## Unleashing the full performance of the All Programmable FPGA while abstracting the hardware details
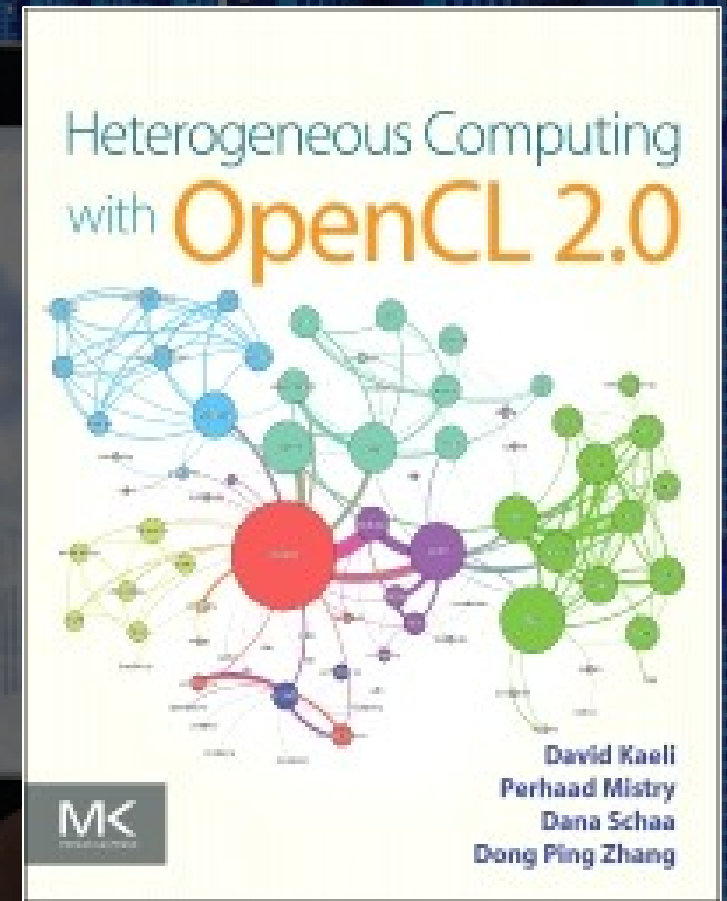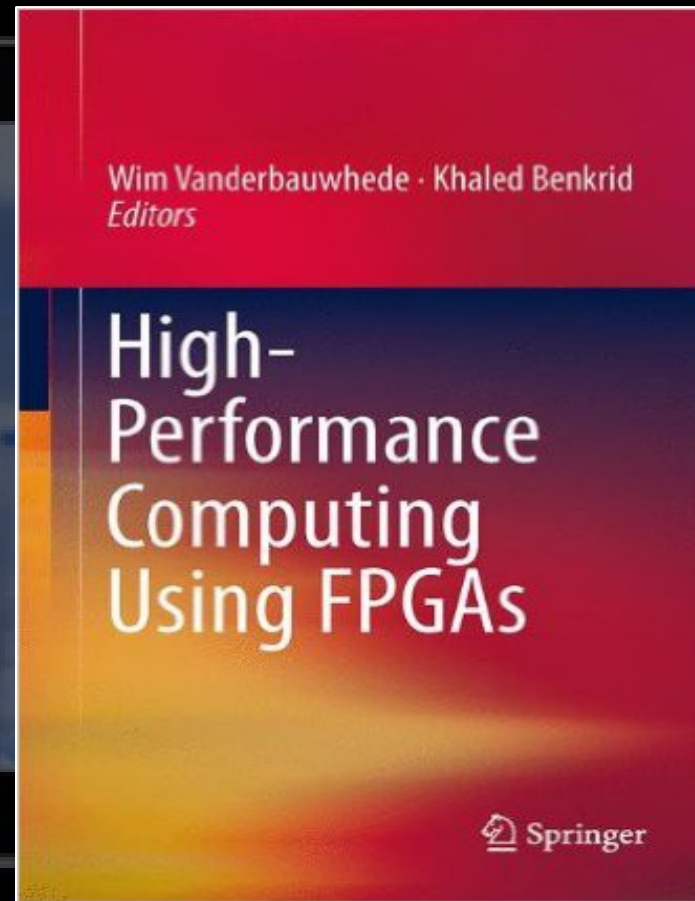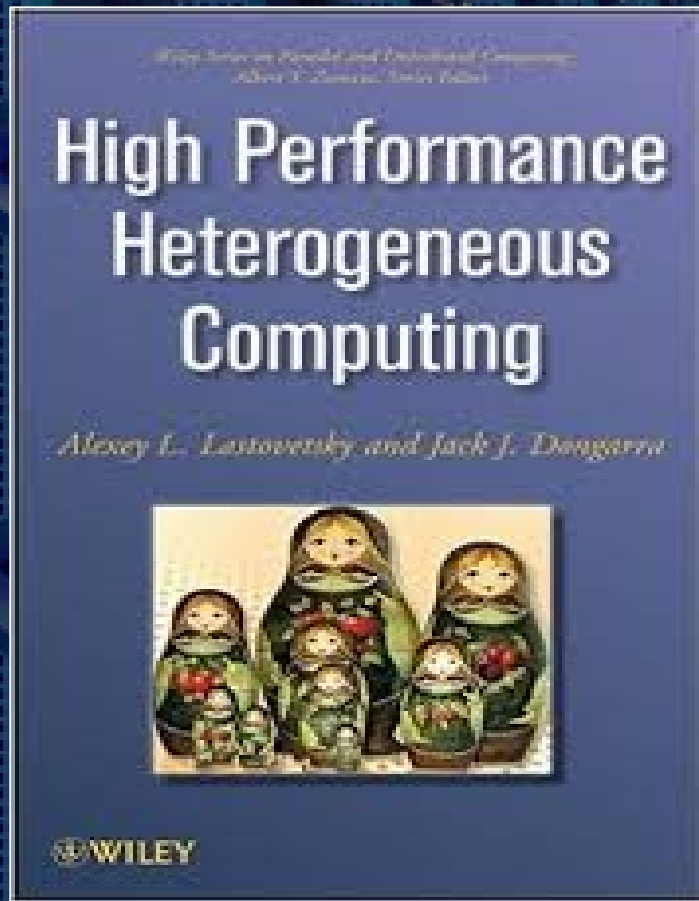
*Ivo Bolsens*

# Traditional Compute Architectures Are Not Scalable



40 Years of Microprocessor Trend Data

Transistors (thousands)

Single-Thread Performance (SpecINT x $10^3$)

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

Requires Heterogeneous Architectures & Acceleration

**ΣXILINX** ➤ ALL PROGRAMMABLE.

# Data Center enabling the Cloud



- HPC
- Deep Learning Training
- Deep Learning Inference
- Image & Video Acceleration
- Data Analytics
- Genomics

**Computing**

- Network Acceleration
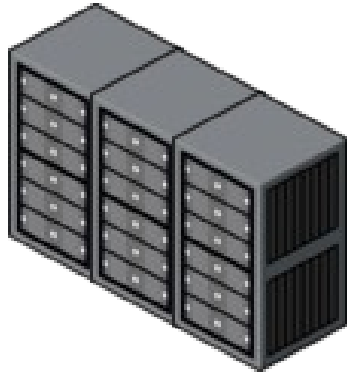- NFV NICs
- SDN Controllers
- Compression

**Networking**

- Flash Arrays Acceleration
- SSDs Acceleration
- NVDIMM
- NVMe over Fabric

**Storage**

# Data Center enabling the Cloud

**GPU**

- HPC
- Deep Learning Training
- Deep Learning Inference
- Image & Video Acceleration
- Data Analytics
- Genomics

**Computing**

**NIC**

- Network Acceleration
- NFV NICs
- SDN Controllers
- Compression

**Networking**

**Storage Controller**

- Flash Arrays Acceleration
- SSDs Acceleration
- NVDIMM
- NVMe over Fabric

**Storage**

 XILINX  ALL PROGRAMMABLE.

