**XILINX**®

# Versal™ AI Edge Series Announcement

Rehan Tahir, Senior Product Line Manager

# What's Happening at the Edge

## The Edge

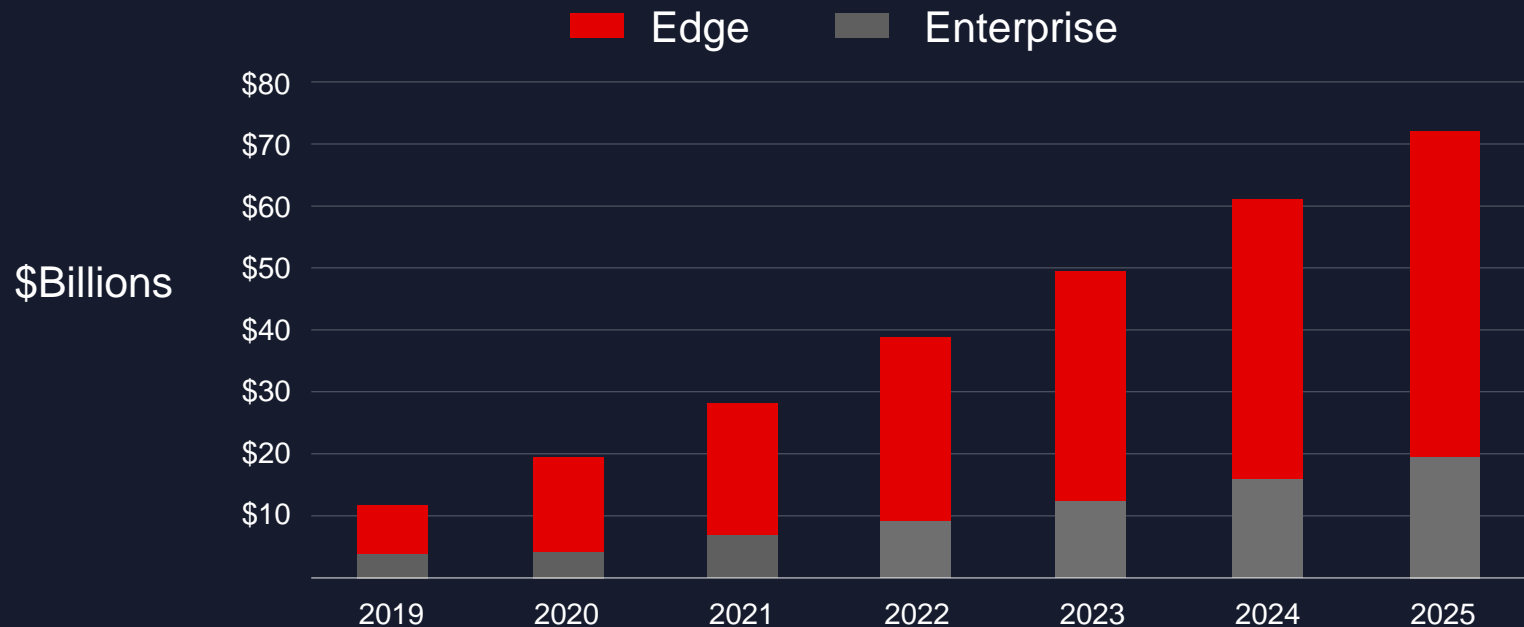| Low Latency | AI Compute | Low Power | Safety and Security |
|---|---|---|---|

XILINX

# Hypergrowth at the Edge

"Edge computing … solves for weaknesses of the cloud"[1]

Edge AI chipset opportunity is 3X that of data center - $65B in 2025[2]

Deep learning chipset revenue, enterprise vs. edge, world markets: 2019–25

■ Edge ■ Enterprise

$Billions

| Year | Value |
|------|-------|
| 2019 | |
| 2020 | |
| 2021 | |
| 2022 | |
| 2023 | |
| 2024 | |
| 2025 | |

