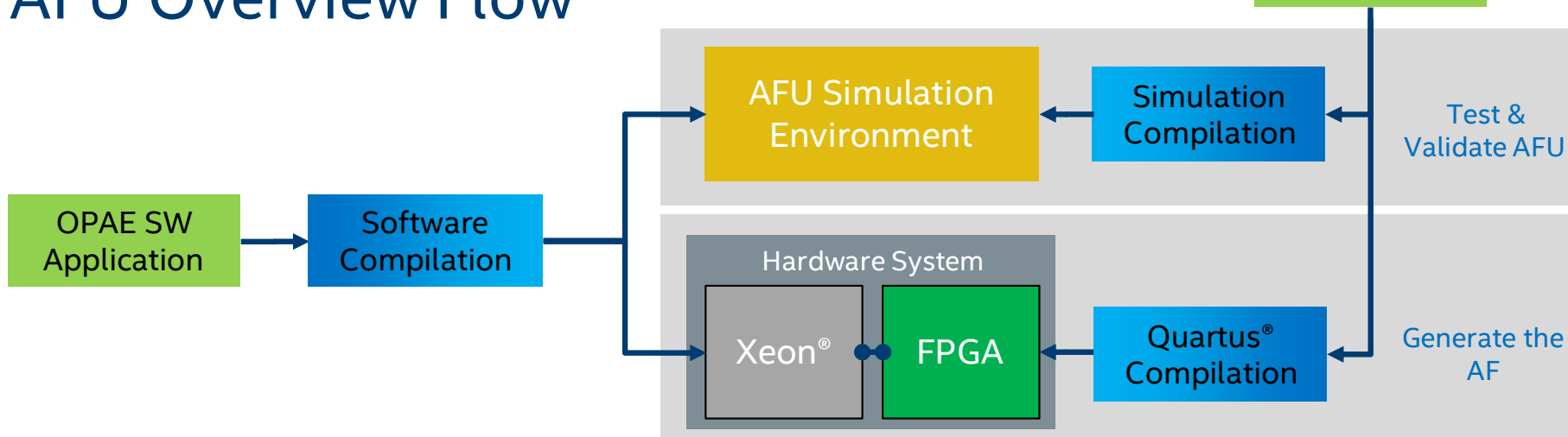


# FPGA Components (Acceleration Stack v1.1)



- Could be other interfaces in the future (e.g. UPI)
- \*\* Available in v1.1 of Acceleration Stack

# AFU Overview Flow



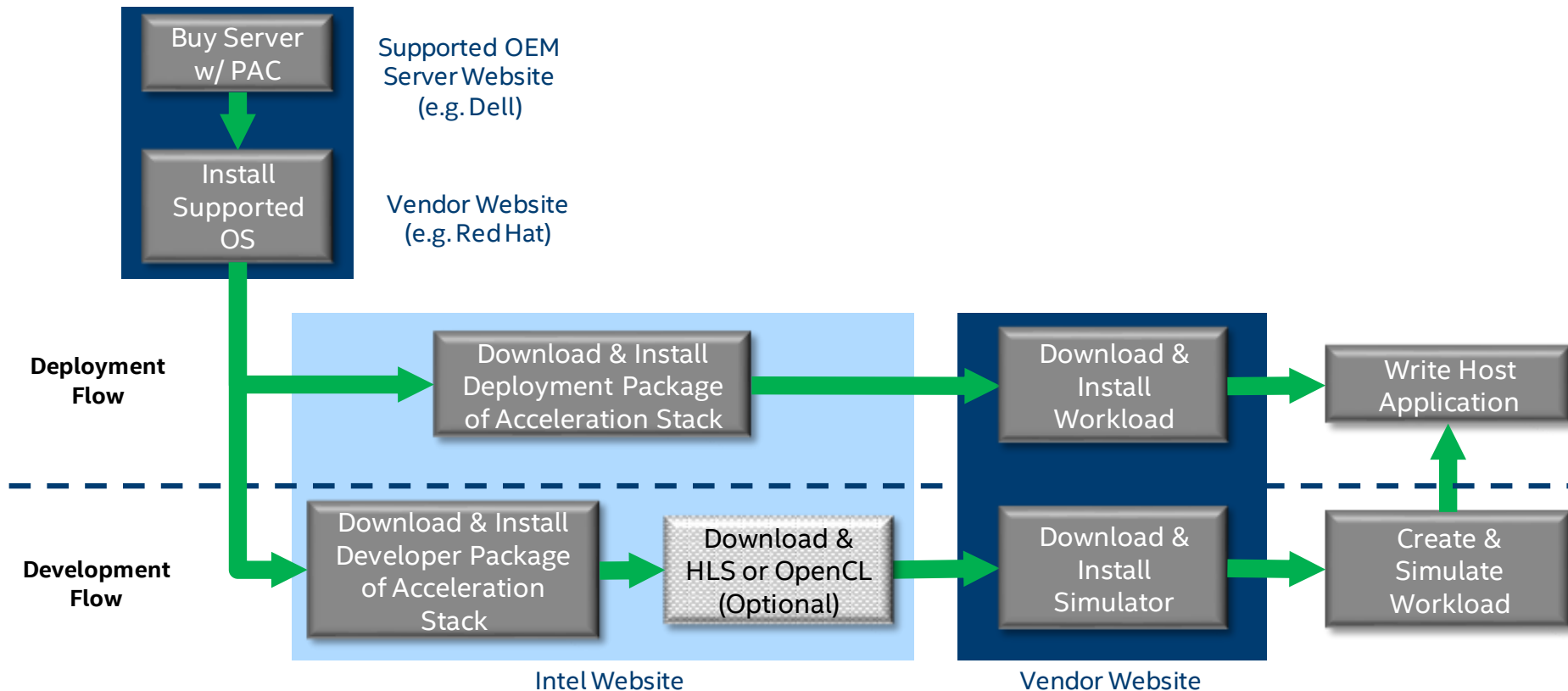
AF Simulation Environment (ASE) enables seamless portability to real HW