# Data Center enabling the Cloud

**GPU**

- HPC
- Deep Learning Training
- Deep Learning Inference
- Image & Video Acceleration
- Data Analytics
- Genomics

**Computing**

**NIC**

- Network Acceleration
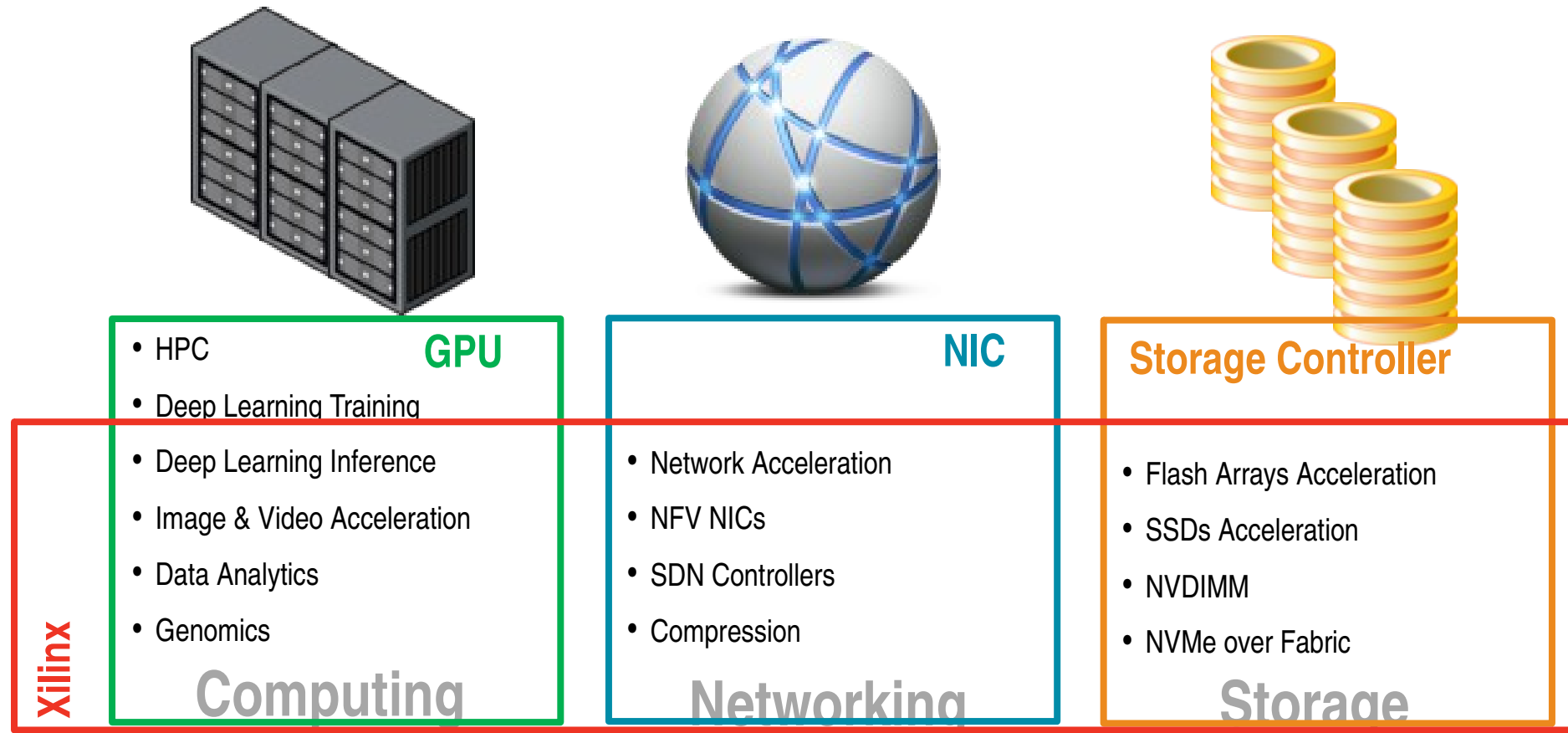- NFV NICs
- SDN Controllers
- Compression

**Networking**

**Storage Controller**

- Flash Arrays Acceleration
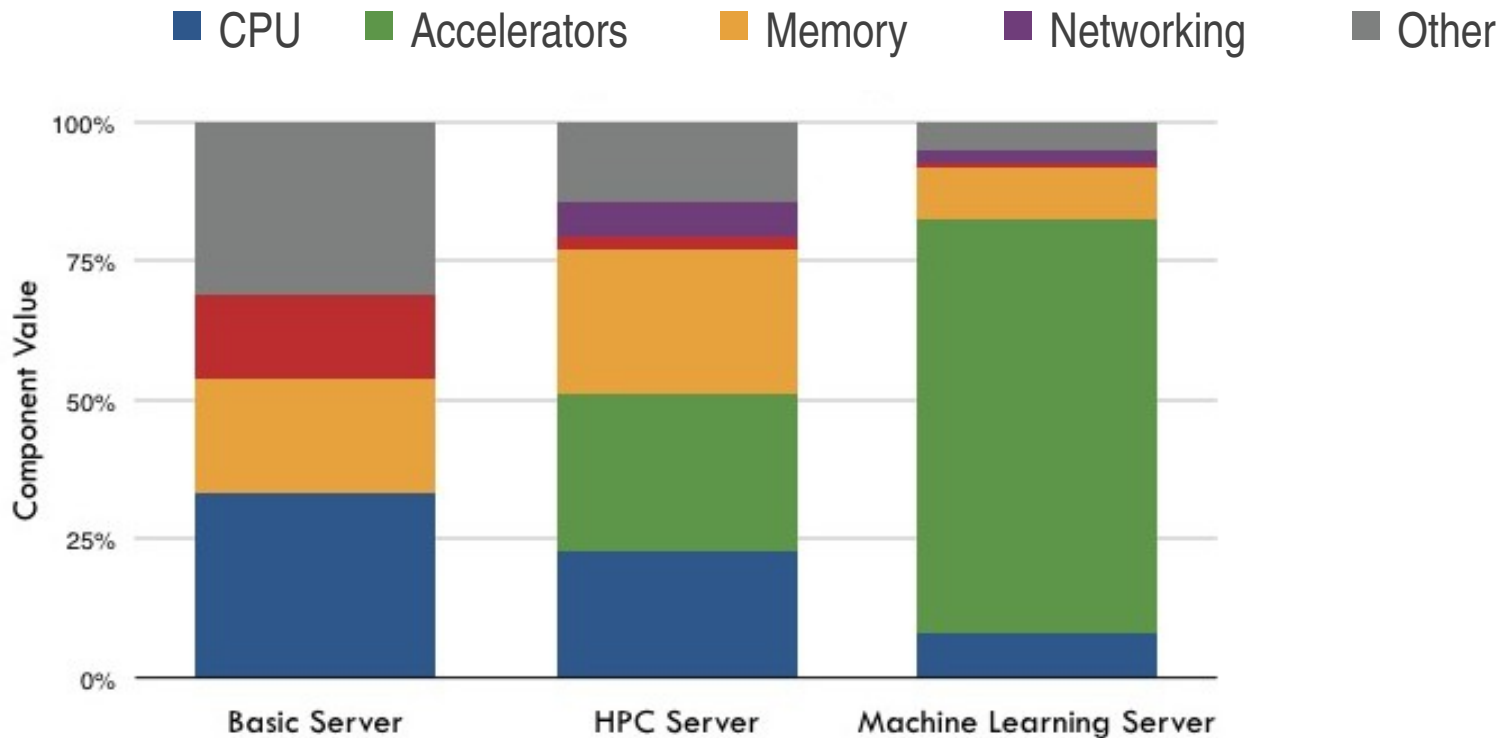- SSDs Acceleration
- NVDIMM
- NVMe over Fabric

