# Heterogeneous Multi-Processing for SW-Defined Multi-Tiered Storage Architectures

Endric Schubert (MLE)
Ulrich Langenbach (MLE)
Michaela Blott (Xilinx Research)

*SDC, 2017*

# Content

**Heterogeneous Multi-Processing for
Software-Defined Multi-Tiered Storage Architectures**

- Who – Xilinx Research and Missing Link Electronics

- Why – Multi-tiered storage needs predictable performance scalability, deterministic low-latency and cost-efficient flexibility / programmability

- What – Tera-OPS processing performance in a single-chip heterogeneous compute solution running Linux

- How – Combine "unconventional" dataflow architectures for acceleration & offloading  with Dynamic Partial Reconfiguration and High-Level Synthesis

© MLE

# Xilinx Research and Missing Link Electronics

# Xilinx – The All Programmable Company



XILINX
ALL PROGRAMMABLE™

**XILINX - Founded 1984**

🟢 Headquarters    🟡 Sales and Support

🔵 Research and Development    🔴 Manufacturing

**$2.38B** FY15 revenue

**>55%** market segment share

**3,500+** employees worldwide
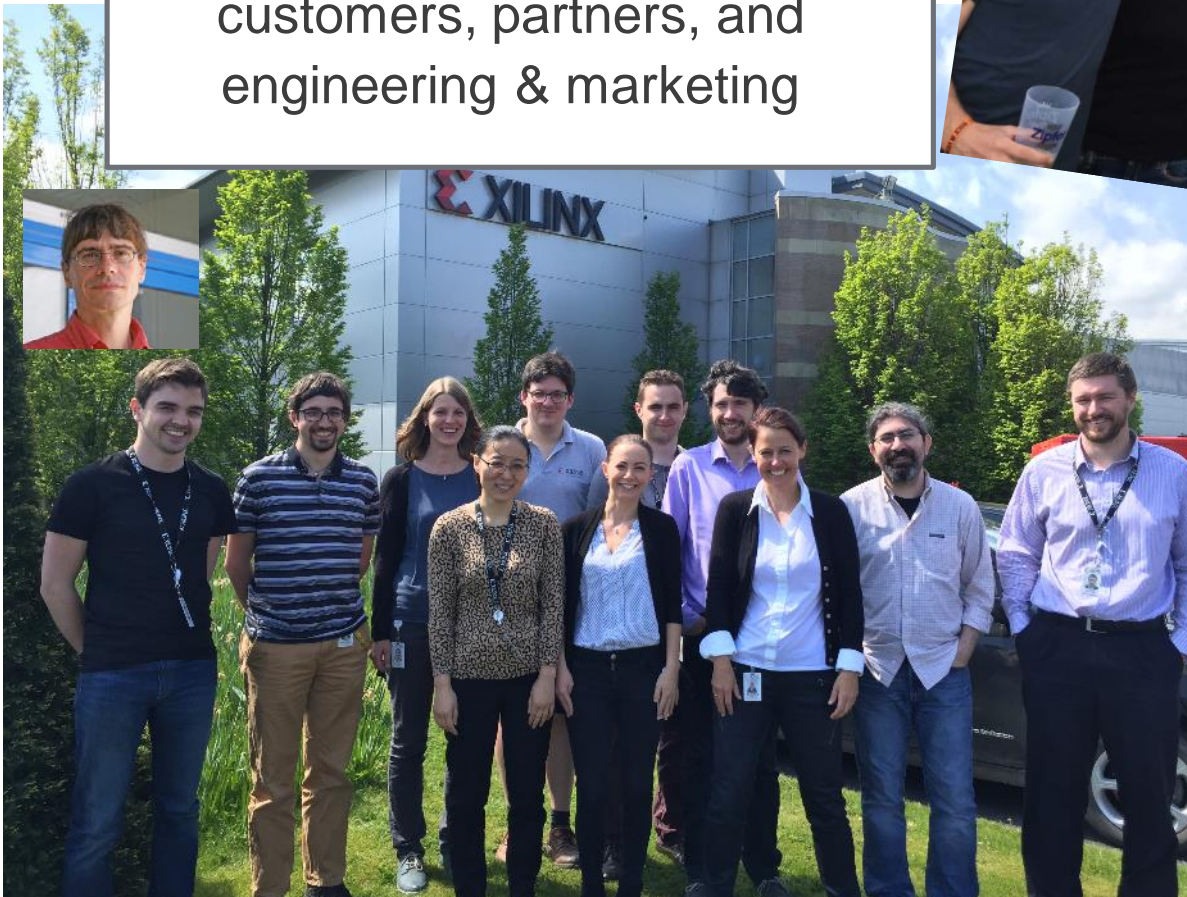
**20,000** customers worldwide

**3,500+** patents

**60** industry firsts

© MLE

XILINX › ALL PROGRAMMABLE™

# Xilinx Research - Ireland

**Applications & Architectures**

Through application-driven technology development with customers, partners, and engineering & marketing

© MLE

# Missing Link Electronics
*Xilinx Ecosystem Partner*

Vision: The convergence of software and off-the-shelf programmable logic opens-up more economic system realizations with predictable scalability!

Mission: To de-risk the adoption of heterogeneous compute technology by providing pre-validated IP and expert design services.

Certified Xilinx Alliance Partner since 2011, Preferred Xilinx PetaLinux Design Service Partner since 2013.

# Missing Link Electronics Products & Services

TCP/IP & UDP/IP Network Protocol Accelerators at 10/25/50 GigE line-rate.

Low-Latency Ethernet MAC form German Fraunhofer HHI.

Patented Mixed Signal systems solutions with integrated Delta-Sigma converters in FPGA logic.

Key-Value-Store Accelerator for hybrid SSD/HDD memcached and object storage.

SATA Storage Extension for Xilinx Zynq All-Programmable Systems-on-Chip.

A team of FPGA and Linux engineers to support our customer's technology projects in the USA and Europe.

**XILINX > ALL PROGRAMMABLE.**

# Motivation

© MLE

# Technology Forces in Storage

- **Software significantly impacts latency and energy efficiency in systems with nonvolatile memory**

- **However, software-defined flexibility is necessary to fully utilize novel storage technologies**

- **Hyper-capacity hyper-converged storage systems need more performance, but within cost and energy envelopes**



*Source: Steven Swanson and Adrian M. Caulfield, UCSD*
*IEEE Computer, August 2013*

© MLE

# The Von Neumann Bottleneck [J. Backus, 1977]

**CPU system performance scalability is limited**



40 Years of Microprocessor Trend Data

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

New Compute Architectures are needed

© MLE

# Spatial vs. Temporal Computing



Sequential Processing
with CPU

Parallel Processing
with Logic Gates

*Source: Dr. Andre DeHon, Upenn: "Spatial vs. Temporal Computing"*

➤ **CPU system performance scalability is limited**

➤ **Spatial computing offers further scaling opportunity**

New Compute Architectures are needed to
take advantage of this

© MLE

# Architectural Choices for Storage Devices



*Source: T.Noll, RWTH Aachen*

# Use Case: Image/ Video Storage

# A Flexible All Programmable Storage Node

© MLE

# Meta Data Extraction, e.g. Image Quality Metrics



Storage Node

- Partially under and over exposed
- Medium contrast

NVMe Drive

NVMe Drive

NVMe Drive

Network

Meta Data Extraction

Monitoring

© MLE

# Processing, e.g. Thumbnailing, Auto-Correction

© MLE

# Semantic Feature Extraction, e.g. Classification

© MLE

# Semantic Search Support



Storage Node

Network

Semantic Feature Search

Monitoring

NVMe Drive

NVMe Drive

NVMe Drive

- Group of People
- Xilinx
- Outdoor Scene

© MLE

# Performance Metrics, e.g. Bandwidth, Latency



Storage Node

NVMe Drive

NVMe Drive

NVMe Drive

Network

- Performance Counter
- Pattern Matching
- ID Generation for Tracing

Monitoring

© MLE

**XILINX** ➤ ALL PROGRAMMABLE.

# Runtime Programmability



Meta Data Extraction

Semantic Feature Extraction

Storage Node

Network

Reconfigurable Processing

Monitoring

NVMe Drive

NVMe Drive

NVMe Drive

# Architectural Concepts

© MLE

# Key Concepts Presented at SDC-2016

➤ **Heterogeneous compute device as a single‑chip solution**

➤ **Direct network interface with full accelerator for protocols**

➤ **Performance scaling with dataflow architectures**

➤ **Scaling capacity and cost with a Hybrid Storage subsystem**

➤ **Software-defined services**

© MLE

# SDC-2016: Single-Chip Solution for Storage



**Data Node**

Processing System with quad core 64b processors (A53)

Network management

Memory management

Petalinux

FPGA fabric (PL)

NVMe interface

Router

Key Value Store Abstraction (memcached)

Hybrid Memory System

TCP/IP stack

Memory controller

DDRx channels

M.2 NVMe drives

DDRx channels

MLE IP | SD Services | Xilinx IP

© MLE

# SDC-2016: Hardware Accelerated Network Stack



Data Node

Processing System with quad core 64b processors (A53)

Network management

Memory management

Petalinux

FPGA fabric (PL)

Router

Key Value Store Abstraction (memcached)

Hybrid Memory System

NVMe interface

Memory controller

TCP/IP stack

**TCP/IP Full Accelerator Supports 10/25/50 GigE line-rates**

DDRx channels

M.2 NVMe drives

DDRx channels

© MLE

# SDC-2016: Dataflow architectures for performance scaling



**Streaming Architecture:** Flow-controlled series of processing stages which manipulate and pass through packets and their associated state