XILINX

# Now Bringing Versal ACAPs to the Edge

▸ Versal™ ACAPs first introduced breakthrough compute for the cloud and network

▸ Now 'miniaturizing' this technology for performance/watt at the edge



CLOUD

NETWORK

EDGE

# New Versal™ Platform for Intelligence at the Edge



**Smart Vision**

**Unmanned Aerial Vehicles**

**Collaborative Robotics**

**ADAS & Automated Drive**

**Endoscopy**

**Ultrasound**

XILINX®
VERSAL™

AI RF
Series

HBM
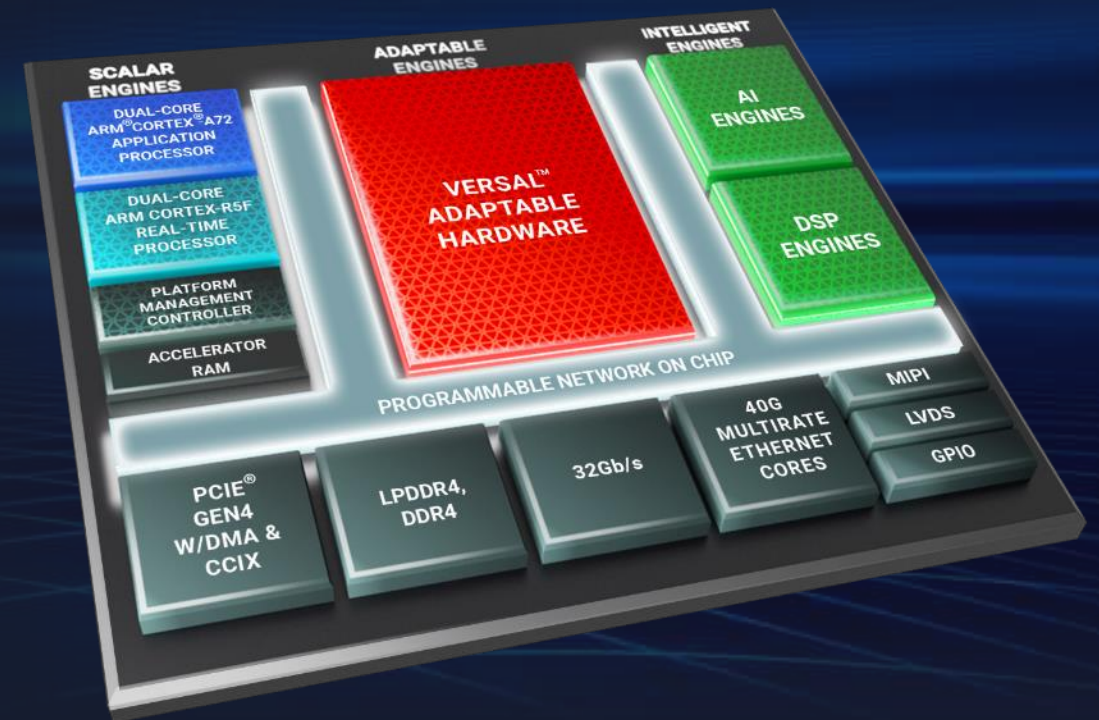Series

AI Core
Series

Premium
Series

AI Edge
Series

Prime
Series

XILINX®

# Versal™ AI Edge: Intelligence Unleashed



- ▶ 4X AI Performance/Watt vs. GPUs[1] with Innovations in AI Engines and Memory Hierarchy

- ▶ 10X Compute Density[2] with Highest Levels of Safety and Security

- ▶ World's Most Scalable and Adaptable Platform for Edge and Endpoint

1: vs. Jetson AGX Xavier, ResNet50 224x224, batch=1, https://developer.nvidia.com/embedded/jetson-agx-xavier-dl-inference-benchmarks
2: Compared to Zynq® UltraScale+™ MPSoCs

# 4X AI Performance/Watt

# Proven AI Engine Architecture

## Array of Compute Core

▸ Flexible compute: fixed- & floating-point vector processors
▸ HW adaptable to evolving algorithms

## Tightly Coupled Memory

▸ Cache-less memory hierarchy
▸ Maximizes bandwidth, ensures determinism & low latency

## Flexible Interconnect

▸ Connect any tile to any tile for custom microarchitecture
▸ High bandwidth

AI Engines Array
(Part of ACAP Device)

AI Engine
Tile

Distributed
Data Memory

Interconnect

**Architected for Adaptability, Low Power, and Low Latency**

XILINX

# Optimized AI Engines-ML for Machine Learning

## Optimized the compute core for ML

▸ Doubled the multipliers, doubled INT8 performance

▸ Native support for INT4 and BFLOAT16

## Doubled the data memory

▸ From 32kB to 64 kB

▸ Improved localization of data

## New Memory Tile

▸ Up to 38 Megabytes across the AI engine array

▸ Higher bandwidth memory access

Optimized
AI Engine-ML Array
(Part of ACAP Device)

Optimized
Compute Core

2X Data
Memory

New
Memory Tile

Memory Tile

## Delivering 4X ML Compute at ½ the Latency[1]

1: AI Engine-ML delivers 2X INT8 compute, 4X INT4 compute, and 16X BFLOAT16 compute vs. AI Engine (per core)
2: Native 32-bit support in AI Engines only

XILINX

# AIE-ML Complements AI Engines for Diverse Workloads

AI Engine
Architecture

Balanced for ML & DSP

▶

4X ML Compute
at ½ the Latency

▶

AI Engine-ML
Architecture

Optimized for ML

← Advanced Signal Processing

Machine Learning →

BEAMFORMING,
RADAR PROCESSING,
HIGH PERF. COMPUTING

AI Engine[2]

AI Engine-ML[1]

CNN, RNN, MLP

1: AI Engine-ML delivers 2X INT8 compute, 4X INT4 compute, and 16X BFLOAT16 compute vs. AI Engine (per core)
2: Native 32-bit support in AI Engines only