- Allows fast verification of OPAЕ software together with AF RTL without HW
  - SW Application loads ASE library and connects to RTL simulation
- For execution on HW, application loads Runtime library and RTL is compiled by Intel® Quartus into FPGA bitstream

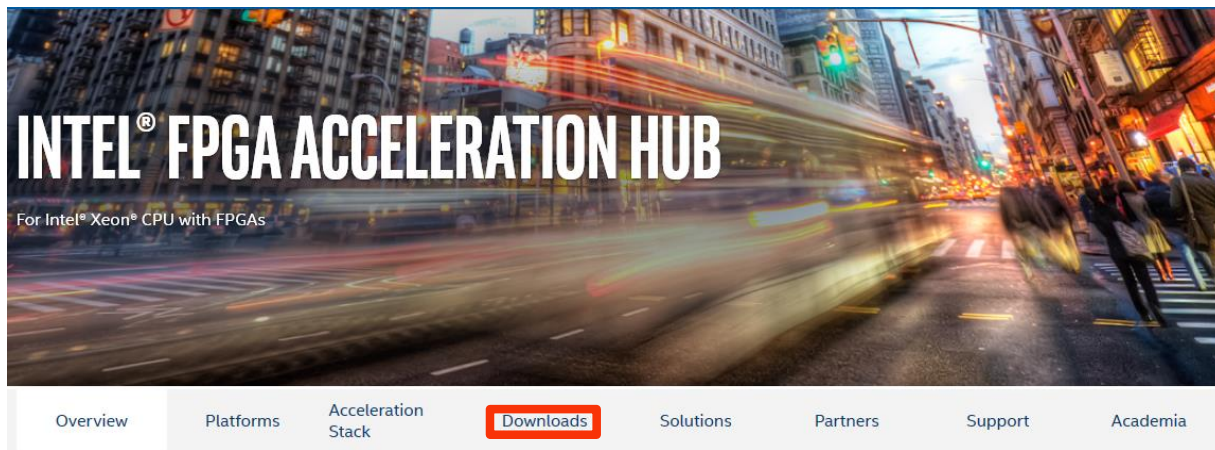
The background is a deep blue gradient. On the left side, there is a stylized, glowing globe or sphere. The globe is composed of a network of white lines and dots, suggesting a digital or data-driven theme. Several circular icons are visible on the globe's surface, including one with three people silhouettes and another with a single person silhouette. A bright, horizontal light streak cuts across the middle of the globe. Radiating from the globe are numerous thin, white lines that extend across the blue background, creating a sense of motion and connectivity.

