# Xeon+FPGA Platform for the Data Center

ISCA/CARL 2015

PK Gupta, Director of Cloud Platform Technology, DCG/CPG
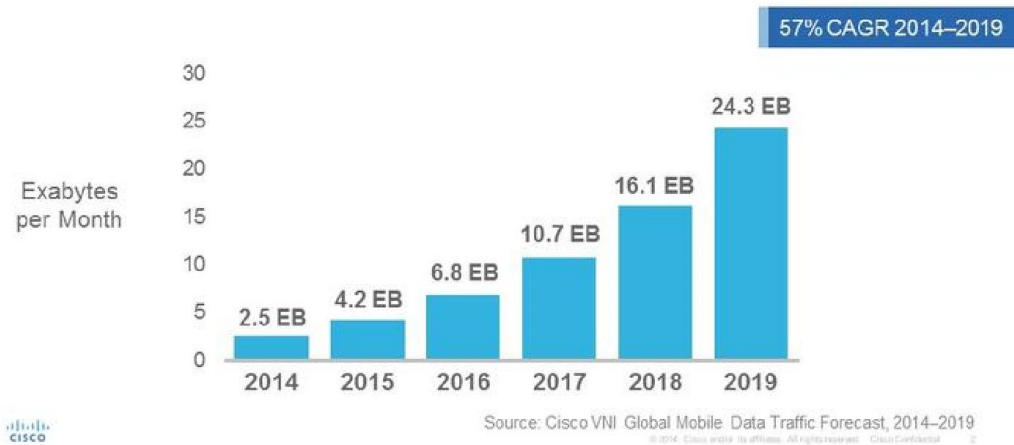
# Overview

- **Data Center and Workloads**

- Xeon+FPGA Accelerator Platform

- Applications and Eco-system
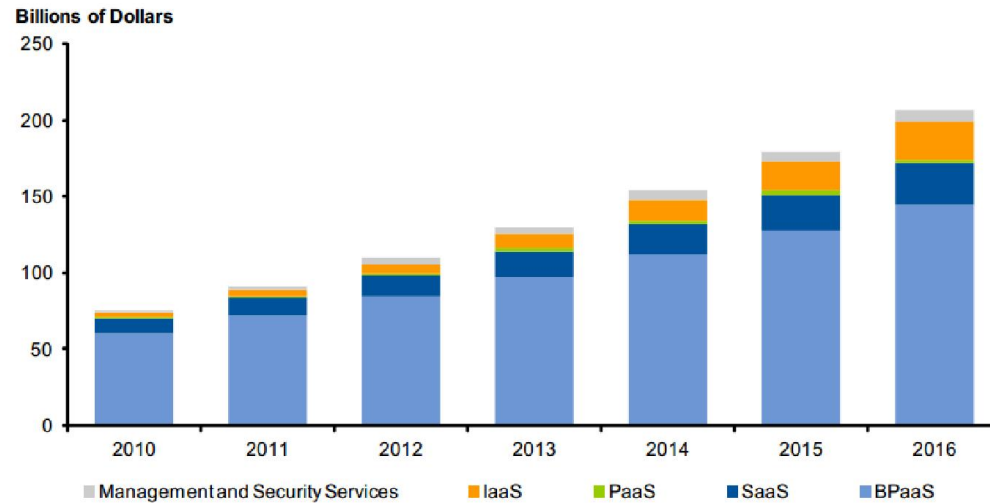
# Exponential growth in mobile….

## Global Mobile Data Traffic Growth / Top-Line
Global Mobile Data Traffic will Increase 10-Fold from 2014—2019

**57% CAGR 2014–2019**

Exabytes per Month

- 2014: 2.5 EB
- 2015: 4.2 EB
- 2016: 6.8 EB
- 2017: 10.7 EB
- 2018: 16.1 EB
- 2019: 24.3 EB

Source: Cisco VNI Global Mobile Data Traffic Forecast, 2014–2019

cisco

# …is driving Data Center growth

**Billions of Dollars**



Source: Gartner (August 2012)

# Leading to search for greater performance efficiencies...



Monthly Costs

- 57% Servers
- 8% Networking Equipment
- 18% Power Distribution & Cooling
- 13% Power
- 4% Other Infrastructure

3yr server & 10 yr infrastructure amortization

# ...Across Data Center Workloads

- Diverse workloads:

  - Cloud Services: Search, Web Servers, ..