- **Now**: **10 Gbps demonstrated with a 64b data path @ 156MHz using 20% of FPGA**
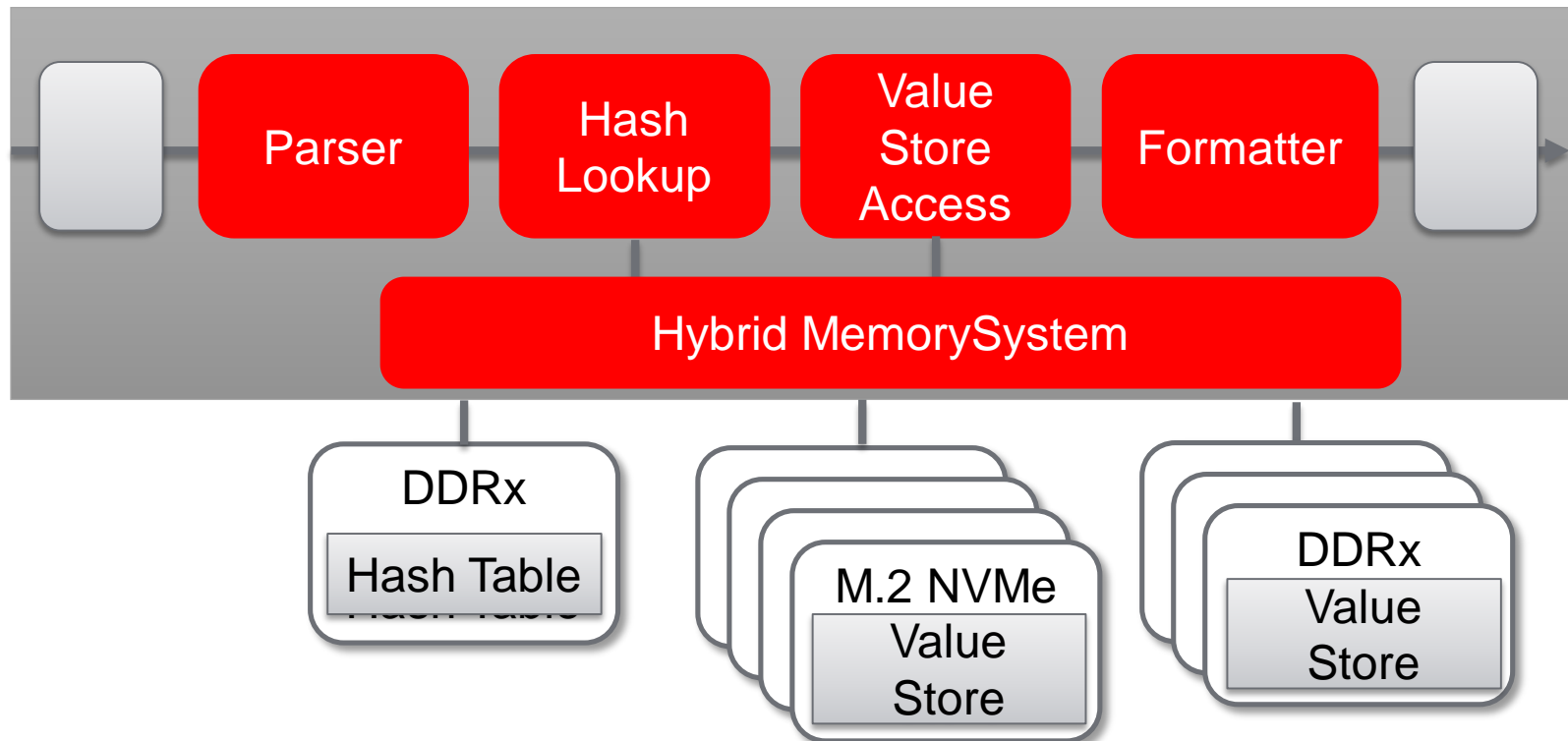- **Next**: **100 Gbps can be achieved by using a 512b @ 200MHz pipeline for example**

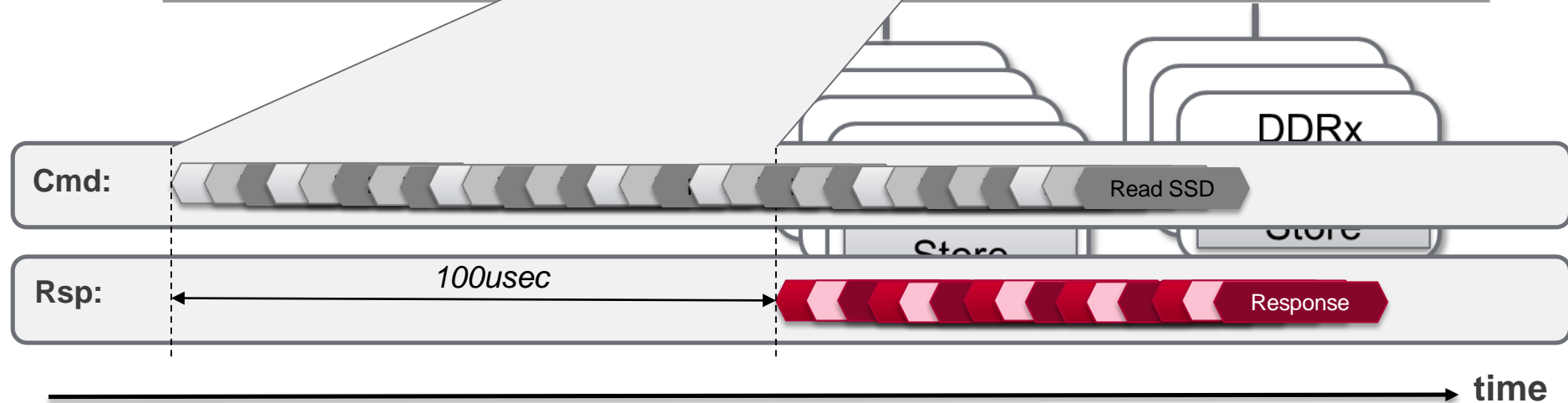*Source: Blott et al: Achieving 10Gbps line-rate key-value stores with FPGAs; HotCloud 2013*
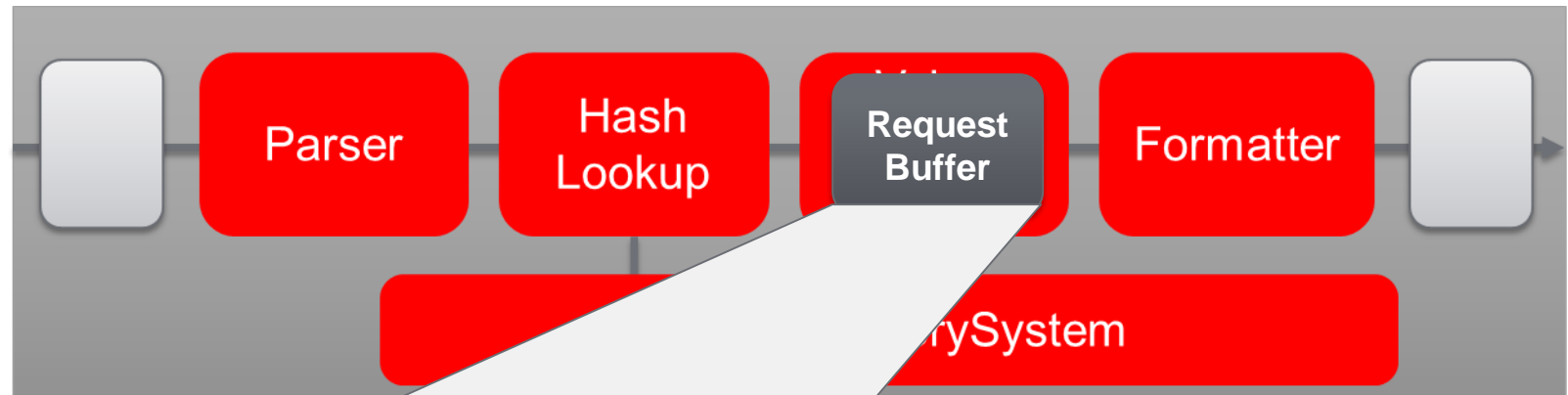
© MLE

# SDC-2016: Scaling Capacity via hybrids

▶ **SSDs combined with DDRx channels can be used to build high capacity & high performance object stores**