**Storage**

**Xilinx**

*Unified HW to address compute, storage, and networking apps*

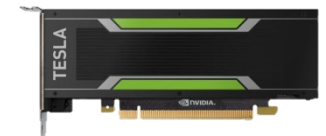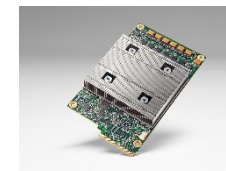**XILINX** ➤ ALL PROGRAMMABLE.

# Heterogeneous Compute Platforms in Datacenter

**How Server Components Change with Machine Learning**
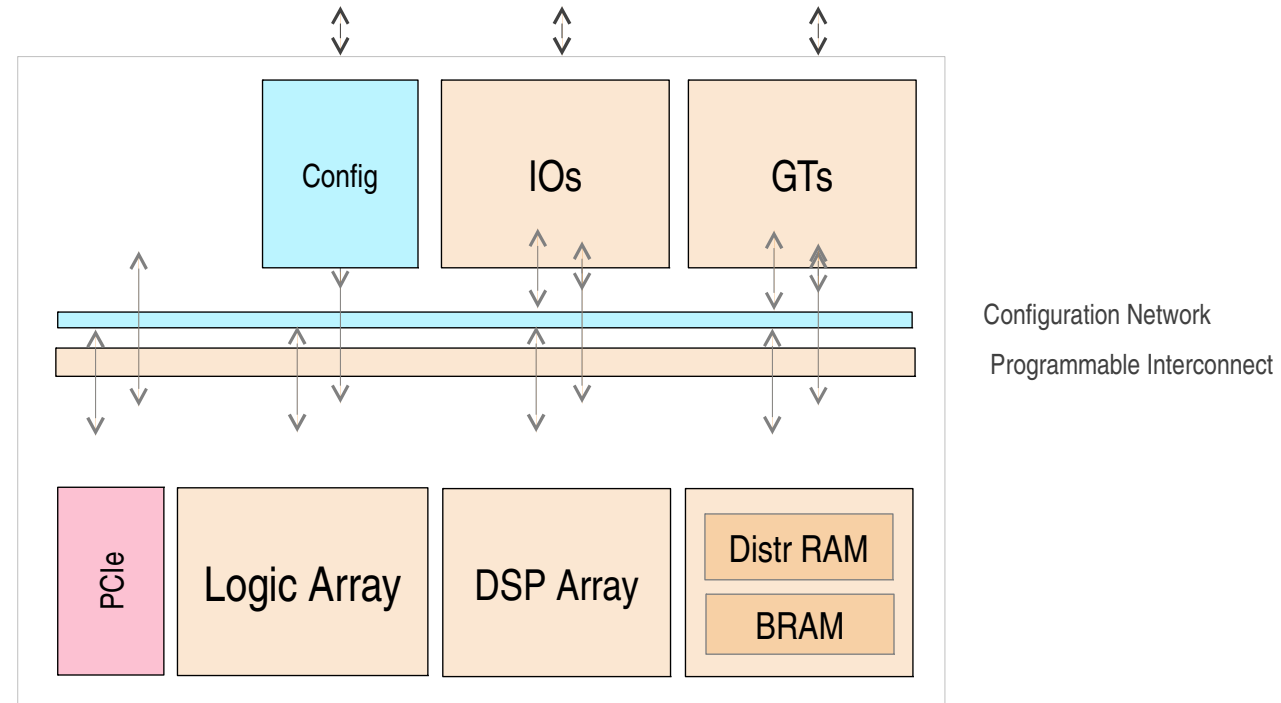
■ CPU  ■ Accelerators  ■ Memory  ■ Networking  ■ Other