  - Analytics: Big Data, Machine Learning, ...

  - Scientific: Genomics, Security, ...

  - Communication: Packet Processing, Virtual Switching, ...

  - Storage: Compression, Deduplication, ...

- Changing dynamics:

  - No single killer app

  - Emerging new apps drive changes in workloads

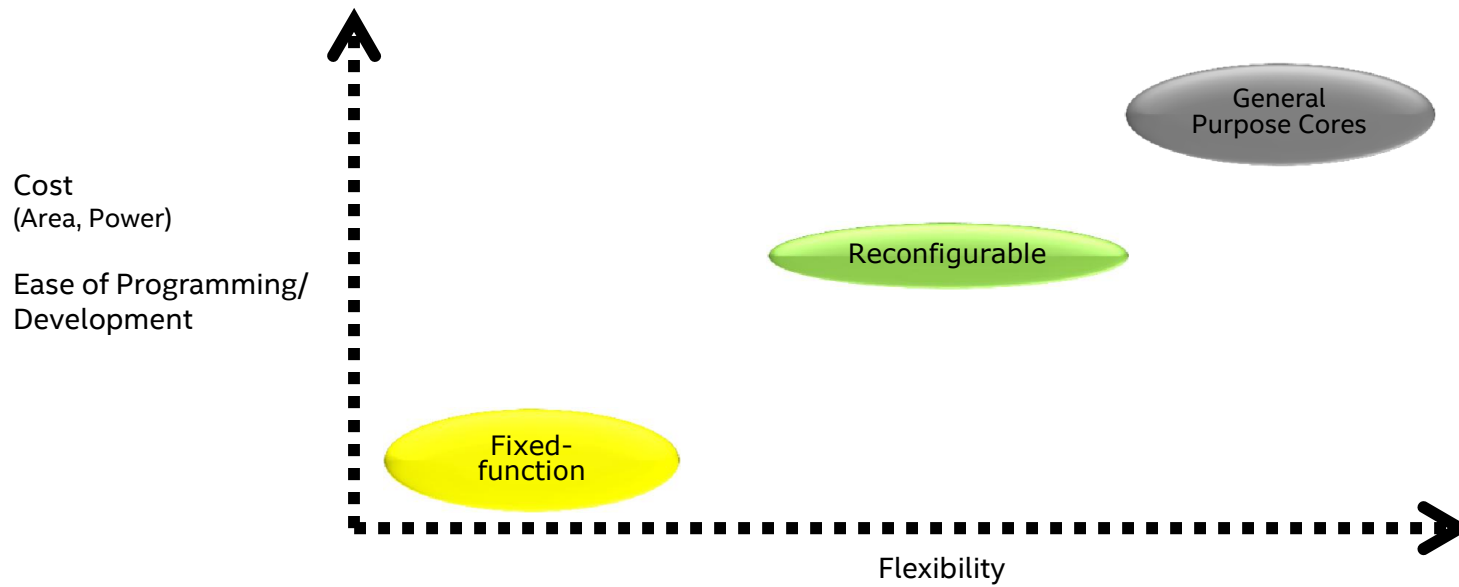# A homogenous compute platform for the Data Center?