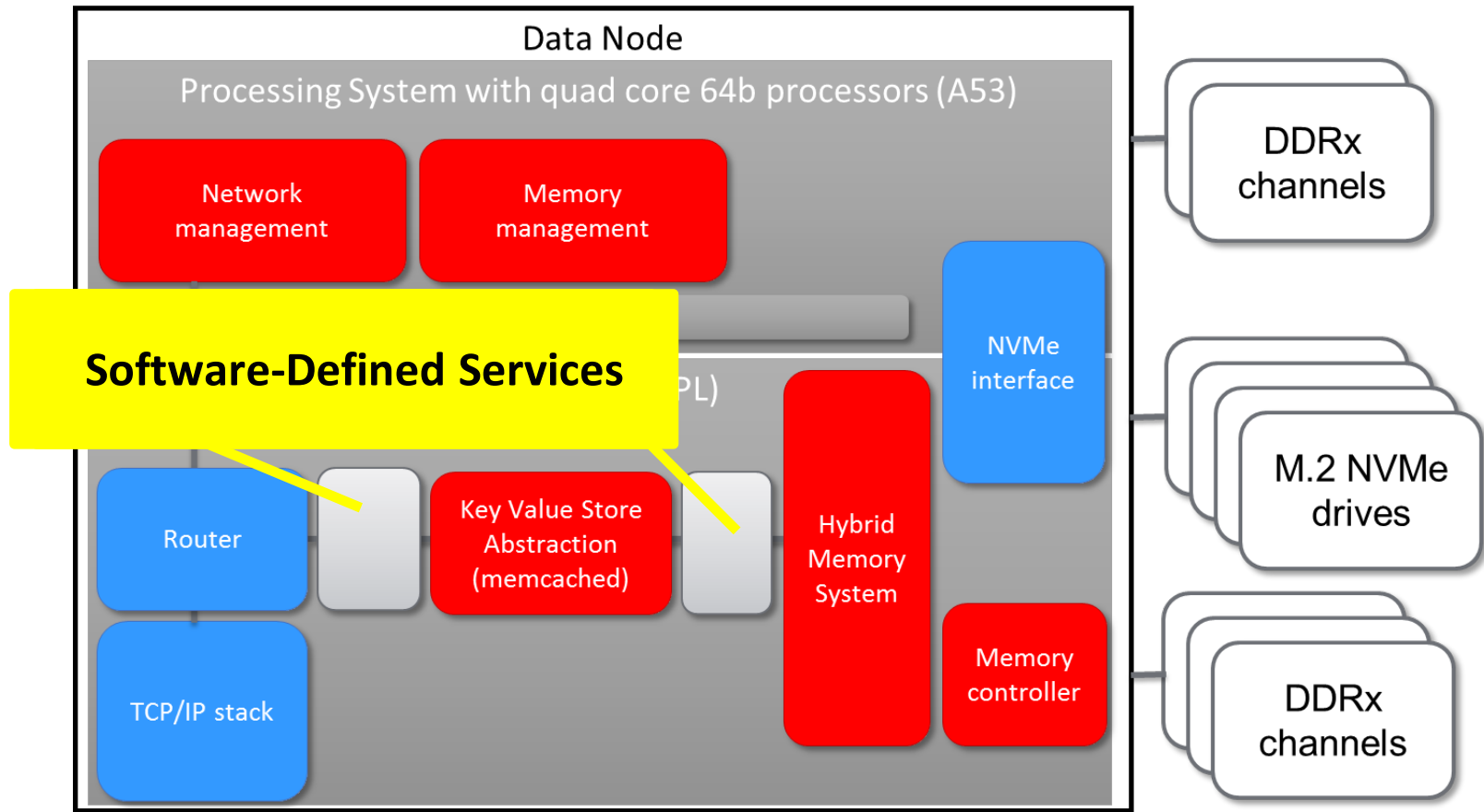▶ **Concepts and early prototype to scale to 40TB & 80Gbps key value stores**



*Source: HotStorage 2015, Scaling out to a Single-Node 80Gbps Memcached Server with 40Terabytes of Memory*

© MLE

# SDC-2016: Handling High Latency Accesses without Sacrificing Throughput



- **Dataflow architectures: no limit to number of outstanding requests**
- **Flash can be serviced at maximum speed**

© MLE

# Software-Defined Services



> **Spatial computing of additional services at no performance cost until resource limitations are reached**

# Software-Defined Services

# Software-Defined Services – Proof-of-Concepts
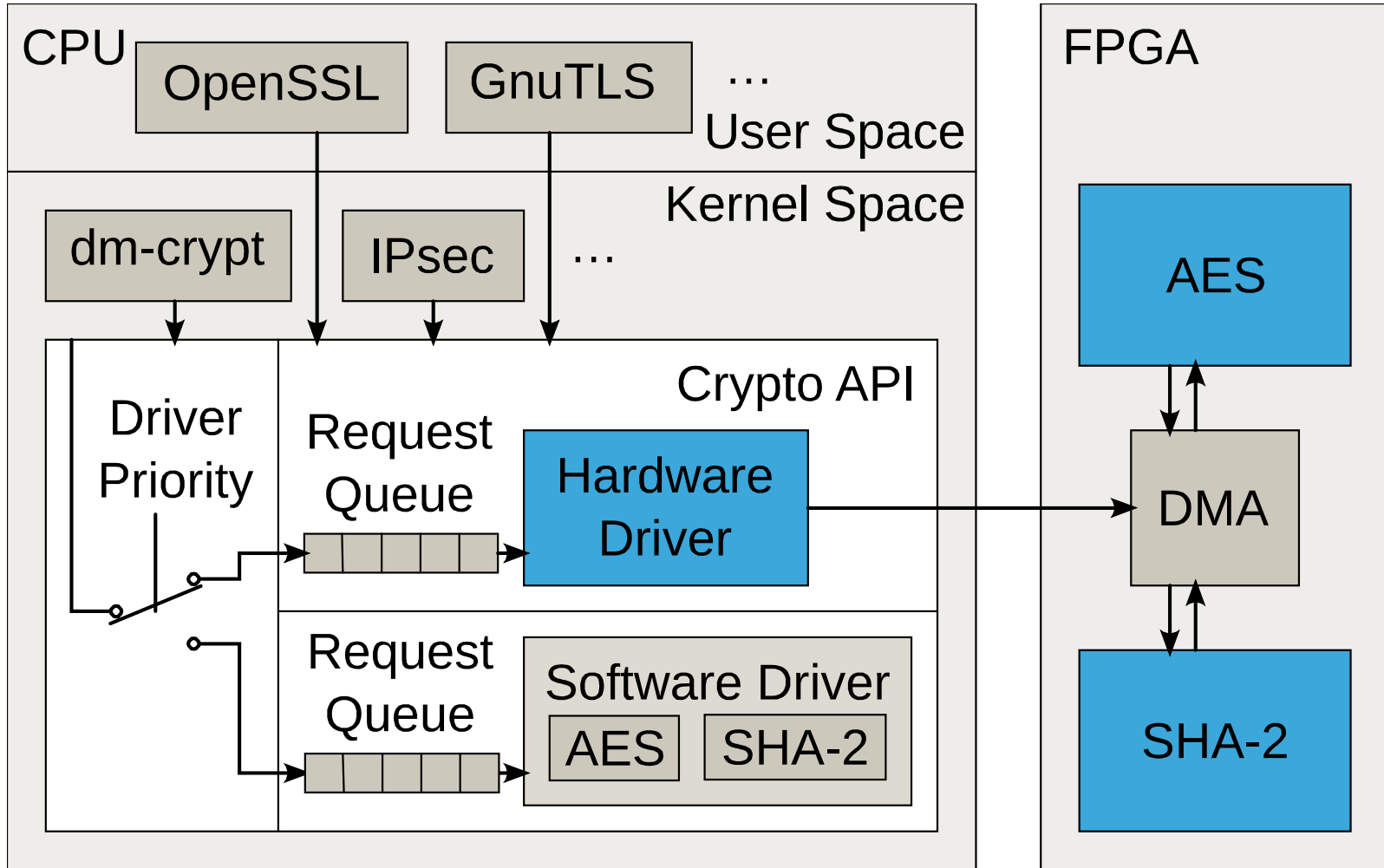
➤ **Offload engines for Linux Kernel Crypto-API**

➤ **Non-intrusive latency analysis via PCIe TLP "Tracers"**

➤ **Inline processing with Deep Convolutional Neural Networks**

➤ **Declarative Linux Kernel Support Partial Reconfiguration**

© MLE

# Software-Defined Services - Example 1) Accelerating the Linux Kernel Crypto-API

- **Crypto-API is a cryptography framework in the Linux kernel used for encryption, decryption, compression, de-compression, etc.**

- **Needs acceleration to support processing at higher line-rates (100 GigE).**

- **Open Source software implementation that follows a streaming dataflow processing architecture**
  - Hardware Interface: AXI Streaming
  - Software/ Hardware Interface: SG-DMA in, SG-DMA out

- **High-Level Synthesis generated accelerator blocks from reference C code**

© MLE

# System Architecture of Crypto-API Accelerator
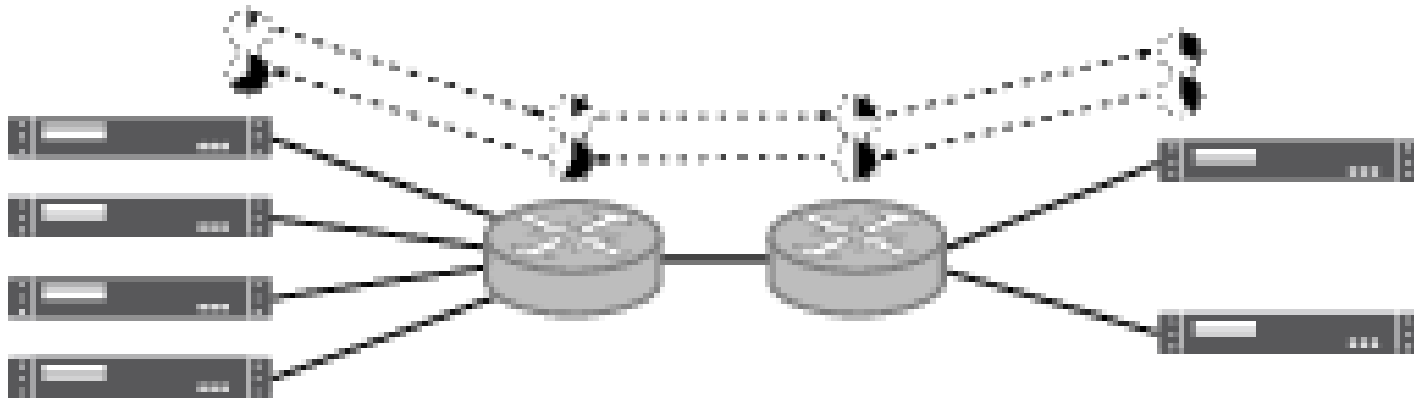
© MLE

# Software-Defined Services - Example 2) Non-Intrusive Latency Analysis via PCIe TLP Tracers

- **Performance analysis and ongoing monitoring of bandwidth <u>and</u> latency in distributed systems is difficult.**
  - Round-trip times
  - Time-outs
  - Throttling

