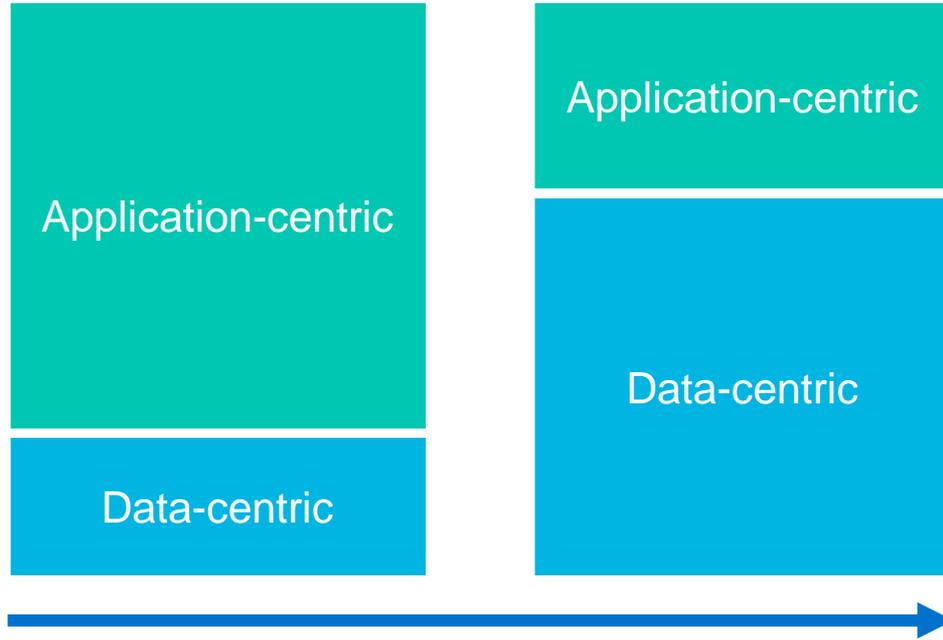# OCTEON® 10 and TVM

Dec 2021

Derek Chickles (dchickles@marvell.com)
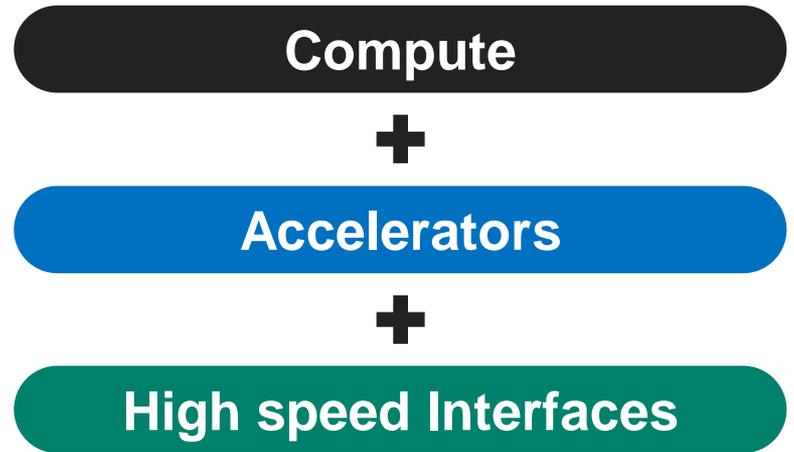
# Workloads are shifting to data-centric compute
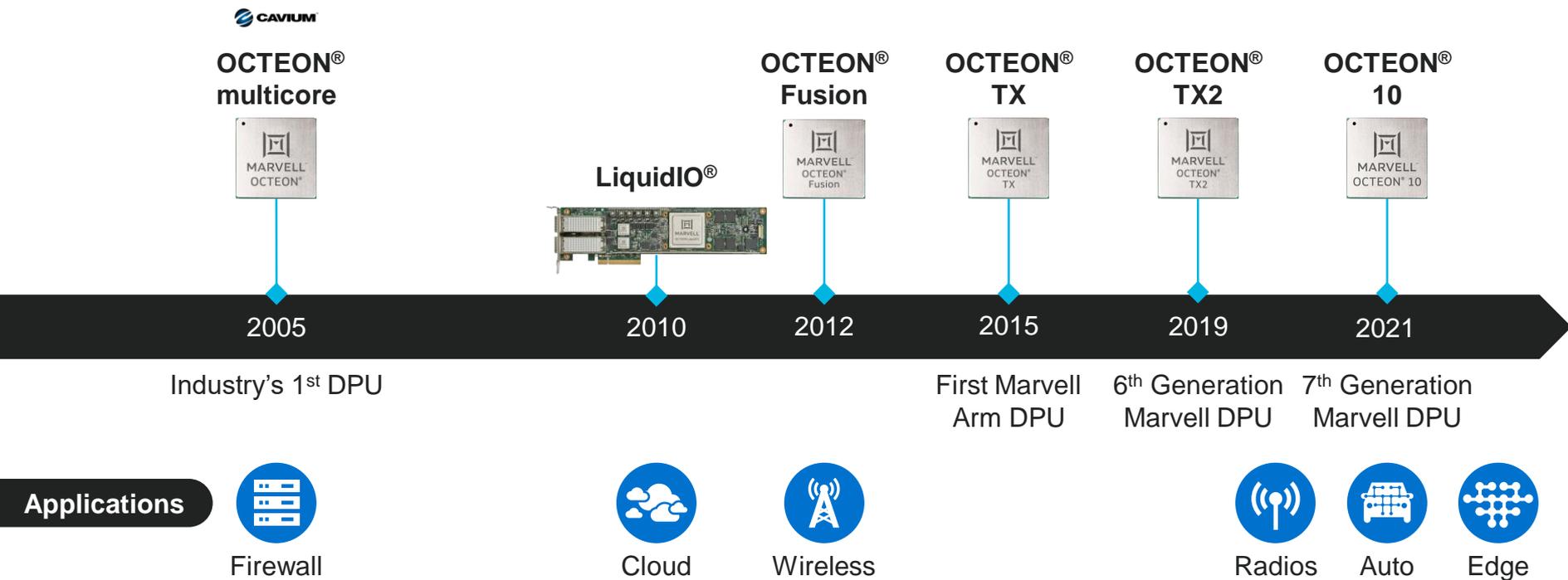


**AI, Networking, security, video and storage virtualization**

# DPU definition

A Data Processing Unit (DPU) is a compute entity that is used to move, process, secure and manage data, as it travels or while at rest, to make it available and optimized for application