# Overview

- Data Center and Workloads

- **Xeon+FPGA Accelerator Platform**

- Applications and Eco-system
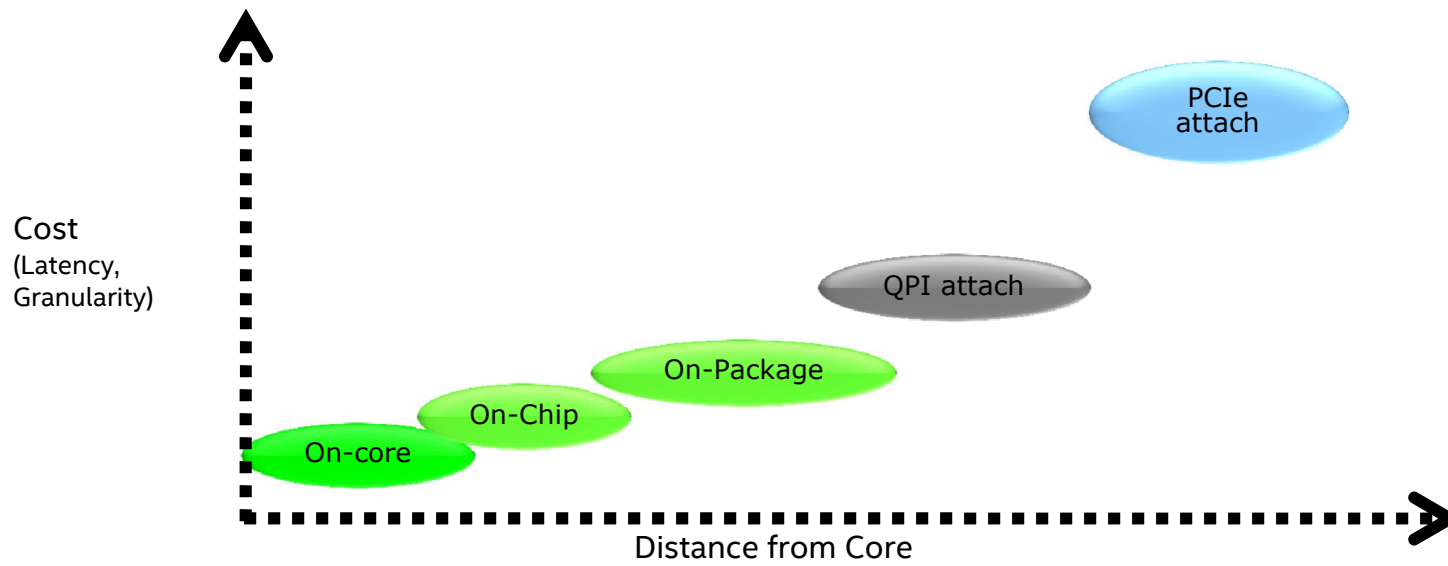
# Motivation for Accelerators

- **Enhanced Performance**: Accelerators compliment CPU cores to meet market needs for performance of diverse workloads in the Data Center:

  - Enhance single thread performance with tightly coupled accelerators or compliment multi-core performance with loosely coupled accelerators via PCIe or QPI attach

- **Move to Heterogeneous Computing**: Moore's Law continues but demands radical changes in architecture and software.

  - Architectures will go beyond homogeneous parallelism, embrace heterogeneity, and exploit the bounty of transistors to incorporate application-customized hardware.

# Accelerator Architecture

Cost
(Area, Power)

Ease of Programming/
Development

Fixed-
function

Reconfigurable

General
Purpose Cores

Flexibility

Performance Efficiency: Performance/Watt, Performance/$
Programming Complexity : Effort, Cost

# Accelerator Attach



Cost
(Latency,
Granularity)