- **When done in software, results get distorted by additional compute burden.**

- **When done in Programmable Logic, it can be (clock cycle) accurate and non-intrusive via adding so-called "Tracers" into the dataflow.**
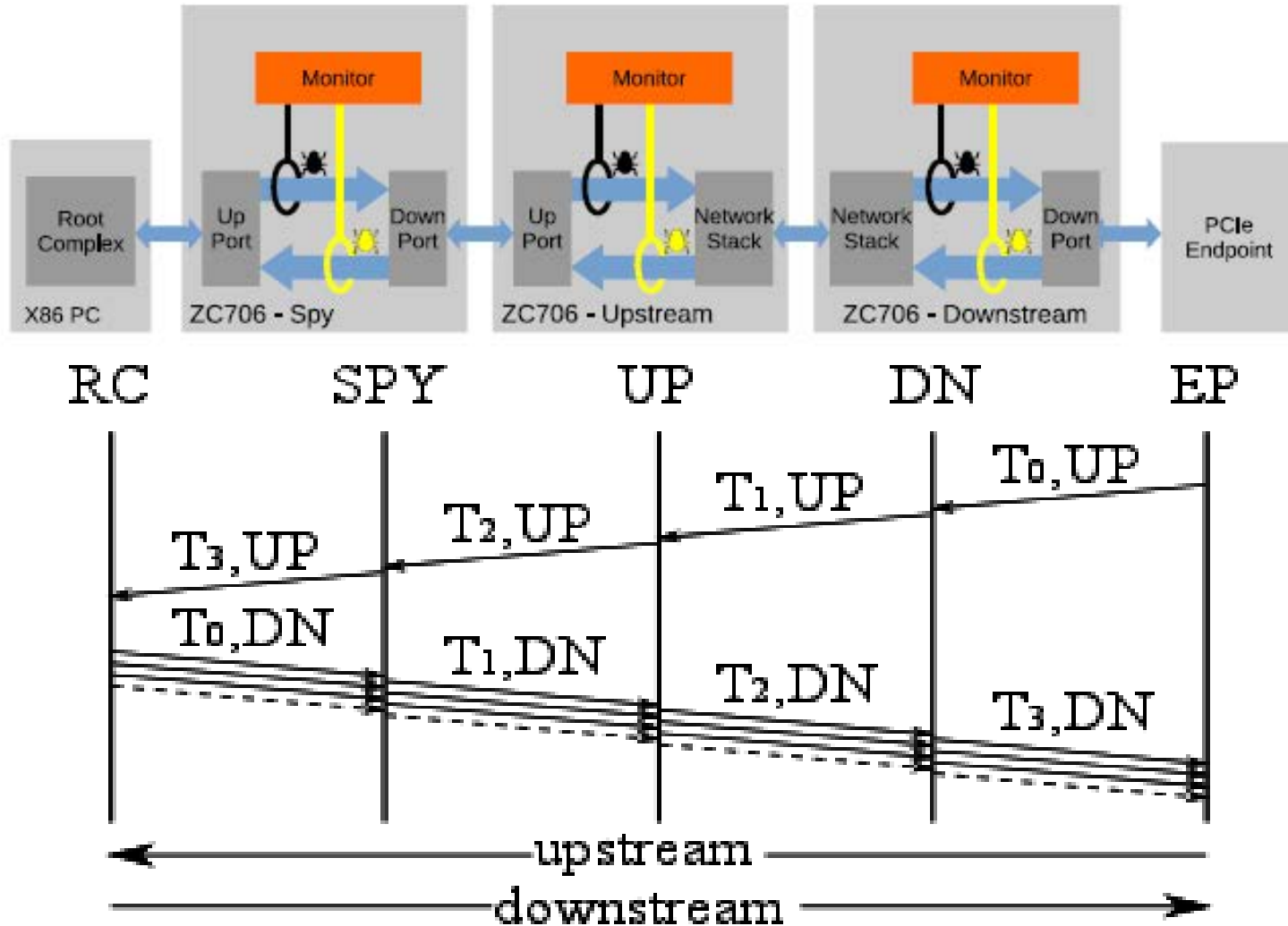
© MLE

# Tracer-Based Performance Analysis

> **Tracers within PCIe Transaction Layer Packets (TLP)**

– Based on addresses/ IDs, detected at PCIe switches and endpoints

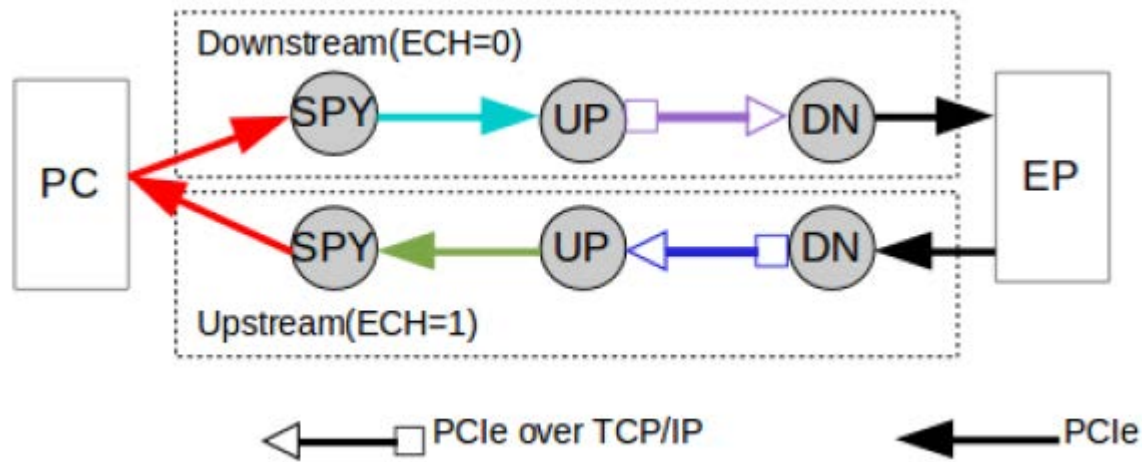– Transparent for transport layer (Ethernet, etc)

# Proof-of-Concept Implementation
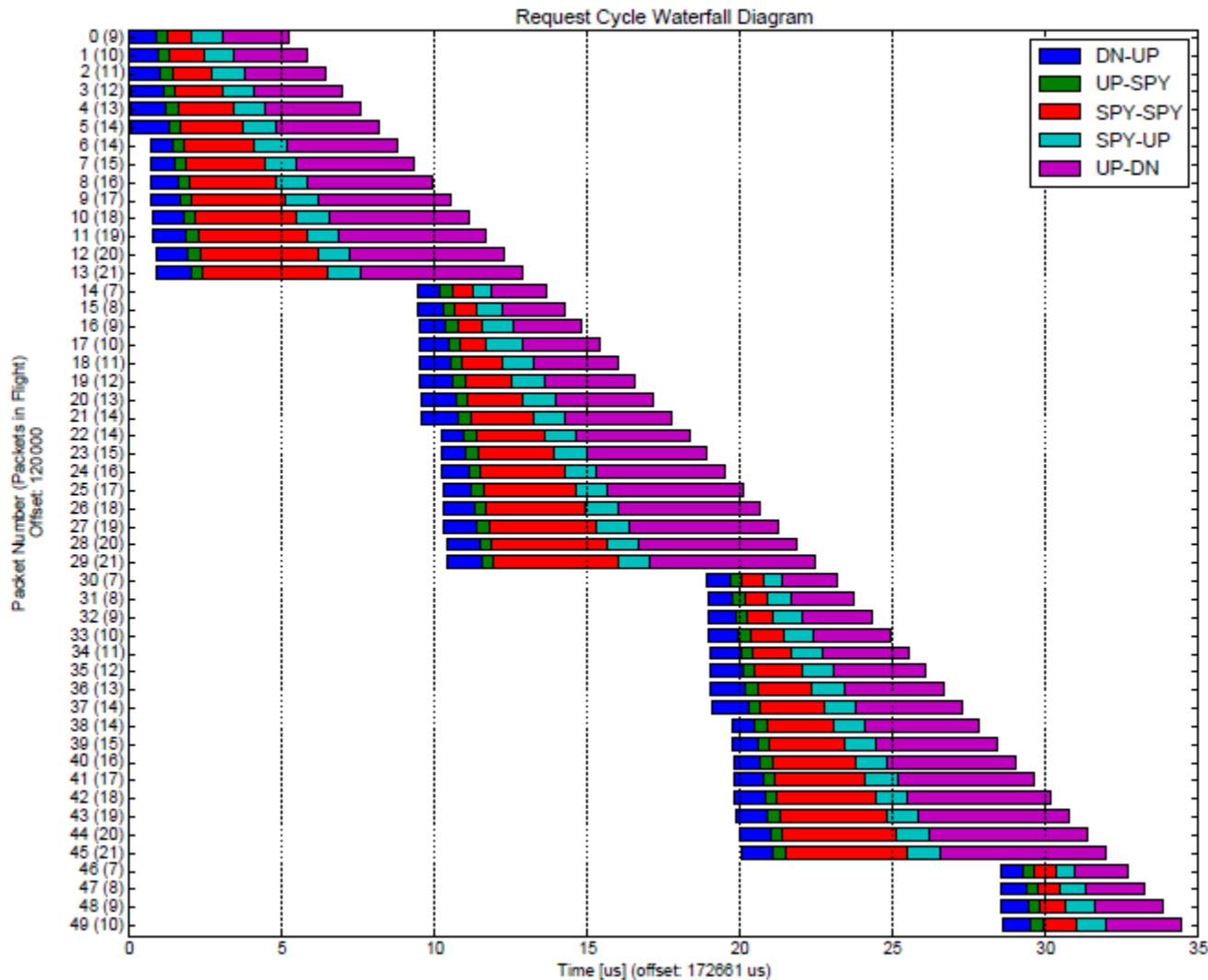
> **Full implementation on network with multiple boards**