# GETTING STARTED

# Out-of-Box Flow for Acceleration Stack



# Download Acceleration Stack on Acceleration Hub

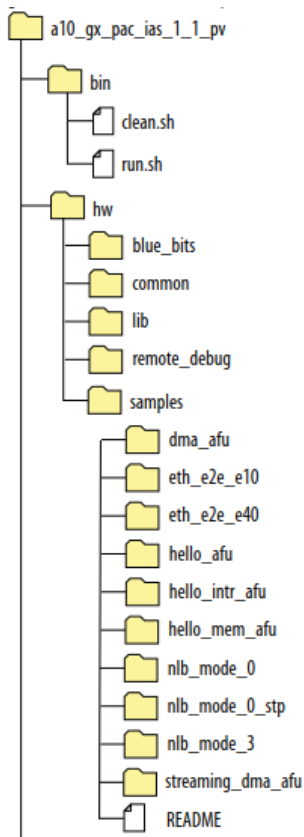


In today's world, the number of connected devices and the amount of data continues to increase every day. The rate at which data arrives from these devices into data centers also continues to increase. By leveraging Intel® FPGAs as accelerators, a wide range of workloads can be enhanced to accommodate this increased data, and new demands for analyzing it.

FPGAs are silicon devices that can be dynamically reprogrammed with a datapath that exactly matches your workloads, such as data analytics, image inference, encryption, and compression. This versatility enables the provisioning of a faster processing, more power efficient, and lower latency service – lowering your total cost of ownership, and maximizing compute capacity within the power, space, and cooling constraints of your data centers.

Traditionally, FPGAs require deep domain expertise to program for, but the Intel® Acceleration Stack for Intel Xeon® CPU with FPGAs simplifies the development flow and enables rapid deployment across the data center. Intel is partnering with FPGA intellectual property (IP) developers, server original equipment manufacturers (OEMs), virtualization platform providers, operating system (OS) vendors, and system integrators to enable customers to efficiently develop and operationalize their infrastructure.

# Example Designs to Get Started



Example	Description
Hello AFU	Simple AFU with direct CCI connection for MMIO access
Hello Intr AFU	Example use of user interrupts
Hello Mem AFU	Example showing using USR Clock to auto close timing in the AFU
DMA AFU	Example DMA AFU to move data between host memory and local FPGA memory. Uses BBB and bridges Avalon to CCI
Streaming DMA AFU	Example DMA AFU to move data between host memory and the AFU directly as a streaming packet
Eth e2e e10	10Gb Ethernet loopback design
Eth e2e e40	40Gb Ethernet loopback design
NLB mode 0	Native LoopBack adaptor (rd/wr) with more features
NLB mode 0 stp	Native LoopBack adaptor with SignalTap remote debug
NLB mode 3	Native LoopBack adaptor (rd/wr)

# Recommended Getting Started Process

Install VM with Linux and install Acceleration Stack Development Flow

Run through quick start guide to validate environment

- Run ASE on Hello AFU
- Regenerate hello\_afu .gbs file and host software application

If applicable, run the OpenCL example

Examine the DMA\_AFU and/or Streaming\_DMA\_AFU example to understand how to move data using DMA with AFU

Create your own by modifying one of the example designs (e.g. DLA through OpenVINO™ toolkit)



# Follow-On Training:

## Online Training Course

- Introduction to the Acceleration Stack for Intel® Xeon w/ FPGA
- OpenCL™ Development with the Acceleration Stack
- RTL development and acceleration with the Acceleration Stack
- Application Development on the Acceleration Stack for Intel® Xeon® CPU with FPGAs
- Introduction to High-Level Synthesis (7 courses)
- Introduction to Parallel Computing w/ OpenCL on FPGAs
- Deploying Intel FPGAs for Inferencing with OpenVINO Toolkit
- Programmers' Introduction to the Intel® FPGA Deep Learning Acceleration Suite

## Instructor Led Training Courses

- Introduction to High-Level Synthesis with Intel® FPGAs
- High-Level Synthesis Advanced Optimization Techniques
- Introduction to OpenCL
- Optimizing OpenCL™ for Intel® FPGAs (16 Hours Course)

<https://www.altera.com/support/training/overview.html>

# Summary

Acceleration Stack for Intel Xeon CPU with FPGAs enables seamless integration of the FPGA accelerators into clusters

DLA Suite comes with OpenVINO as bitstreams for the PAC cards that runs on the acceleration stack

The OPAE software layer makes adoption of the FPGA into the host application code seamless and

The acceleration stack makes building accelerator functions much faster and board agnostic through the use of a standardized FIM

# Legal Disclaimers/Acknowledgements

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration. Check with your system manufacturer or retailer or learn more at [www.intel.com](http://www.intel.com).

Intel, the Intel logo, Intel Inside, the Intel Inside logo, MAX, Stratix, Cyclone, Arria, Quartus, HyperFlex, Intel Atom, Intel Xeon and Enpirion are trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries.

OpenCL is the trademark of Apple Inc. used by permission by Khronos

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

© Intel Corporation