XILINX

# Innovations in Memory Hierarchy: Accelerator RAM

## 4MB of On-Chip RAM for Massive Bandwidth

Avoid DDR to store AI compute data or safety-critical code
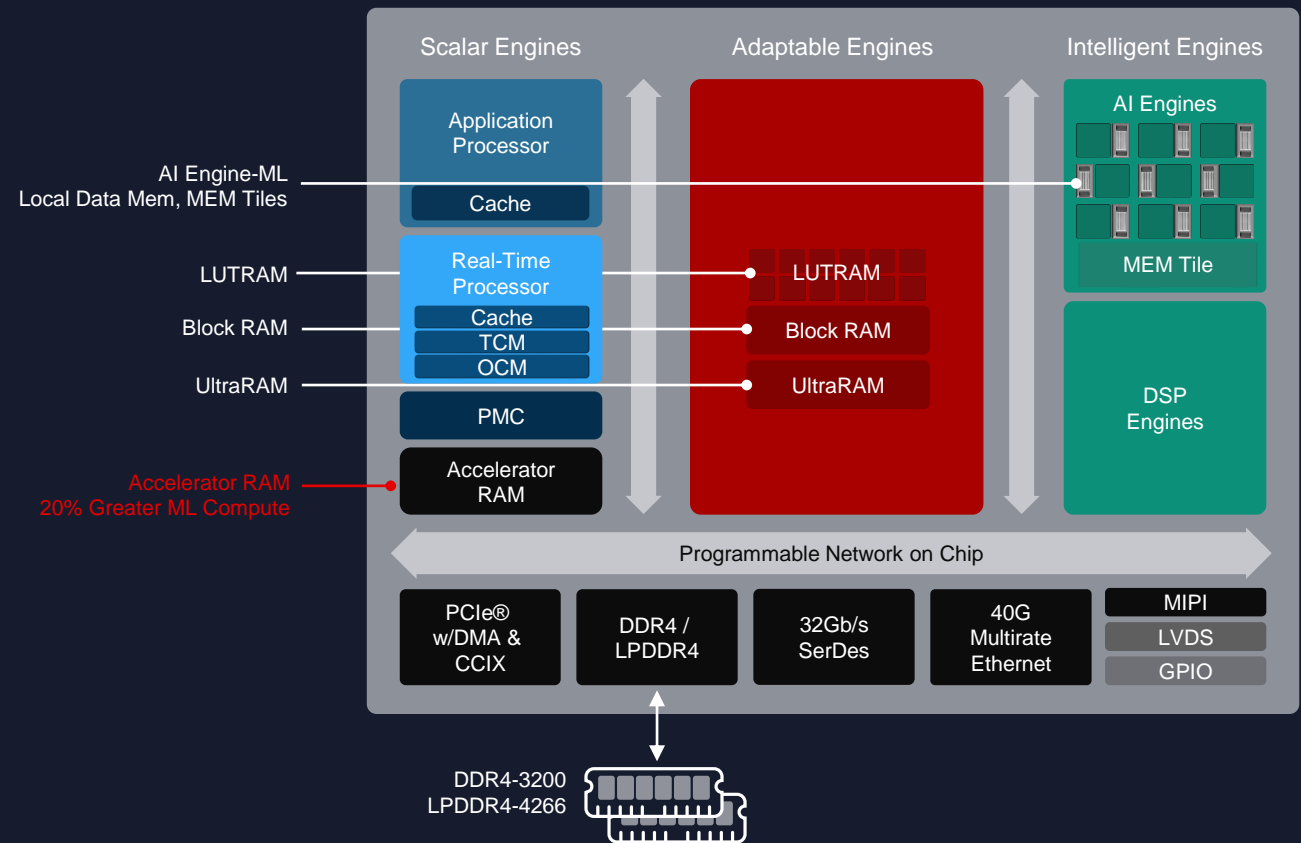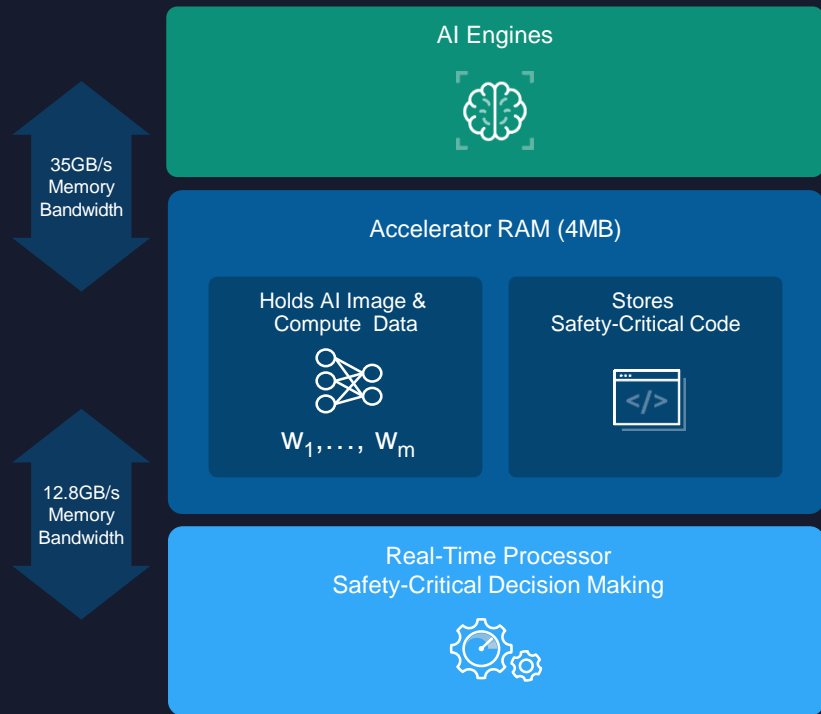
**AI Engines**