© MLE

© MLE

# Latency Monitoring with Tracers - Results

© MLE

# Software-Defined Services - Example 3) Inline Processing w/ Neural Networks

- **Deep Convolutional Neural Networks (CNN) have demonstrated values in classification, recognition and data-mining.**

- **However, CNN can be very compute intensive, when done at single or double float precision.**

- **Recent approaches involve reduced precision (INT8, or even less), as well as dataflow-oriented compute architectures.**
  - Taps into tremendous compute power within Programmable Logic


- **What if, CNN can be run close to the data, within the storage node?**

© MLE

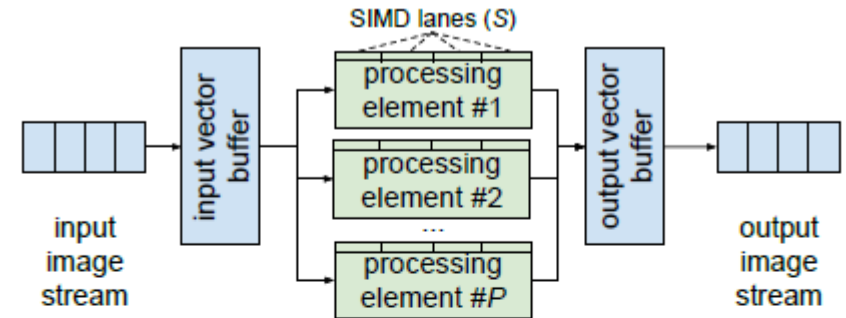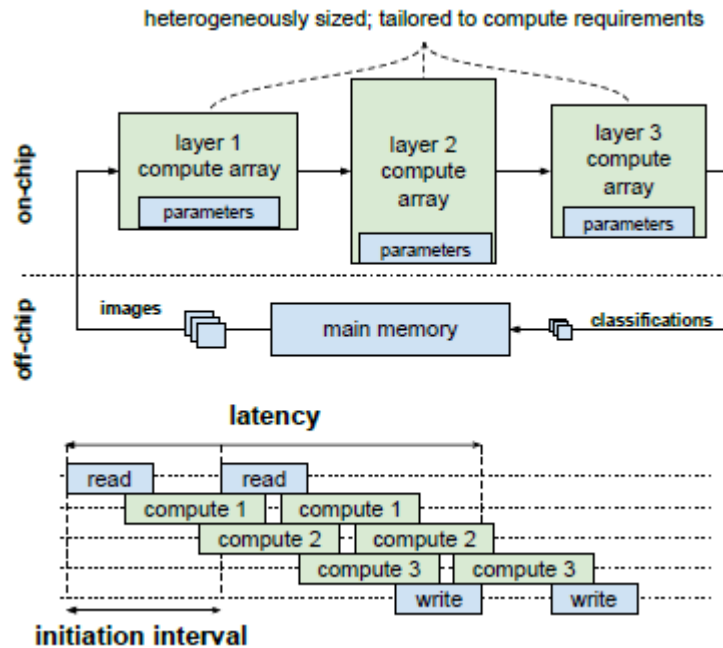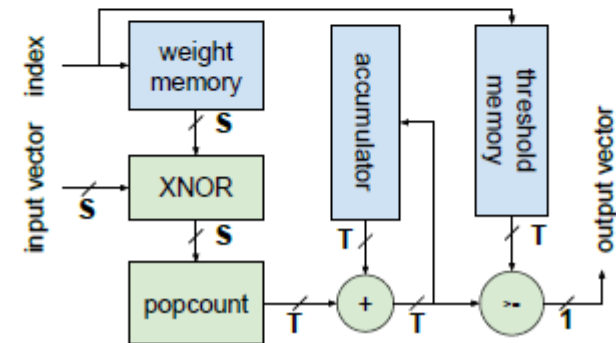# Streaming Dataflow Processing in BNN Inference



Figure 5: Overview of the MVTU.

Courtesy "FINN: A Framework for Fast, Scalable Binarized Neural Network Inference",
Umuroglu, Fraser, Blott et al., 25th Symp. on FPGA, 2017

© MLE

# BNN Results

Table 3: Summary of results from FINN 200 MHz prototypes.

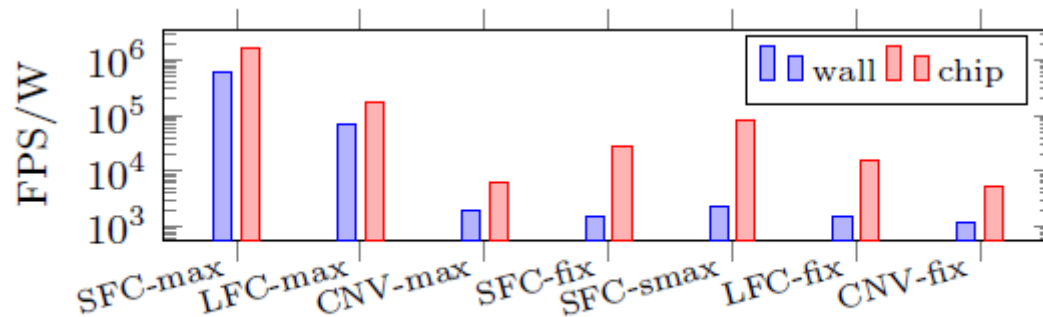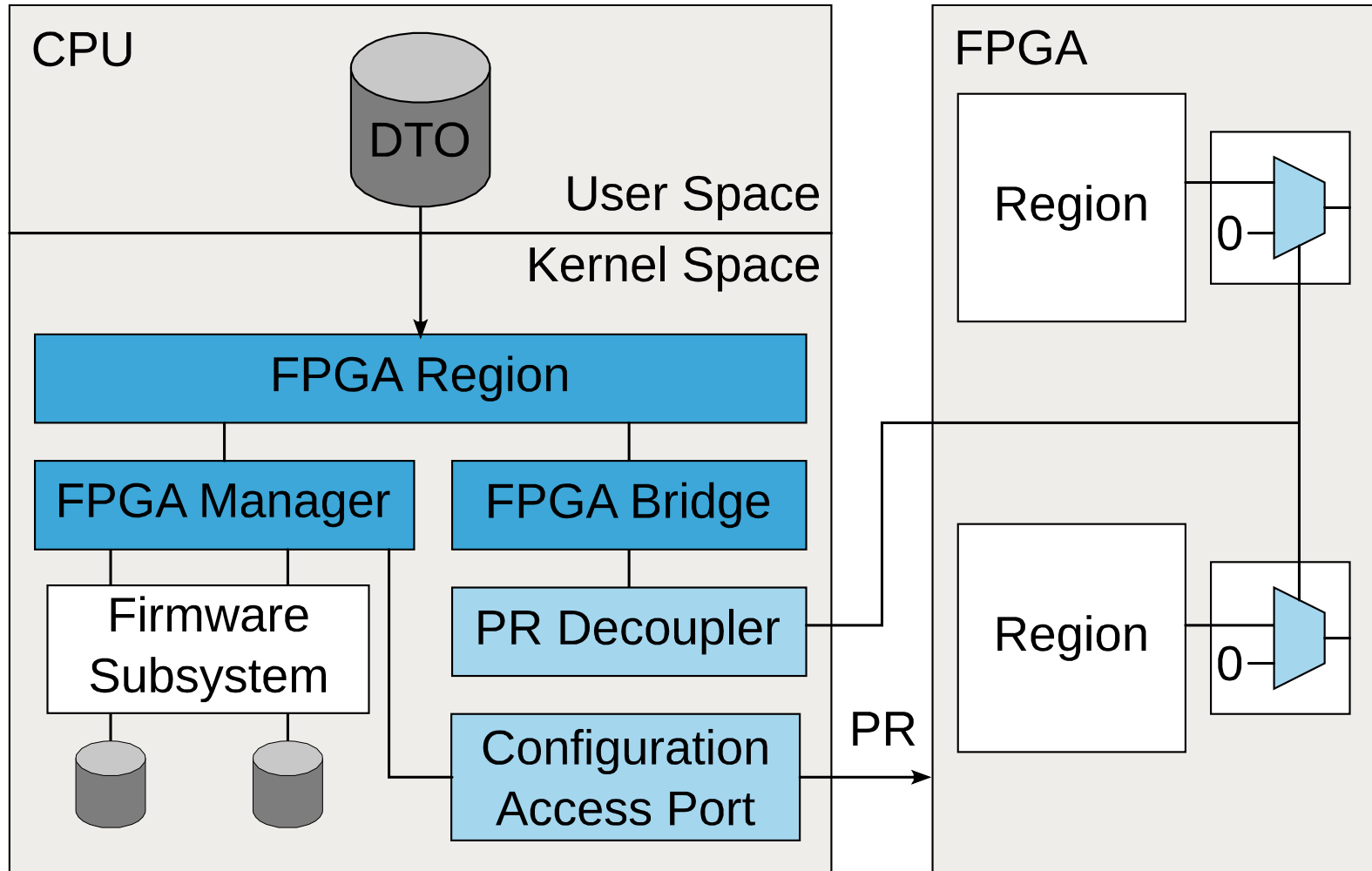| Name | Thr.put (FPS) | Latency ($\mu s$) | LUT | BRAM | $P_{chip}$ (W) | $P_{wall}$ (W) |
|---|---|---|---|---|---|---|
| SFC-max | 12361 k | 0.31 | 91131 | 4.5 | 7.3 | 21.2 |
| LFC-max | 1561 k | 2.44 | 82988 | 396 | 8.8 | 22.6 |
| CNV-max | 21.9 k | 283 | 46253 | 186 | 3.6 | 11.7 |
| SFC-fix | 12.2 k | 240 | 5155 | 16 | 0.4 | 8.1 |
| LFC-fix | 12.2 k | 282 | 5636 | 114.5 | 0.8 | 7.9 |
| CNV-fix | 11.6 k | 550 | 29274 | 152.5 | 2.3 | 10 |