**Compute**

**+**

**Accelerators**

**+**

**High speed Interfaces**

# OCTEON: The original DPU platform

CAVIUM

**OCTEON®**
**multicore**

**OCTEON®**
**Fusion**

**OCTEON®**
**TX**

**OCTEON®**
**TX2**

**OCTEON®**
**10**

**LiquidIO®**

2005       2010    2012      2015      2019      2021

Industry's 1st DPU

First Marvell
Arm DPU

6th Generation
Marvell DPU

7th Generation
Marvell DPU

**Applications**

Firewall      Cloud    Wireless      Radios    Auto    Edge

4

# OCTEON 10
## Industry Firsts

Compute leadership with Arm Neoverse N2 cores

Based on TSMC 5nm process

VPP hardware acceleration

Integrated hardware ML engine

Integrated 1terabit switch

Advanced inline crypto accelerators

**MARVELL OCTEON® 10**

**Compute leadership with industry-leading performance per Watt**

# OCTEON 10 innovations



DDR 5 Memory Controllers

Up to 400GE Ethernet

16 x 50G Ethernet Switch

Arm v9 N2 — 64K I / d cache — 1MB L2

Arm v9 N2 — 64K I / d cache — 1MB L2

L3 cache 2MB per core

Inline Crypto Processor

System Virtualization

Inline ML Processor
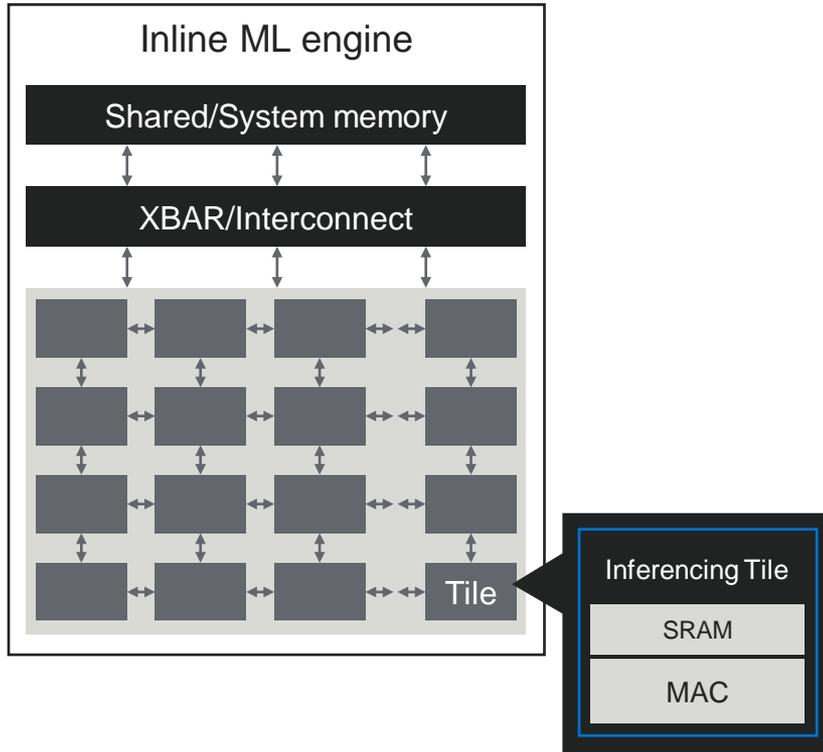
Vector Packet Processing