Source: ARK Investment Management LLC | ⊙

**Machine Learning**
ASICs (TPU by Google)
FPGA (Xilinx, Intel)
GPU (NVidia Tesla P4)

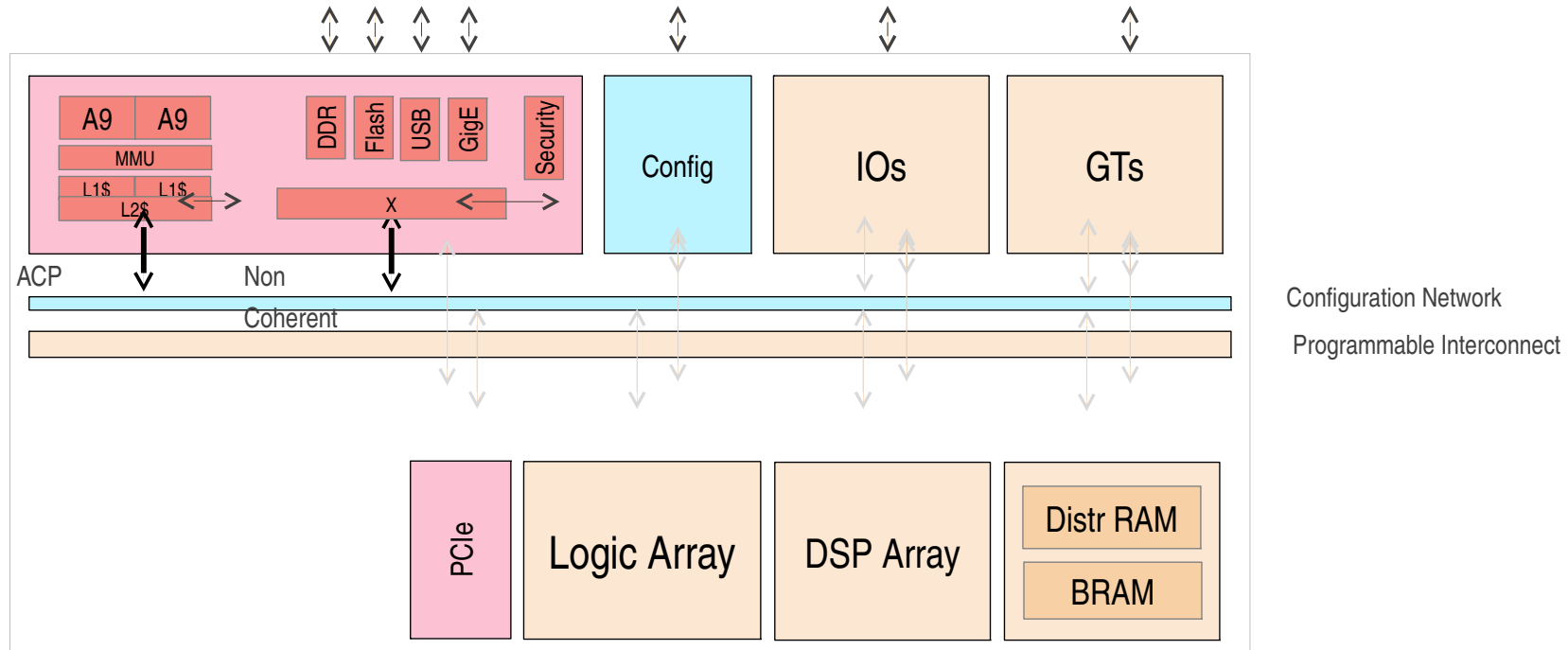**⚡ XILINX ➤** ALL PROGRAMMABLE.™

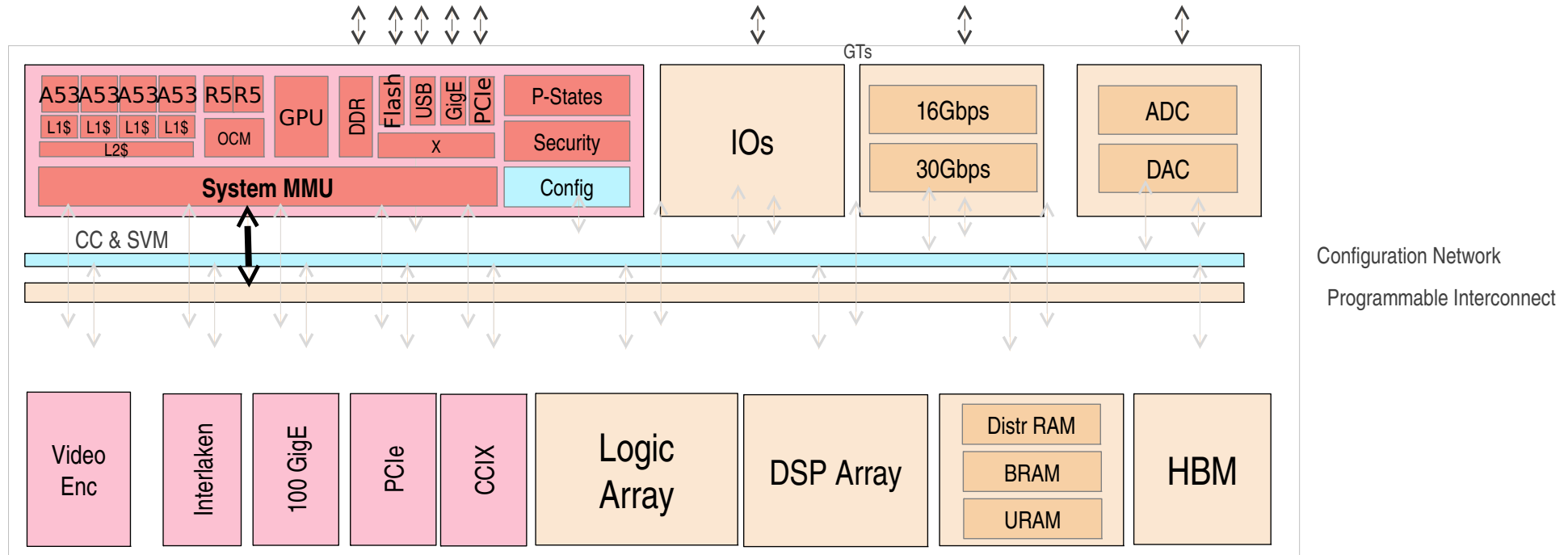# FPGA Silicon Architecture



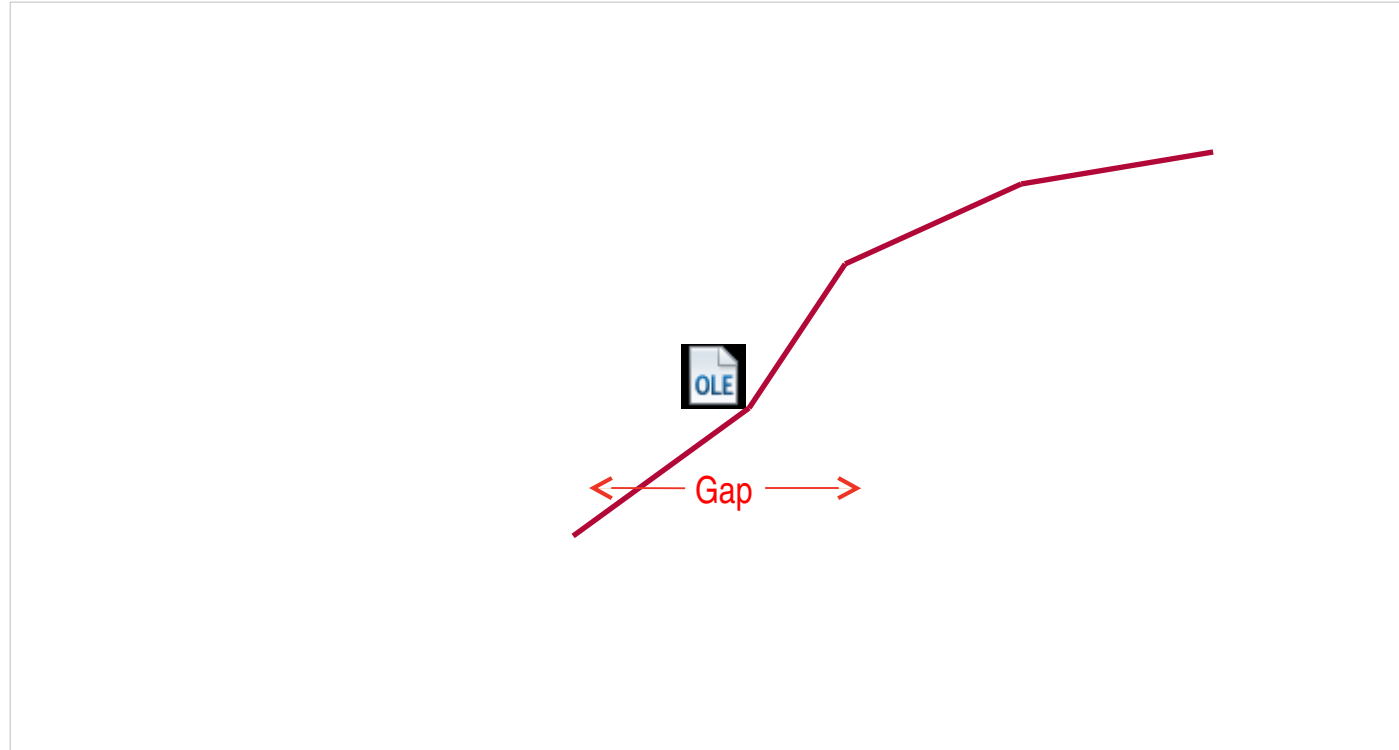Standard FPGA

# FPGA Silicon Architecture



Zynq 7000 (28nm)
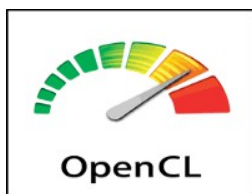
# All Programmable FPGA Silicon Architecture



Zynq MPSoC  (16nm)

# Efficiency Programming Experience



FPGA: Performance Advantage, But Productivity Gap

# Programming Heterogeneous Parallel Platforms



ASIC Refugees

HW/SW co-design

Software Programming

Logic
BRAM
DSP

RTL
Logic
VHDL

Dual A9

acc1

acc2

C/C++   RTL   RTL
CPU   Accel
C + VHDL

Quad A53

Dual R5
GPU
H.265

acc1

acc2

C/C++   RTL   C/C++
CPU   Accel
C + High Level Synthesis

Quad A53

Dual R5
GPU
H.265

acc1

acc2

C/C++   C/C++   C/C++
CPU   Accel
OpenCL

OpenCL

Heterogeneous Parallel Programming

SYCL

Bring power of C++ to OpenCL

SPIR

Intermediate Representation

HLx

SDx Environments

XILINX  ALL PROGRAMMABLE

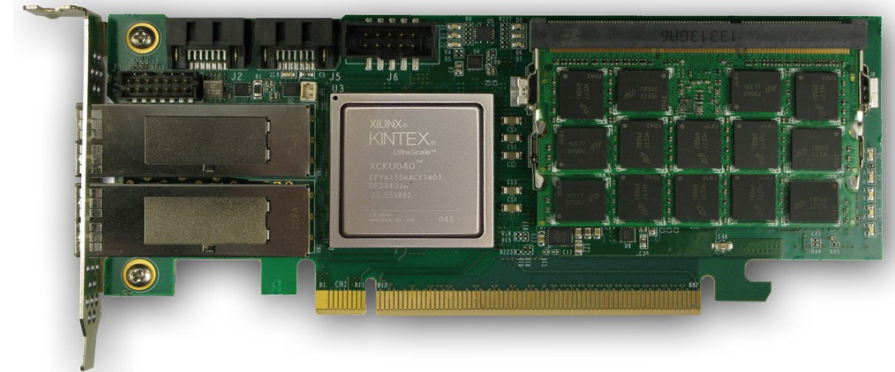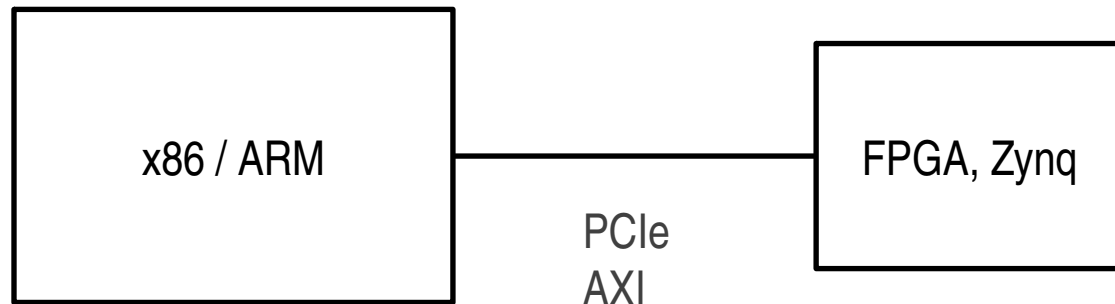# Towards single source C++

➤ No special memory allocation    (malloc)

➤ No special  data movement       (no copy needed)

➤ No special  accessing accelerator memory      (private memories)