PCIe attach

QPI attach

On-Package

On-Chip

On-core

Distance from Core

Best attach technology might be application or even algorithm dependent

# Coherency and Programming Model
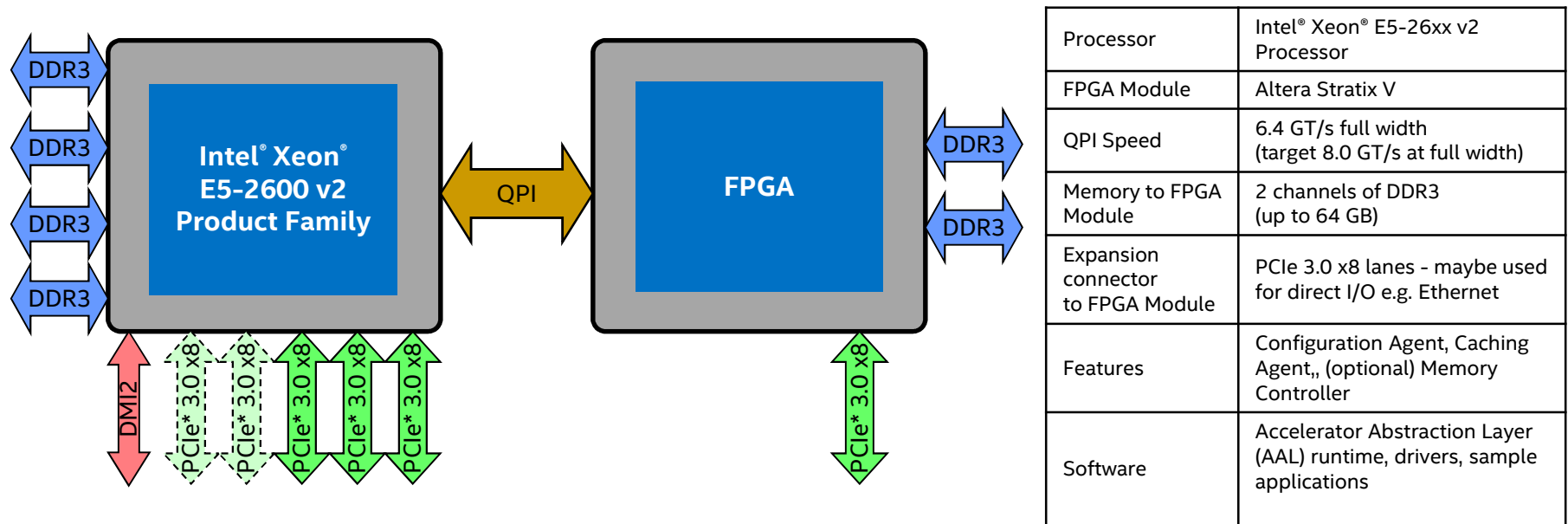
- Data Movement

  - In-line
    - Accelerator processes data fully or partially from direct I/O

  - Shared Virtual Memory :
    - Virtual addressing eliminates need for pinning memory buffers
    - Zero-copy data buffers

- Interaction between Core and Accelerator

  - Off-load

  - Hybrid : algorithm implemented on host and accelerator

(intel)

# Proposed Platform for the Data Center

- FPGA with coherent low-latency interconnect:

  - Simplified programming model

    - Support for virtual addressing

    - Data Caching

  - Enables new classes of algorithms for acceleration with:

    - Full access to system memory

    - Support for efficient irregular data pattern access

  - Remapping of algorithms from off-load model to hybrid processing model
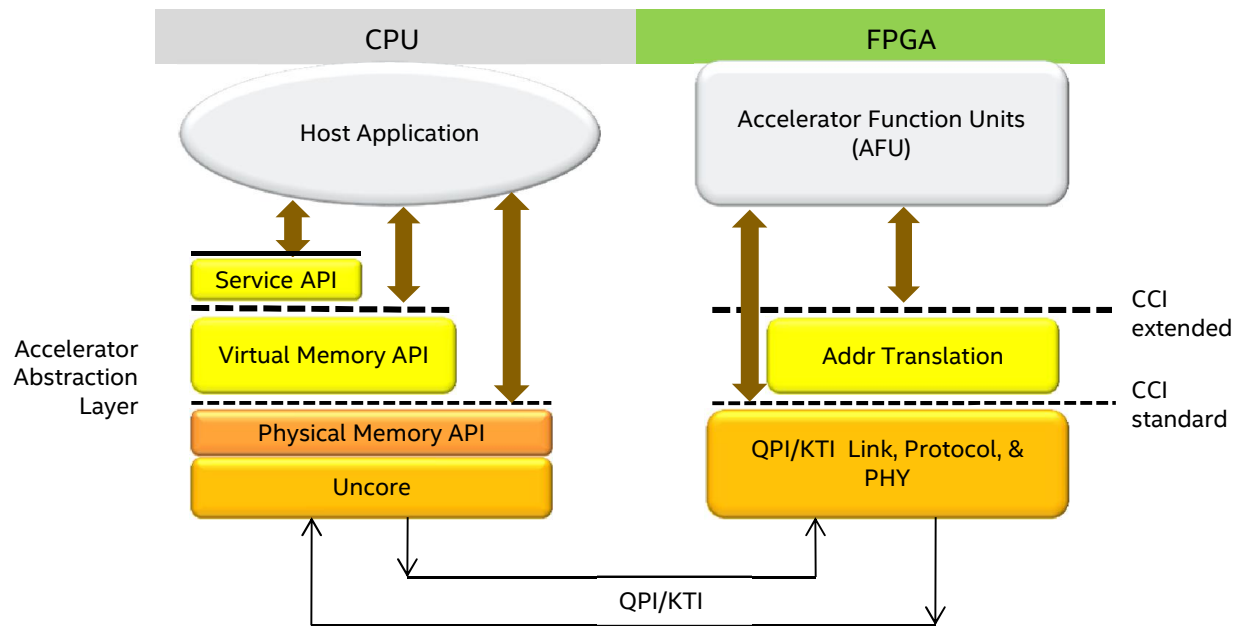
    - Fine grained interactions