Figure 10: Prototype energy efficiency.

Courtesy "FINN: A Framework for Fast, Scalable Binarized Neural Network Inference", Umuroglu, Fraser, Blott et al., 25[th] Symp. on FPGA, 2017

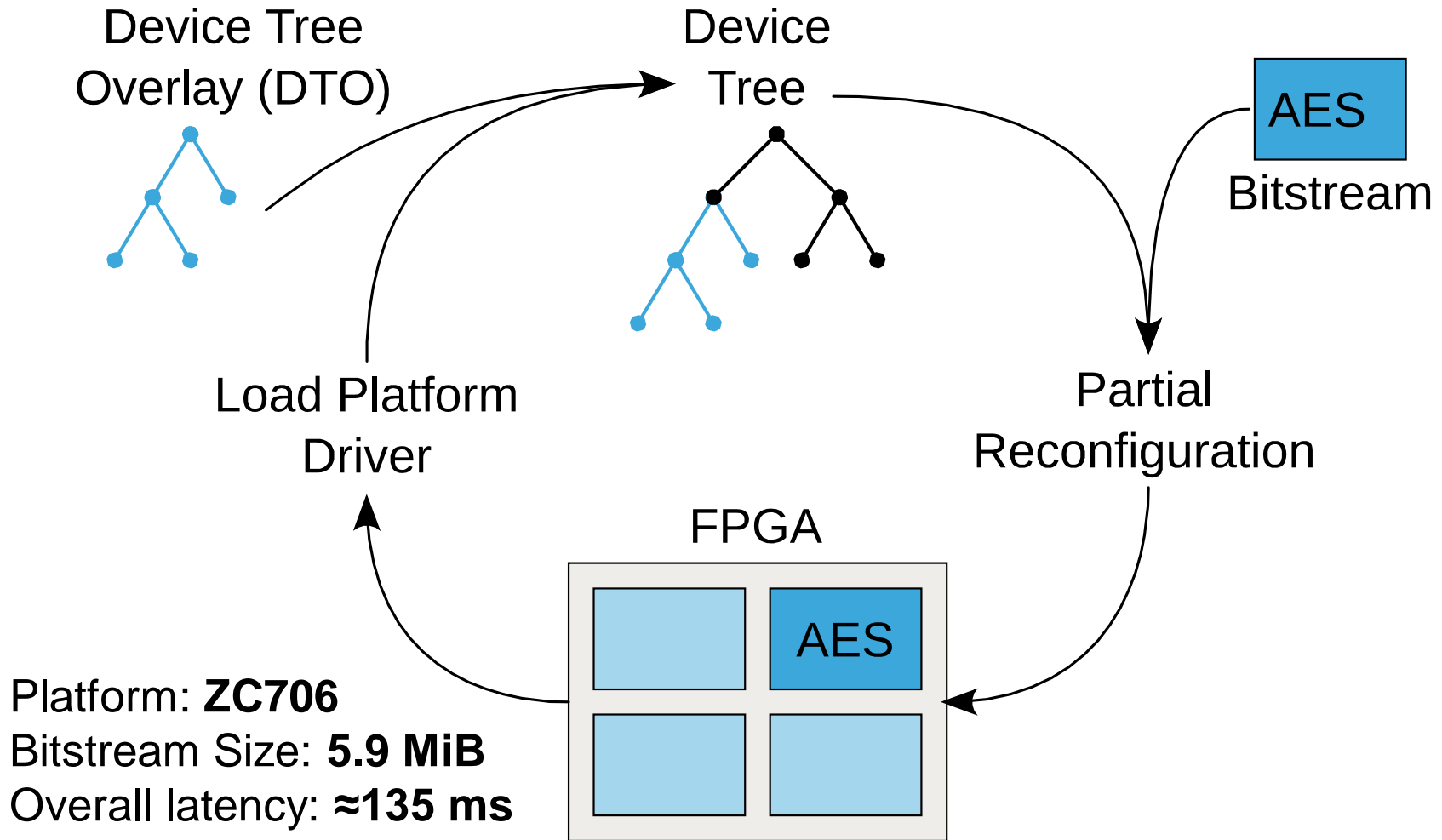# Software-Defined Services – Infrastructure Linux Kernel FPGA Framework

> **Supports both full and partial reconfiguration of FPGAs**

> **Adds a device tree interface for controlling the partial reconfiguration process**

> **Handles all FPGA internal processes**

> **Abstract device and vendor neutral interface**

© MLE

# Linux FPGA Framework Architecture

© MLE

# A Declarative Partial Reconfiguration Framework

Device Tree Overlay (DTO)

Device Tree

AES
Bitstream

Load Platform Driver

Partial Reconfiguration

FPGA

AES

Platform: **ZC706**
Bitstream Size: **5.9 MiB**
Overall latency: **≈135 ms**

# Conclusion & Outlook

© MLE

# Conclusion

> **Trend towards unconventional architectures**
- A diversification of increasingly heterogeneous devices and systems
- Convergence of networking, compute and storage within single nodes
- CPU-only processing runs out of steam

> **Key concepts for demonstrating Software-Defined Services**
- Offload engines for Linux Kernel Crypto-API
- Non-intrusive latency analysis via PCIe TLP "Tracers"
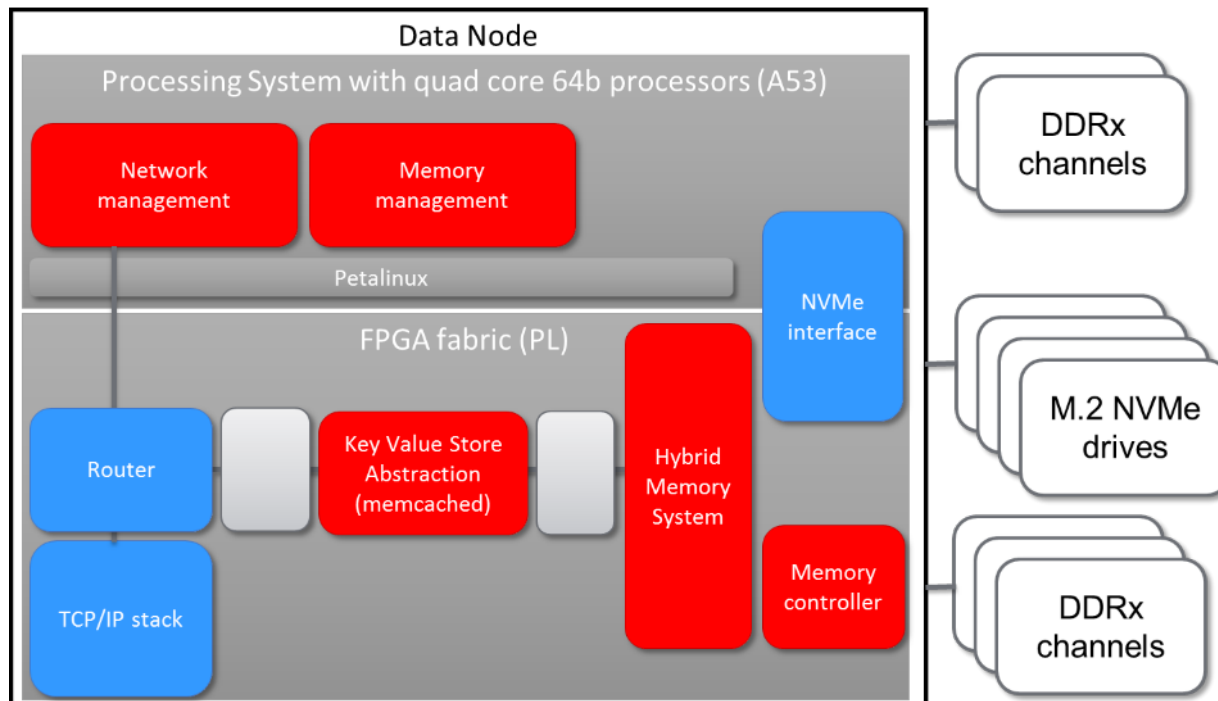- Inline processing with Deep Convolutional Neural Networks