**≡ XILINX ➤** ALL PROGRAMMABLE.™

# x86 / ARM Host: Non-Coherent



```
┌─────────────────────┐                    ┌─────────────────┐
│                     │                    │                 │
│                     │        PCIe        │                 │
│     x86 / ARM       │────────────────────│   FPGA, Zynq    │
│                     │        AXI         │                 │
│                     │                    │                 │
└─────────────────────┘                    └─────────────────┘
```

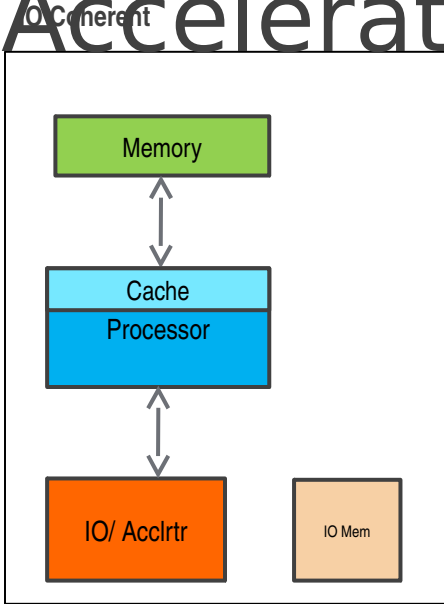| | |
|---|---|
| Memory allocation special? | Yes |
| Data movement special? | Yes |
| Accelerator memory access special? | Yes |

**Product Description**
The ADM-PCIE-KU3 is a high performance reconfigurable Half-Length, low profile x16 PCIe form factor board based on the Xilinx Kintex UltraSCALE range of Platform FPGAs. The ADM-PCIE-KU3 features two independent channels of DDR3 memory capable of 1600MT/s (fitted with two 8GB SODIMMs), high speed I/O, SATA connections, Dual QSFP ports supporting 10G Ethernet, voltage/temperature/current control and monitoring, passive air-cooled heat sink.

**Key Features & Benefits**
- Dual QSFP High Speed Communications ports
- Dual SATA High Speed Data Storage ports
- PCI Express x16 Interface
- TWO SODIMM slots

XILINX ➤ ALL PROGRAMMABLE.