# IVB+FPGA Software Development Platform

**Software Development for Accelerating Workloads using Xeon and coherently attached FPGA in-socket**

| | |
|---|---|
| Processor | Intel® Xeon® E5-26xx v2 Processor |
| FPGA Module | Altera Stratix V |
| QPI Speed | 6.4 GT/s full width (target 8.0 GT/s at full width) |
| Memory to FPGA Module | 2 channels of DDR3 (up to 64 GB) |
| Expansion connector to FPGA Module | PCIe 3.0 x8 lanes - maybe used for direct I/O e.g. Ethernet |
| Features | Configuration Agent, Caching Agent,, (optional) Memory Controller |
| Software | Accelerator Abstraction Layer (AAL) runtime, drivers, sample applications |

**Intel® Xeon® E5-2600 v2 Product Family** — QPI — **FPGA**

DDR3, DDR3, DDR3, DDR3

DMI2, PCIe* 3.0 x8, PCIe* 3.0 x8, PCIe* 3.0 x8, PCIe* 3.0 x8, PCIe* 3.0 x8

DDR3, DDR3

PCIe* 3.0 x8

**Heterogeneous architecture with homogenous platform support**

# Programming Interfaces



| CPU | FPGA |
|---|---|

**Host Application**

**Accelerator Function Units (AFU)**

Service API

Accelerator Abstraction Layer

Virtual Memory API

Addr Translation

CCI extended

Physical Memory API

CCI standard

Uncore

QPI/KTI Link, Protocol, & PHY

QPI/KTI

**Programming interfaces will be forward compatible from SDP to future MCP solutions**
**Simulation Environment available for development of SW and RTL**