35GB/s Memory Bandwidth

**Accelerator RAM (4MB)**

Holds AI Image & Compute Data

$w_1, \ldots, w_m$

Stores Safety-Critical Code

12.8GB/s Memory Bandwidth

**Real-Time Processor**
Safety-Critical Decision Making

## Part of the Adaptable Memory Hierarchy

Select the right memory for bandwidth requirement

Scalar Engines

Adaptable Engines

Intelligent Engines

Application Processor

Cache

AI Engine-ML
Local Data Mem, MEM Tiles

AI Engines

MEM Tile

LUTRAM

Real-Time Processor

LUTRAM

Block RAM

Cache
TCM
OCM

Block RAM

UltraRAM

UltraRAM

PMC

DSP Engines

Accelerator RAM
20% Greater ML Compute

Accelerator RAM

Programmable Network on Chip

PCIe® w/DMA & CCIX

DDR4 / LPDDR4

32Gb/s SerDes

40G Multirate Ethernet

MIPI

LVDS

GPIO

DDR4-3200
LPDDR4-4266

XILINX

# Up to 4X Performance/Watt vs. GPUs

| Intelligent Edge Sensor | Autonomous System or Edge Aggregation | CPU Accelerator |
|---|---|---|

**1.9X**
Performance/Watt

**3.3X**
Performance/Watt

**4.2X**
Performance/Watt

Images/sec/watt

Jetson Xavier NX [1]
Versal AI Edge (VE2102)

Jetson AGX Xavier (15W Mode) [2]
Versal AI Edge (VE2302)

Jetson AGX Xavier (MAX N-Mode) [3]
Versal AI Edge (VE2802)

ResNet50 224x224

ResNet50 224x224, batch=1

ResNet50 224x224, batch=1

1: Jetson NX Xavier: https://mlcommons.org/en/inference-edge-10, batch size not provided
2: Jetson AGX Xavier run in a mid-performance & power configuration, categorized as "15 W-Mode": https://developer.nvidia.com/embedded/jetson-agx-xavier-dl-inference-benchmarks
3: Jetson AGX Xavier MAX N-Mode and Versal™ VE2802 ACAP represent the highest performing device configuration in their respective portfolios
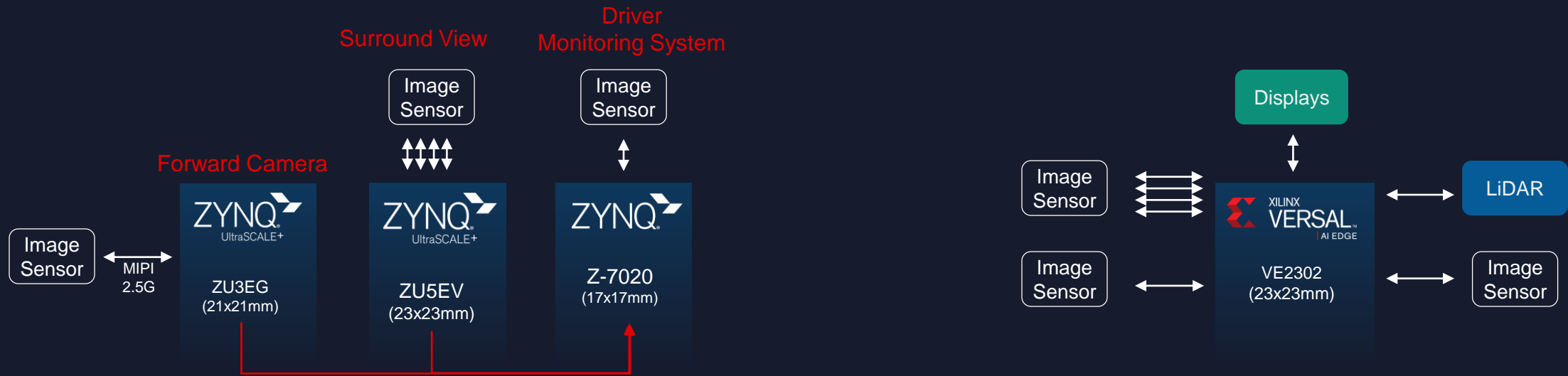Jetson Xavier device power estimated by subtracting published memory & I/O power from total module power
All charts are normalized