> **Results:**
- On commercially available hardware
- Available for collaboration or in-house development
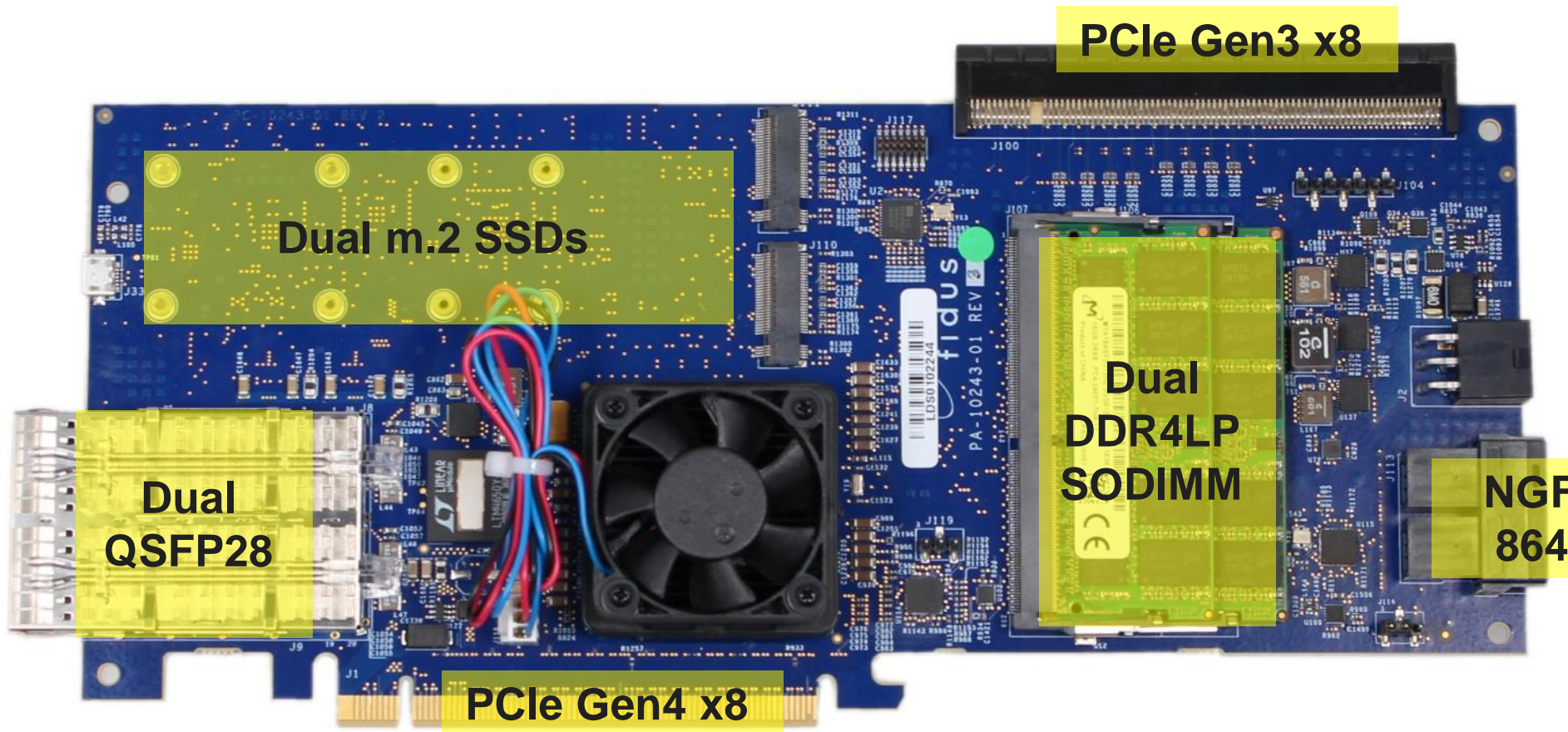
© MLE

# Single-Chip Implementation

> **Xilinx Zynq UltraScale+ MPSoC (XCZU19EG)**
- ARM Cortex A-53 quad-core, ARM Coretx R5 dual-core, 1,968 DSP slices
- 1.1 million system logic cells, 34Mbit BRAM, 36Mbit UltraRAM
- 5x PCIe Gen3/4, 4x 100GigE, 44x 16.3Gbps, 28x 32.72Gbps

© MLE

# Commercially Available Development System

> **Sidewinder-100 from Fidus Systems**

> **Accelerator IP and Linux BSP from MLE**



PCIe Gen3 x8

Dual m.2 SSDs

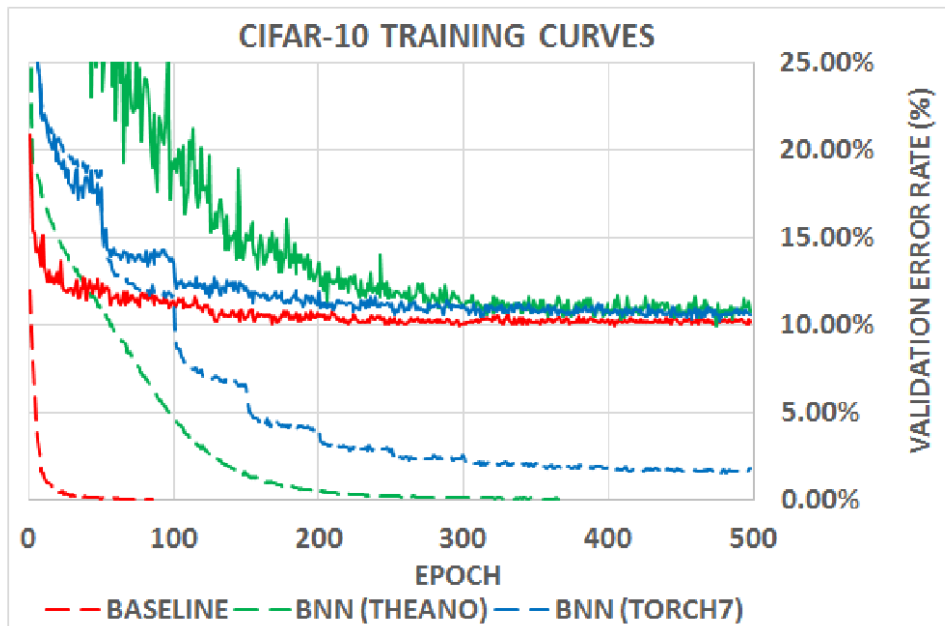Dual DDR4LP SODIMM

Dual QSFP28

NGF 864

PCIe Gen4 x8

© MLE

# Backup

# Reduced Precision Neural Networks

**➤ Binarized Neural Networks (BNN):**
**Training with float, CNN Inference runs at reduced precision**

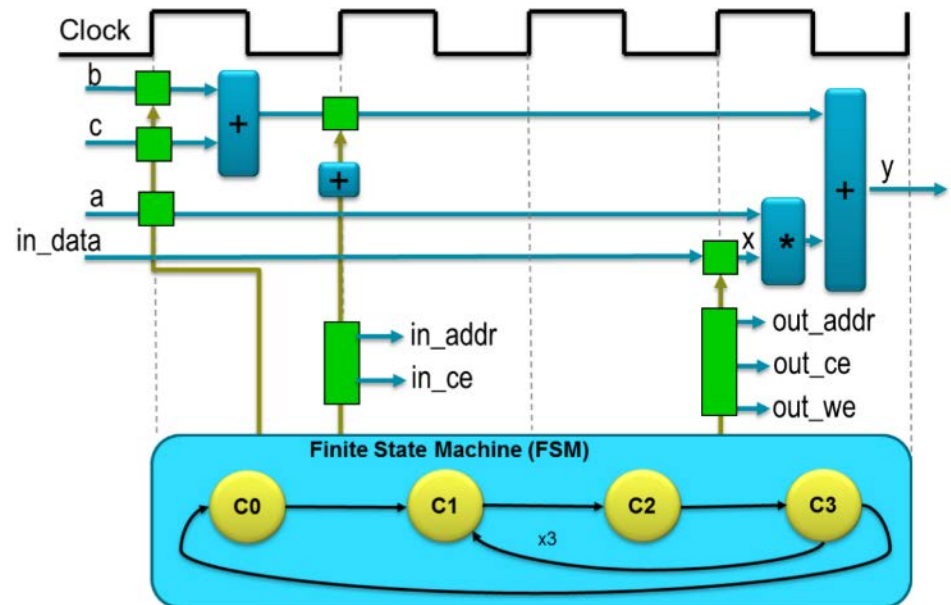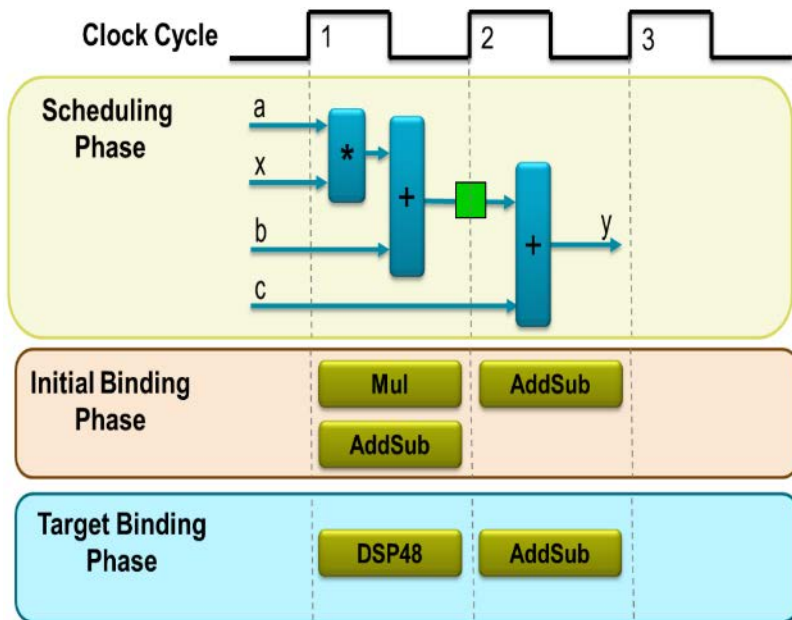– Less data (Mbytes) for parameters, less compute burdon.