# Programming Interfaces : OpenCL



**CPU** | **FPGA**

OpenCL Host Code — OpenCL Application — OpenCL Kernel Code

OpenCL RunTime → CFG → OpenCL Kernels

Accelerator Abstraction Layer

Service API

Virtual Memory API

Physical Memory API

VirtMem

Physical Memory API

CCI Extended

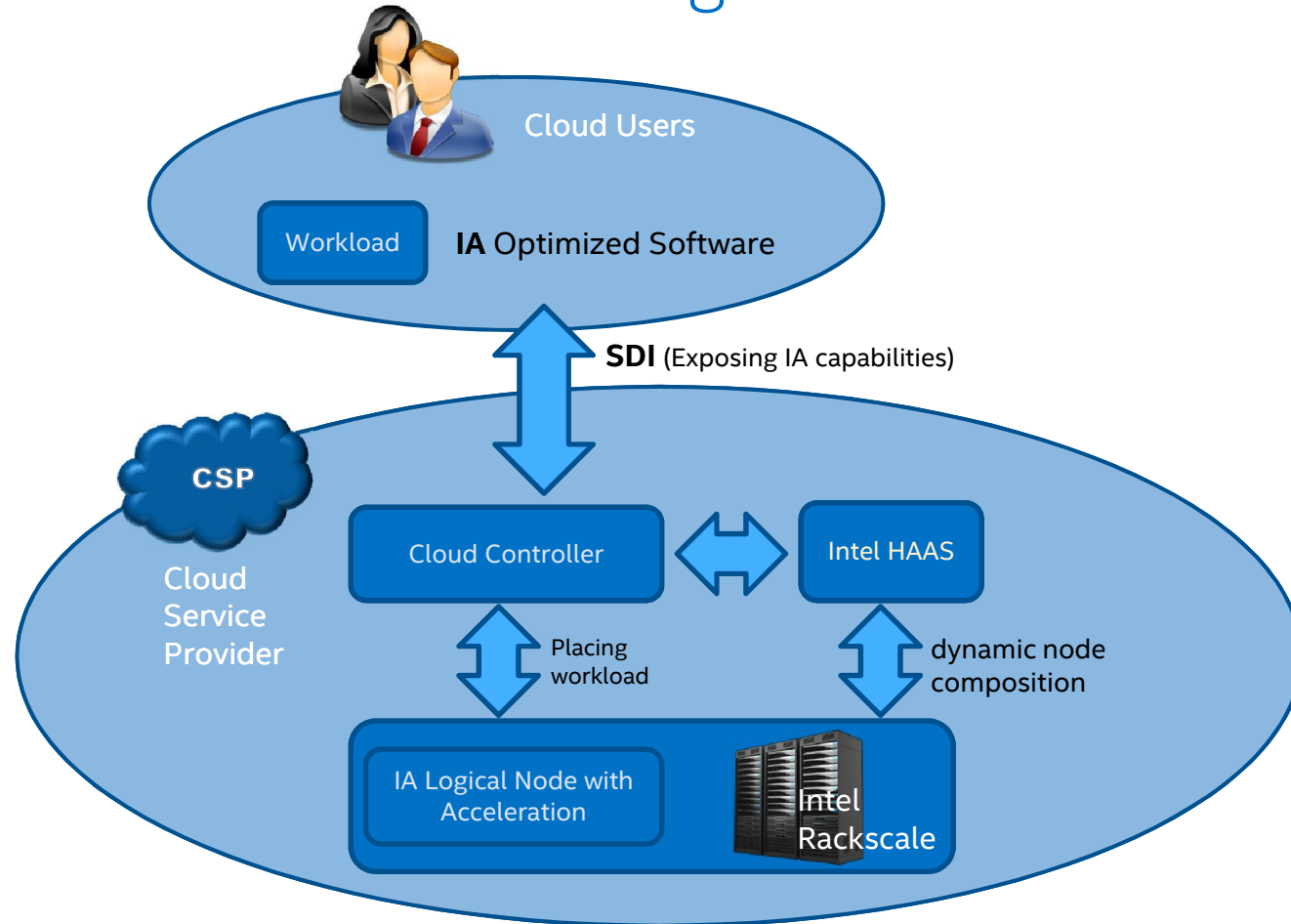CCI Standard

QPI/UPI/P Ci

System Memory

Unified application code abstracted from the hardware environment
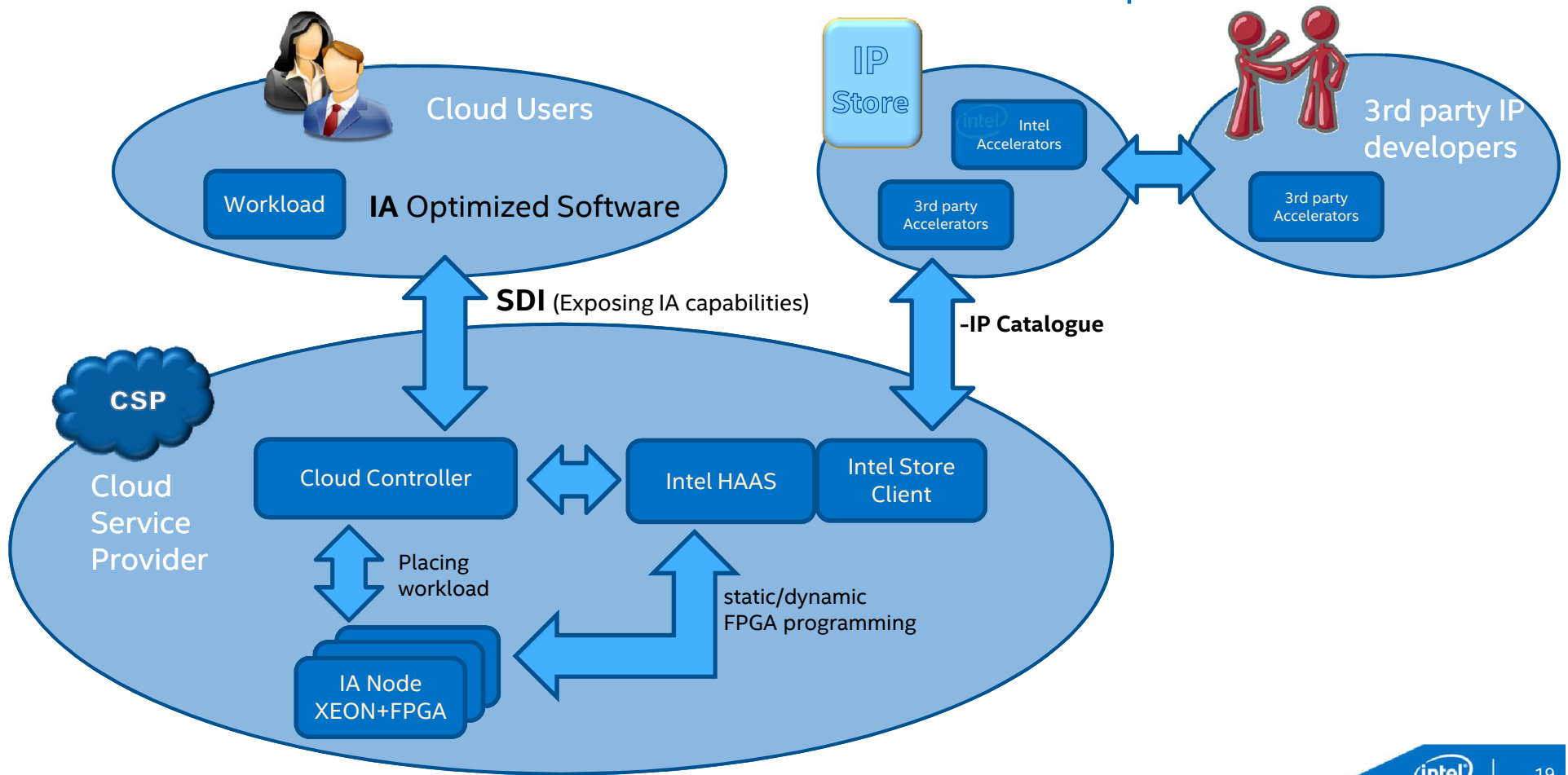Portable across generations and families of CPUs and FPGAs