XILINX

# 10X Compute Density with Highest Levels of Safety and Security
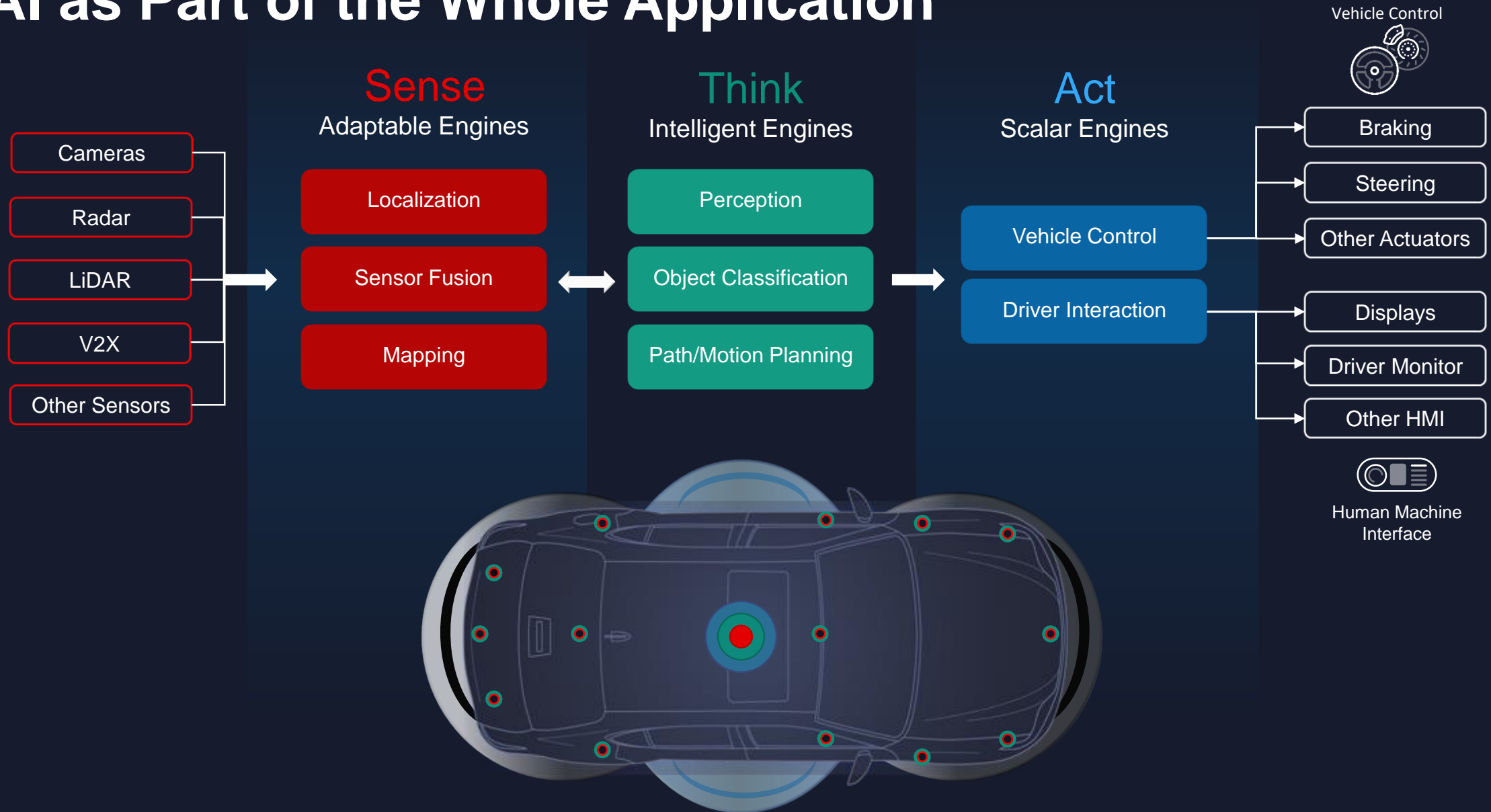
XILINX

# 10X Compute Density: Level 3 Semi-Automated Driving

Surround View

Driver Monitoring System

Image Sensor

Image Sensor

Forward Camera

Image Sensor

MIPI 2.5G

**ZYNQ** UltraSCALE+
ZU3EG (21x21mm)

**ZYNQ** UltraSCALE+
ZU5EV (23x23mm)

**ZYNQ**
Z-7020 (17x17mm)

Image Sensor

Image Sensor

Displays

**XILINX VERSAL** AI EDGE
VE2302 (23x23mm)

LiDAR

Image Sensor

| | Previous Gen Adaptive SoC | | Versal™ ACAP |
|---|---|---|---|
| Compute | 6x cameras (2MP, 4MP) + AI (4TOPs) | ▶ 4.4X ▶ | 6x cameras (2MP, 8MP) + AI (17.4TOPs) |
| Area | 3 devices = 1,259mm$^2$ | ▶ 58% Less ▶ | 1 device = 529mm$^2$ |
| Power | ZU3(6W) + ZU5(10W) + Z-7020(5W) | ▶ ~1X* ▶ | ~20W |

*Power levels are typical, approximate, and estimated at room temp

**XILINX**

# AI as Part of the Whole Application



Sense — Adaptable Engines: Localization, Sensor Fusion, Mapping

Think — Intelligent Engines: Perception, Object Classification, Path/Motion Planning