PCIe 5.0

- 5nm TSMC process
  - Enables fanless designs

- **First inline DPU ML Engine**

- Hardware VPP acceleration

- Inline crypto processor

- Arm Neoverse N2 cores
  - Highest SPECint in industry

- PCIe 5.0, DDR5 support

- Integrated with 16x 50GE switch

- 56G SerDes

# OCTEON DPU platform



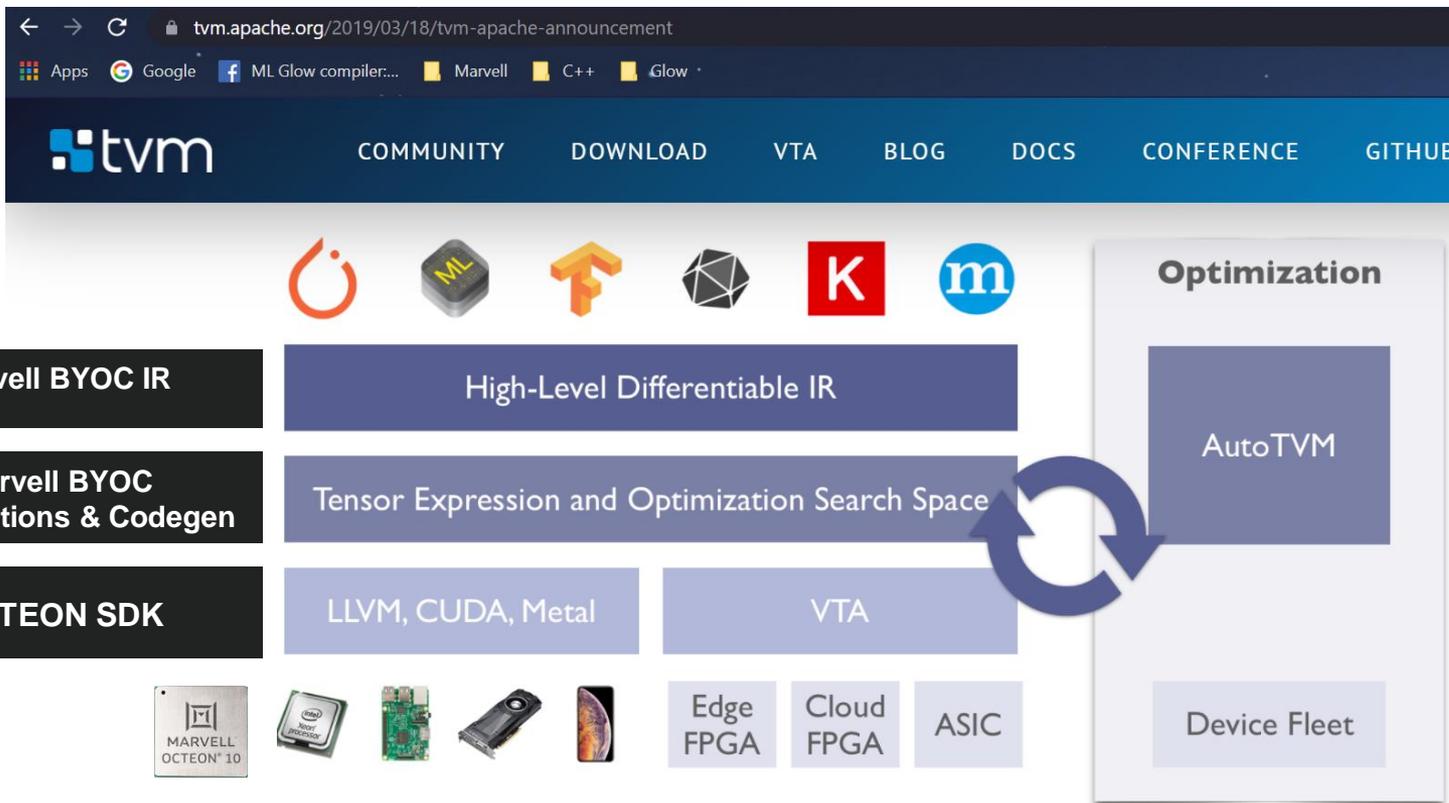| Software | **User Applications** | | |
|---|---|---|---|
| | **OCTEON DPU Open Software Platform** | | |
| | Optimized Stacks: Networking, Storage, Security | Virtualization and Containers | Standard APIs DPDK, SPDK, VPP |
| Silicon | **OCTEON DPU** | | |
| | Arm cores | Ethernet / PCIe / Memory Controllers | Software-enabled Accelerators |

# Integrated ML engine



- **Best-in-class DPU inferencing**
  - Directly in the data pipeline
  - Each ML tile contains private SRAM
  - Ultra Low Power

- **Up to 100x performance vs SW**
  - Supports INT8, FP16
  - Accelerated Tanh and Sigmoid activation functions

- **Use cases**
  - Threat detection
  - Context-aware service delivery
  - QoS
  - Beamforming optimization
  - Predictive maintenance

# TVM: Mrvl-BYOC with Marvell Specializations, Software, and Hardware

# Marvell TVM Integration

| | |
|---|---|
| **1** | Contributions have started (bug fixes) |

| | |
|---|---|
| **2** | Marvell BYOC up for Review: FP16 compile-time flow |

| | |
|---|---|
| **3** | Next year: Quantized INT8 flow, run-time flows and more |

# Introducing Joe!

Chien-Chun (Joe) Chou

- Principal ML Software Engineer at Marvell
  - Leading the TVM initiative to deliver end-to-end ML solution for Marvell OCTEON
  - His background:
    - ML compiler frontend/code-gen/backend development and optimization (with GLOW, TVM, ONNX, etc.)
    - ML and SoC hardware architecture
    - Edge and Automotive ML
- Marvell contact: cchou1@marvell.com
- github contact: https://github.com/ccjoechou

# Thank You