# CCIX : Cache Coherent Interface for Accelerators



**IO Coherent**

**Coherent Accelerator**

**Peer- Peer**

**IO Coherency**

- Allows DMA with IO as master
- IO agent sees limited memory range
- One-way coherency. Processor memory not coherent with IO Memory
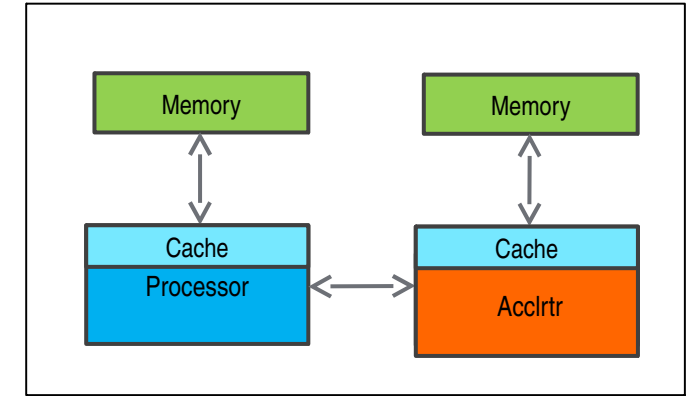- Interface optimized for large transfers
- E.g.: PCIE

**Coherent Accelerator**

- Caching (always), Home Node (limited) capability
- Typically a standard interface
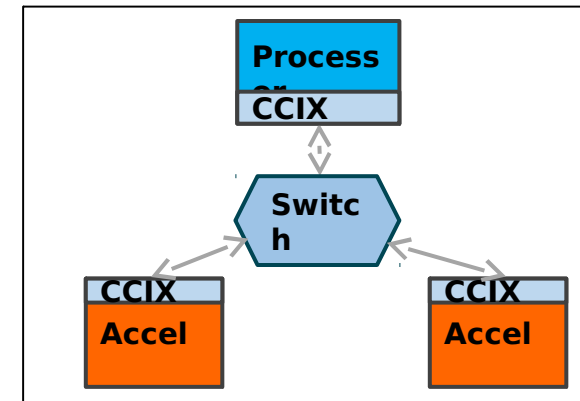- Protocol Bridging function
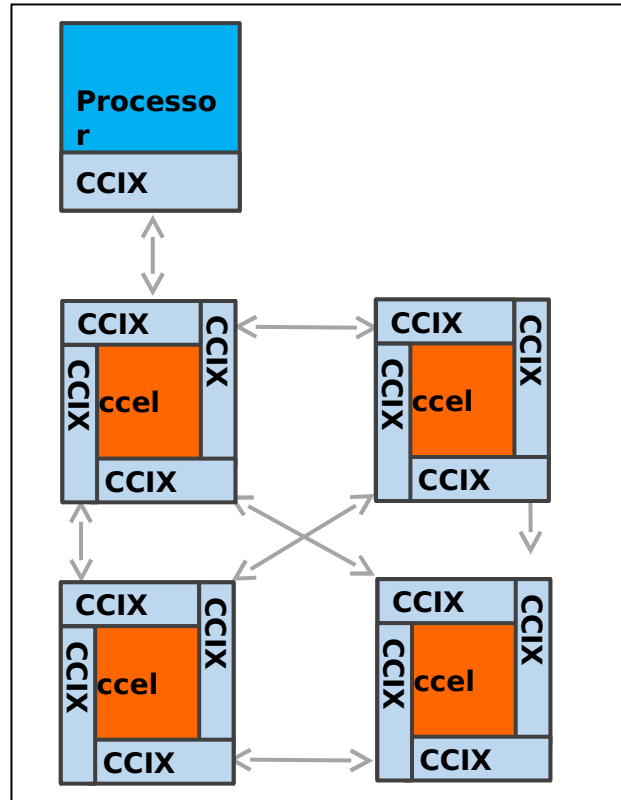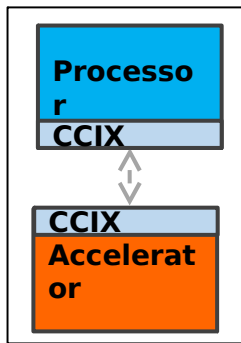- E.g.: IBM CAPI, nVidia nvLink, OpenCAPI

**Coherent Peer-peer**

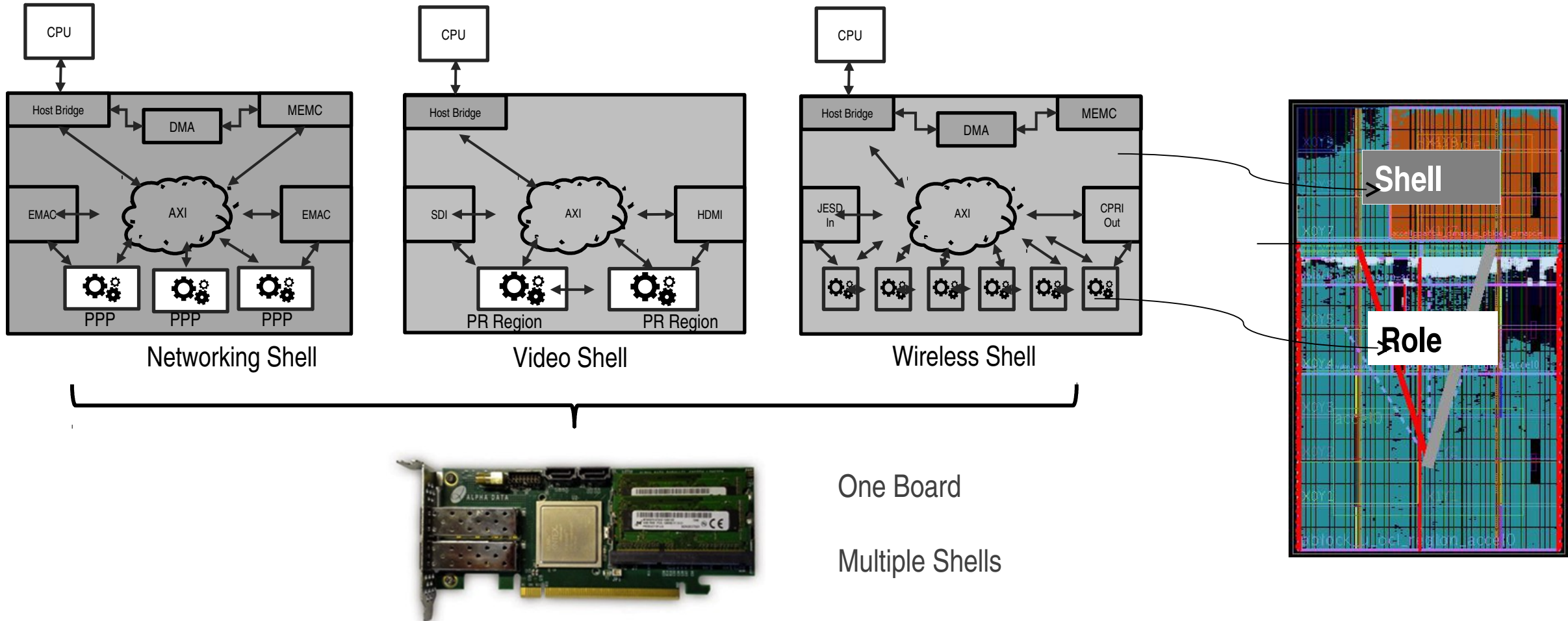- All agents can cache other agents memory
- Each agent is peer and home node
- Mostly proprietary interface
- E.g.: CCIX