Act — Scalar Engines: Vehicle Control, Driver Interaction

Inputs: Cameras, Radar, LiDAR, V2X, Other Sensors

Vehicle Control outputs: Braking, Steering, Other Actuators

Driver Interaction outputs: Displays, Driver Monitor, Other HMI

Human Machine Interface

XILINX

# Whole Application Acceleration for Real-Time Systems
## From Sensor to AI to Real-Time Control

EXECUTION TIME

| Sense | Think (AI) | Act |
|-------|-----------|-----|

Sense

Think

Act

Actuator
Response

© Copyright 2021 Xilinx

XILINX
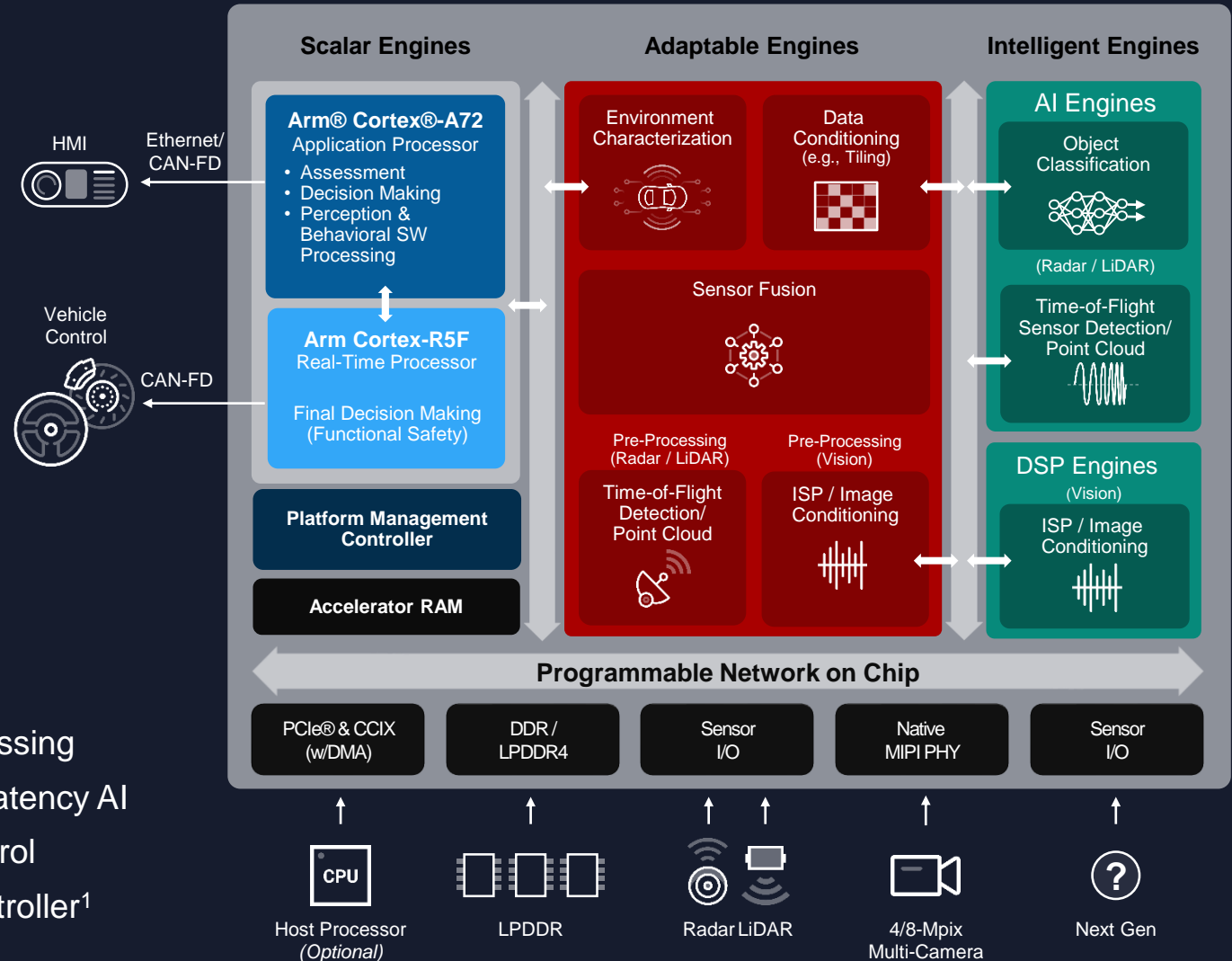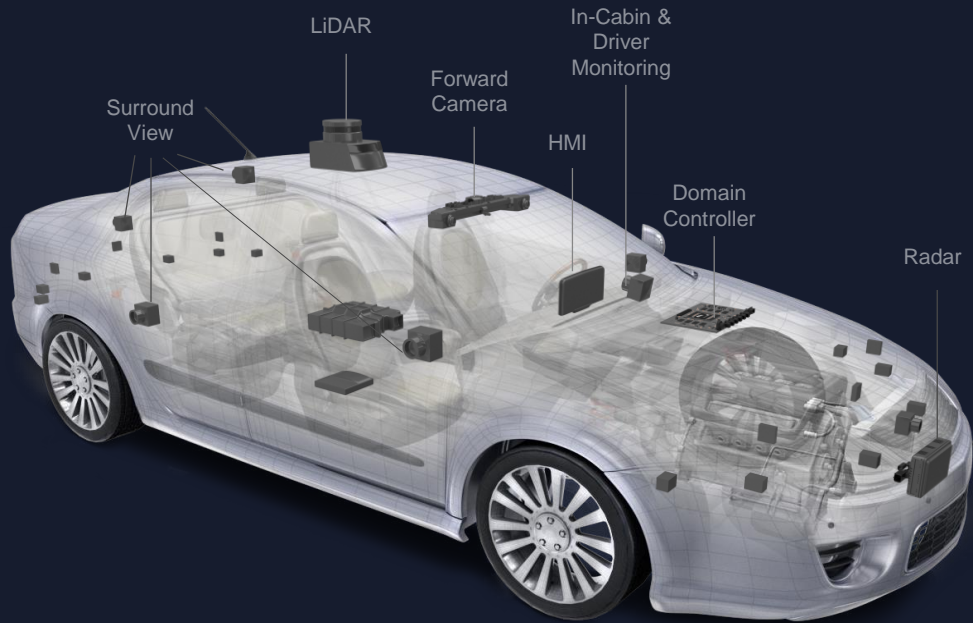
# Versal™ AI Edge ACAP in ADAS and Automated Driving