CIFAR-10 TRAINING CURVES

© MLE

# Working Principles of High-Level Synthesis

❯ Design automation runs scheduling and resource binding to generate RTL code comprising data paths plus state machines for control flow
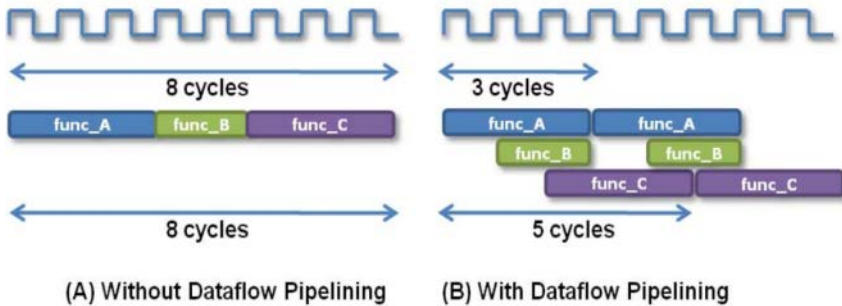
© MLE

# Benefits of HLS-Based C/C++ FPGA Design

> Automated performance optimizations via parallelization at dataflow level

> Automatic interface synthesis and driver code generation for HW/SW connectivity



```
void top (a,b,c,d) {
    ...
    func_A(a,b,i1);        func_A
    func_B(c,i1,i2);       func_B
    func_C(i2,d)           func_C

    return d;
}
```
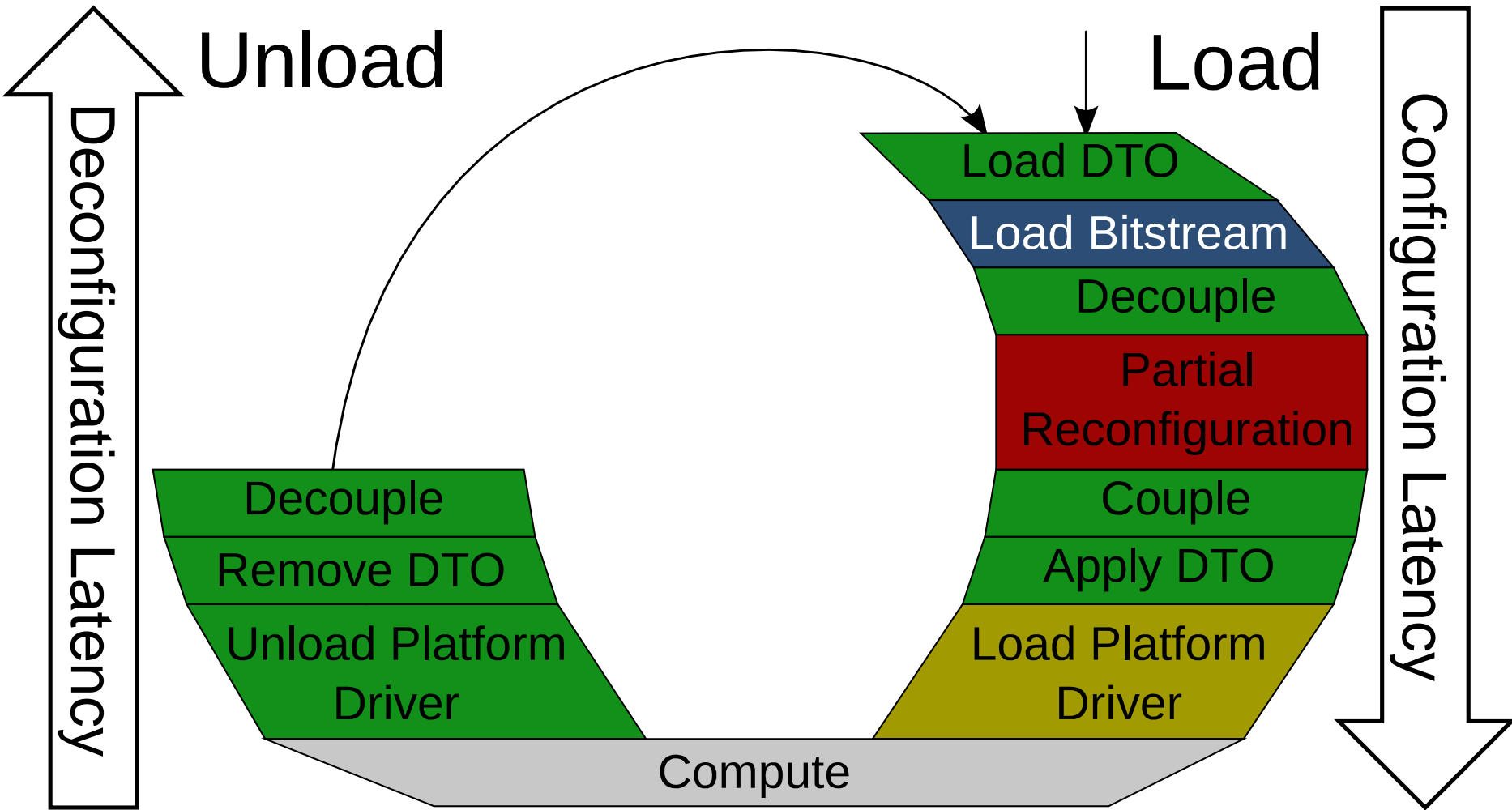
(A) Without Dataflow Pipelining

(B) With Dataflow Pipelining

© MLE
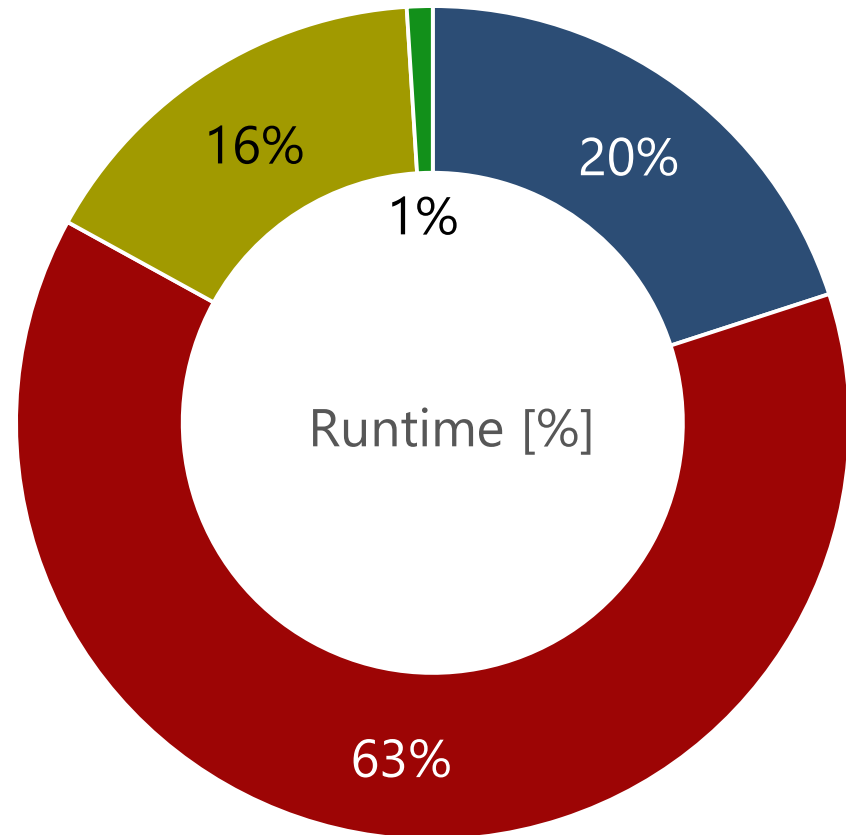
# Reconfiguration Performance

© MLE

# Scheduling Latency - Profiling Results

- **Measurement of example system (AES accelerator on ZC706 board)**
- **Measured latencies via *ftrace* function entry and exit timestamps**

- **Bitstream Size: 5.9 MiB**
- **Overall latency: ≈135 ms**

  - Load Bitstream

  - Partial Reconfiguration

  - Load Platform Driver

  - Rest incl. Framework



Runtime [%]

20%
1%
16%
63%

mle
missing link electronics

XILINX ALL PROGRAMMABLE.

# Contact

- **Endric Schubert**
  **Email: endric@mlecorp.com**


- **Ulrich Langenbach**
  **Email: ulrich@mlecorp.com**


- **Missing Link Electronics**
  **www.missinglinkelectronics.com**
  **Ph US: +1-408-475-1490**
  **Ph GER: +49-731-141149-0**

© MLE