## CCIX : Open Source

**XILINX** ➤ ALL PROGRAMMABLE.

# CCIX: Accelerator – CPU Configurations

# Domain Specific Platform Infrastructure



Networking Shell
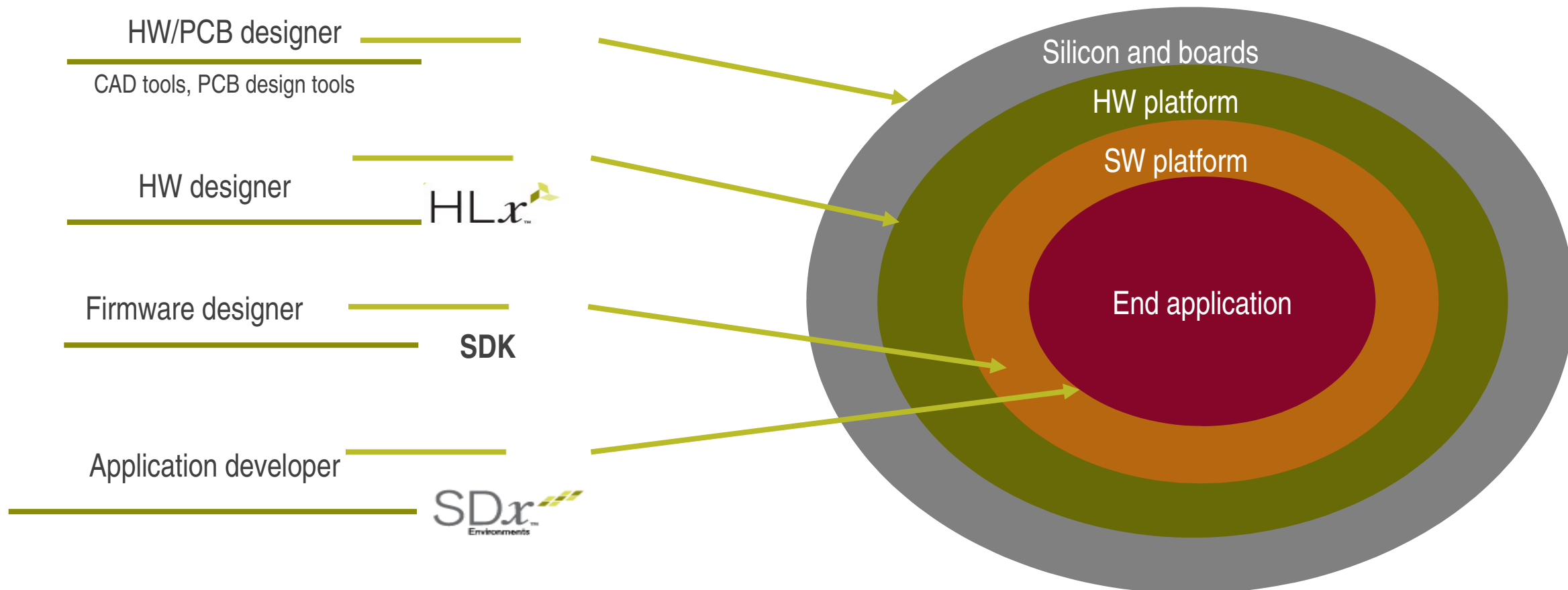
Video Shell

Wireless Shell

One Board

Multiple Shells

Shell

Role

**Shell : Domain specific infrastructure (gray)**

**Role : 'Donut holes' in FPGA or on CPU executing programmable functions (white)**

**XILINX** ➤ ALL PROGRAMMABLE.

# Use Model and Personas

HW/PCB designer

CAD tools, PCB design tools

HW designer

HLx

Firmware designer

**SDK**

Application developer

SDx
Environments

Silicon and boards

HW platform

SW platform

End application

XILINX ➤ ALL PROGRAMMABLE.

# Hardware : HLx – High-Level Design

UltraFast High-Level Design Methodology

HLS
High-Level IP
Creation

C, C++ or SystemC

Vivado™ HLS

VHDL or Verilog

C Libraries

IP Sub-systems + Config. Reference Designs

Automated IP Assembly

**Typically 5-15x productivity improvement via:**

❯ Creation of HW-optimized functions in C/C++

❯ Accelerated verification (>1000X  RTL)

❯ Automated, intelligent assembly (15x manual)
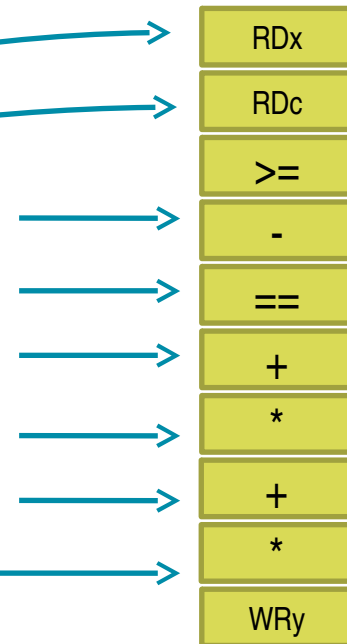
# Hardware : HLS Control & Datapath Synthesis