## Accelerating the Whole Application from Sensor to AI to Real-Time Control

▸ Adaptable Engines for sensor fusion and pre-processing

▸ Intelligent Engines for signal conditioning and low-latency AI

▸ Scalar Engine for decision making and vehicle control

▸ Scalable compute from edge sensor to domain controller[1]

1: Diagram demonstrates capabilities of architecture; does not represent a single chip AD system

© Copyright 2021 Xilinx

XILINX

# Fully Automotive-Qualified and Safety Certified



Architected to Meet Stringent ISO 26262 Requirements

XILINX

# Supporting Multiple Safety Standards

## Versal™ AI Edge ACAP

### ISO 26262
Automotive

### IEC 61508
Safety across
All Industries

### DO-254/178
Avionics HW/SW

**IEC 61511**
Process Industry

**IEC 61800**
Electrical Drives

**ISO 13849**
Machinery Control

**IEC 62061**
Machinery

**EN 60601**
Medical

**XILINX**

# Collaborative Robotics: AI-Based Systems Need to be Safe and Secure

**Real-Time Precision and Control to Augment AI**
Deterministic response, AI to navigate
unpredictable movement of workers

**Environmental Awareness and Perception**
Sensor fusion for perception, self-learning
to improve capabilities over time

**Predictive Maintenance**
Analyze sensor data for actionable insights
to reduce downtime

**Safety and Security are Connected Matters**
Cyber-attack creates safety and data privacy risks,
robotic systems require IEC 62443 compliance