# Overview

- Data Center and Workloads

- Xeon+FPGA Accelerator Platform
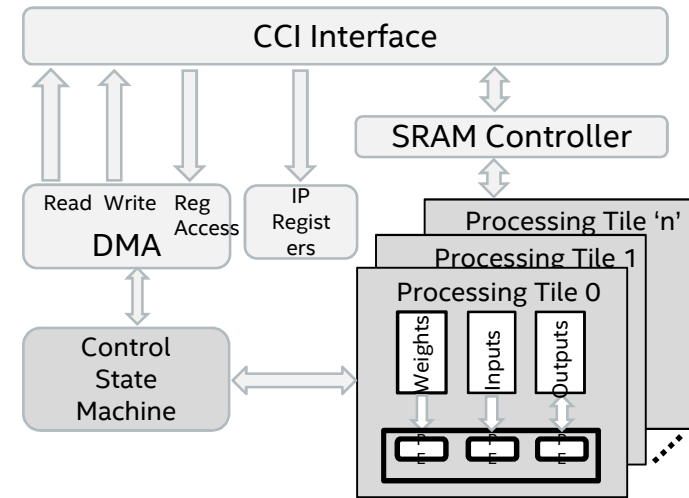
- **Applications and Eco-system**

# XEON+FPGA in the Cloud : integration with SDI and RSA
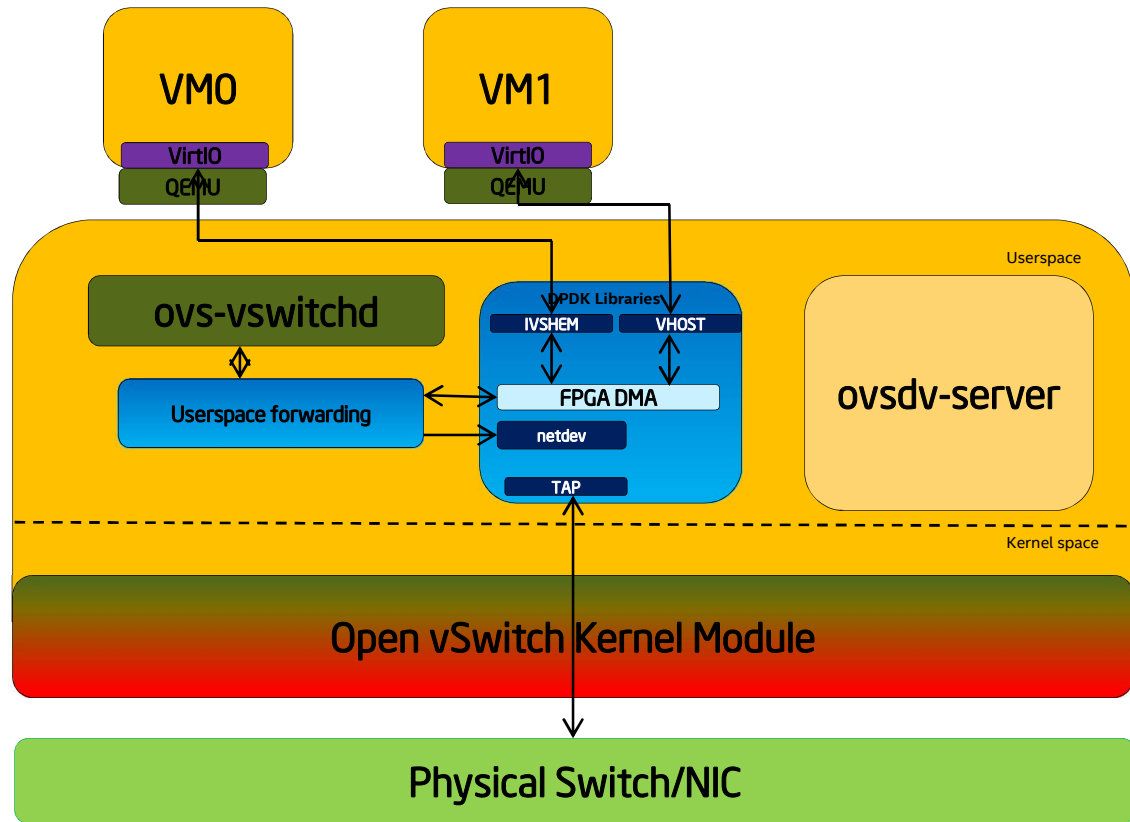
Cloud Users

Workload  **IA** Optimized Software

**SDI** (Exposing IA capabilities)

**CSP**

Cloud Service Provider

Cloud Controller ⟷ Intel HAAS

Placing workload

dynamic node composition

IA Logical Node with Acceleration

Intel Rackscale

# XEON+FPGA in the Cloud: IP Store Concept

**Cloud Users**

Workload  **IA** Optimized Software

**IP Store**

Intel Accelerators

3rd party Accelerators

**3rd party IP developers**

3rd party Accelerators

**SDI** (Exposing IA capabilities)

–IP Catalogue

**CSP**

Cloud Service Provider

Cloud Controller

Intel HAAS

Intel Store Client

Placing workload

static/dynamic FPGA programming

IA Node XEON+FPGA

# Example Usage : Deep Learning Framework for Visual Understanding

# Example Usage: Accelerating Open VSwitch w/DPDK

**VM0**
VirtIO
QEMU

**VM1**
VirtIO
QEMU

Userspace

ovs-vswitchd

DPDK Libraries
IVSHEM    VHOST

Userspace forwarding    FPGA DMA

netdev

ovsdv-server

TAP

Kernel space

**Open vSwitch Kernel Module**

**Physical Switch/NIC**

- Offload DMA Engine to FPGA :
  - Frees up CPU cycles to perform more useful work
  - Reduce cache pollution.
- Add support for Packet Classification, ACL, and other functions including Direct I/O in FPGA

# Example Usage: High Frequency Trading Accelerator



CPU

Host Application

FPGA

Feed Parser

Ethernet PHY & MAC

Trading Logic / Statistics

Order Generation

# Academic Research

**Call for Proposals: Intel-Altera Heterogeneous Architecture Research Platform Program**
Submitted by Nicholas Carter

**Intel-Altera Heterogeneous Architecture Research Platform (HARP) Program**

Intel® Corporation and Altera® Corporation are pleased to announce the Heterogeneous Architecture Research Platform (HARP) program, which will provide faculty with computer systems containing Intel microprocessors and an Altera Stratix® V FPGA module that incorporates Intel® QuickAssist Technology in order to spur research in programming tools, operating systems, and innovative applications for accelerator-based computing systems.

Q & A