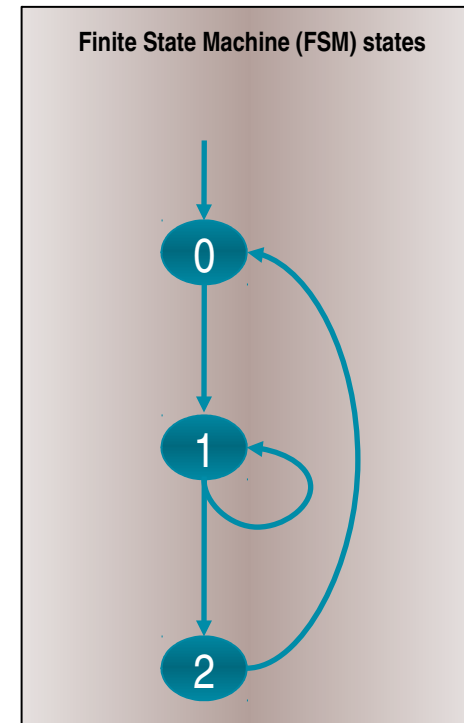**Code**

```
void fir (
  data_t *y,
  coef_t c[4],
  data_t x
) {

  static data_t shift_reg[4];
  acc_t acc;
  int i;

  acc=0;
  loop: for (i=3;i>=0;i--) {
    if (i==0) {
      acc+=x*c[0];
      shift_reg[0]=x;
    } else {
      shift_reg[i]=shift_reg[i-1];
      acc+=shift_reg[i]*c[i];
    }
  }
  *y=acc;
}
```
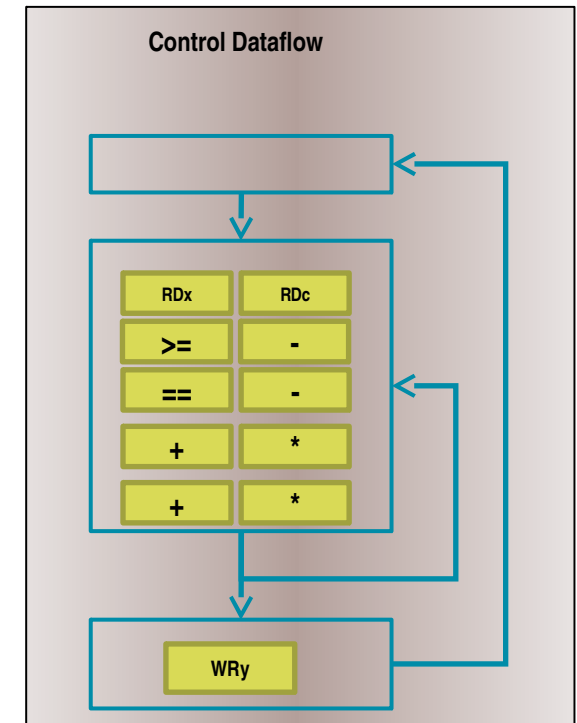
**Operations**

RDx
RDc
>=
-
==
+
*
+
*
WRy

**Control Behavior**

Finite State Machine (FSM) states

0
1
2

**Control & Datapath Behavior**

Control Dataflow

| RDx | RDc |
| >= | - |
| == | - |
| + | * |
| + | * |

WRy

**From any C code example ..**

**Operations are extracted…**

**The control is known**

**A unified control dataflow behavior is created.**

**XILINX** ➤ ALL PROGRAMMABLE.

# Hardware : IPI  Automated IP Integration



**IP Assembly Example:**

**Zynq Processor Subsystem**

**+ Video Subsystem**

**+ 6 IP Blocks**

Video Processing IP Subsystem

**4700 lines of VHDL**
**(top-level connectivity only)**

# Creation of fixed Shell infrastructure

# Hardware Abstraction : Runtime Layers

**API**

OpenCL

Provides API view of platform

**Xilinx Runtime (XRT)**

Core runtime services: buffer management, accelerator scheduling

**Hardware Abstraction Layer (HAL)**

XDMA
libxcldrv.so

MPSoC/Zynq
libzynqdrv.so

Common hardware abstraction

**Linux Kernel Driver**

XDMA
xdma.ko, xclmgmt.ko

MPSoC/Zynq
zoclsvm.ko

Heavy lifting: hardware programming, DMA, Linux VM interaction

**XILINX** ➤ ALL PROGRAMMABLE.

# Overall Platform and SDx Flow



Application developer

Software Platform
(Board Support Package)

HW platform
(DSA/...)

HLx

Vivado HLS,
SysGen

Vivado IPI

Vivado RTL

Hardware
persona

Firmware
Designer

SDK

Linux

Yocto

Hardware description (.hdf)

XILINX ALL PROGRAMMABLE.

# Example : Xilinx Machine Learning Stack

Customers

OpenSource

Xilinx

Xilinx/Partner

Reference Networks (e.g. AlexNet) & Custom Networks

CNN Network Design Tools & Training Frameworks (e.g. CAFFE)

CNN Compiler & Runtime

APIs & Libraries (Primitive operators)

SDx

HW-SW development platform

**XILINX** ➤ ALL PROGRAMMABLE.

# The Future : Single Source SYCL SPIR

OpenCL
SYCL

OpenMP 4+

C++ Based
Frontends

DSELs

#pragmas

Heterogenous
Parallel & Concurrent
Pipes & fine-grain dataflow

SPIR-V

SDx Backend
Interface

Runtime

Platforms
DSA

SDK

Debug
Profile

Compiler
HLS

SDx Backend
Infrastructure

XILINX ➤ ALL PROGRAMMABLE.

# From the Cloud to IOT



**Desktop**

**Mobile**

**Server**

**Embedded**

**Data Center**

IOT

CLOUD

**Real-time**
**Deterministic**

**The Edge/Fog**

**Performance**
**Scale**

**XILINX** ➤ ALL PROGRAMMABLE.

# Productivity languages and Efficiency languages

**Application**

**Implementation**

Applications,
Programming Frameworks

**Productivity Languages**
Python, Scala, ..

**Efficiency Languages**
C, C++, OpenCL

Hardware Systems
OS, hypervisor, drivers

**Programming ZYNQ/MPSoC in a productivity language**



COMMUNICATIONS OF THE ACM

HOME | CURRENT ISSUE | NEWS | BLOGS | OPINION | RESEARCH | PRACTICE | CAREERS | ARCHIVE | VIDEOS

Home / Blogs / BLOG@CACM / Python is Now the Most Popular Introductory Teaching... / Full Text

BLOG@CACM

## Python is Now the Most Popular Introductory Teaching Language at Top U.S. Universities

Number of top 39 U.S. computer science departments that use each language to teach introductory courses

Python · Java · MATLAB · C · C++ · Scheme · Scratch

Analysis done by Philip Guo (www.pgbovine.net) in July 2014, last updated 2014-07-29

**XILINX** **ALL PROGRAMMABLE.**

# Python-based Open Source Platform : PYNQ

Overlays:
- IO programming
- OpenCV
- CNN

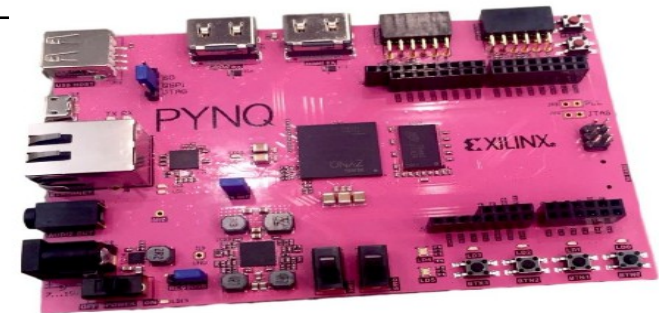| Web Server | |
|---|---|
| Python VM | Pynq Libraries |
| Ubuntu Server | Overlays |
| ARM | FPGA |

**Architecture emphasizes :**

- a software-centric approach
- based on open, de facto standards
- platform, OS and browser agnostic
- minimal learning curve
- no proprietary methodologies

**SW running natively on Zynq**

XILINX ALL PROGRAMMABLE.

# Summary

**HW designers:**

**C-based IP development + high-level IP assembly**

**SW developers:**

**FPGA-based acceleration using SDx**

**XILINX**

**Committed to major investments in next generation silicon and tools that will revolutionize programming All Programmable FPGA**

**XILINX** ➤ ALL PROGRAMMABLE.