Human Safety

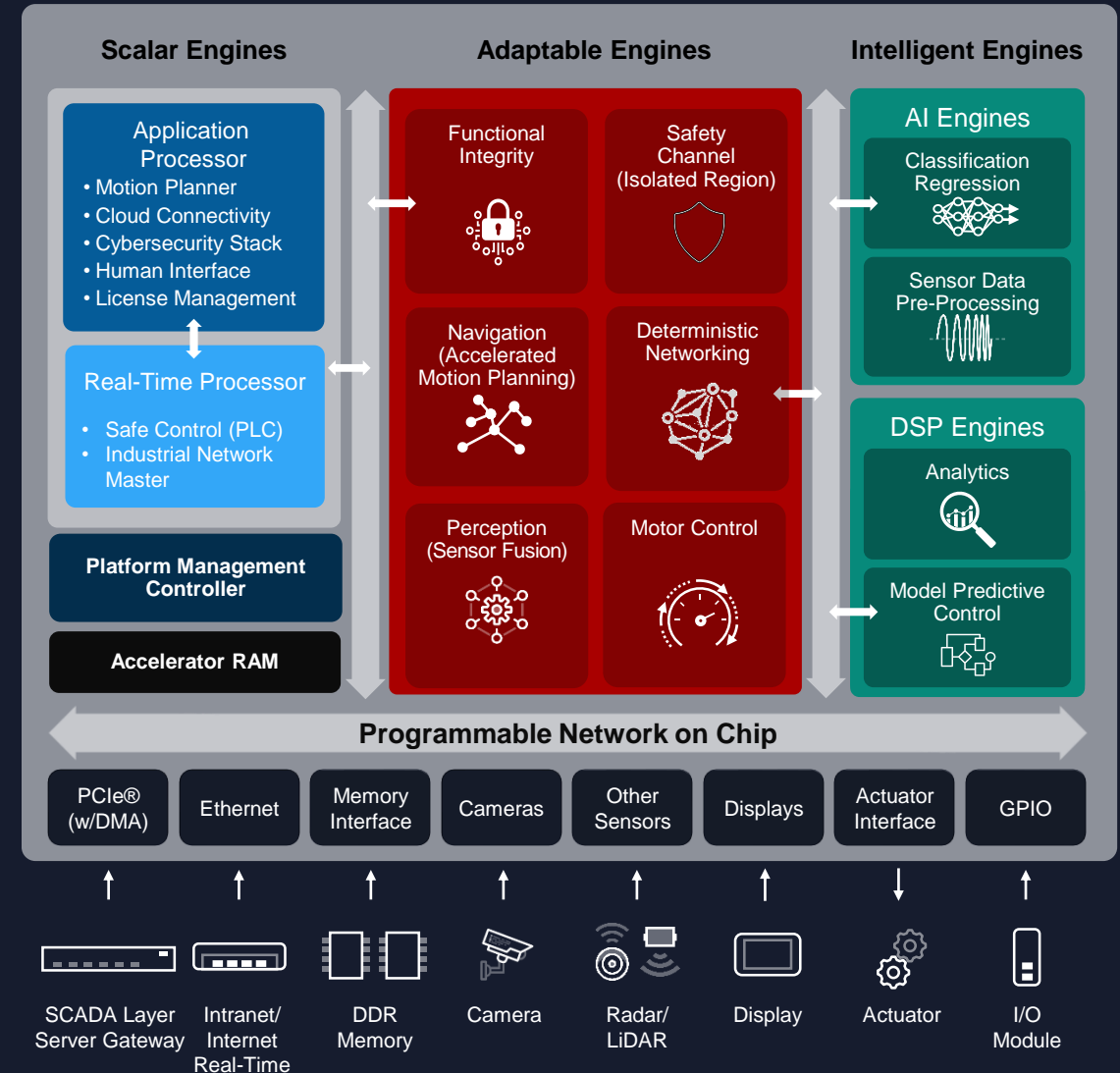Real-Time
Control

Sensor
Data

XILINX

# Whole Application Acceleration for Collaborative Robotics

**Robotic Perception Systems for Real-Time Control, Safety Critical, and Predictive Maintenance**

▶ Adaptable Engines for perception, control/networking, navigation

▶ AI to augment control for dynamic execution, predictive maintenance

▶ Scalar Engines for cybersecurity (IEC 62443), safety control, UI



**Scalar Engines**

Application Processor
- Motion Planner
- Cloud Connectivity
- Cybersecurity Stack
- Human Interface
- License Management

Real-Time Processor
- Safe Control (PLC)
- Industrial Network Master

Platform Management Controller

Accelerator RAM

**Adaptable Engines**

Functional Integrity

Safety Channel (Isolated Region)

Navigation (Accelerated Motion Planning)

Deterministic Networking

Perception (Sensor Fusion)

Motor Control

**Intelligent Engines**

AI Engines
- Classification Regression
- Sensor Data Pre-Processing

DSP Engines
- Analytics
- Model Predictive Control

Programmable Network on Chip

| PCIe® (w/DMA) | Ethernet | Memory Interface | Cameras | Other Sensors | Displays | Actuator Interface | GPIO |

| SCADA Layer Server Gateway | Intranet/ Internet Real-Time | DDR Memory | Camera | Radar/ LiDAR | Display | Actuator | I/O Module |

22

© Copyright 2021 Xilinx

XILINX

# AI-Enabled Multi-Mission Payloads for UAVs

*AI with Software Defined Radio (SDR), Signal Intelligence (SIGINT), Image/Video Processing*

## Vision AI for Real-Time Analysis and Response

Autonomous flight control, optimize navigation paths
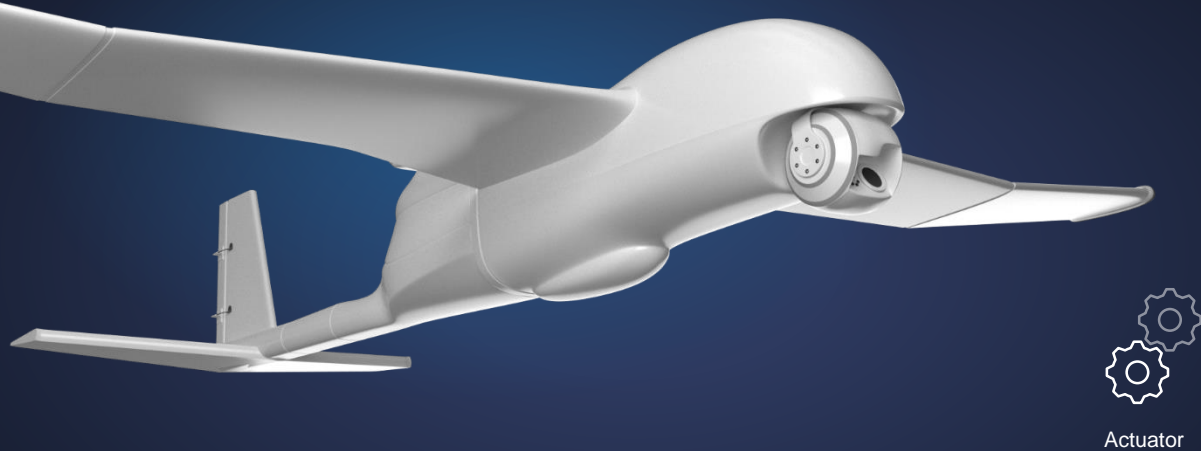
## Cognitive RF

Optimizing radio communication and protecting against malicious intrusion

## Diverse and Emerging Forms of AI

AI is rapidly evolving in tactical applications and vendors will need to adapt over time
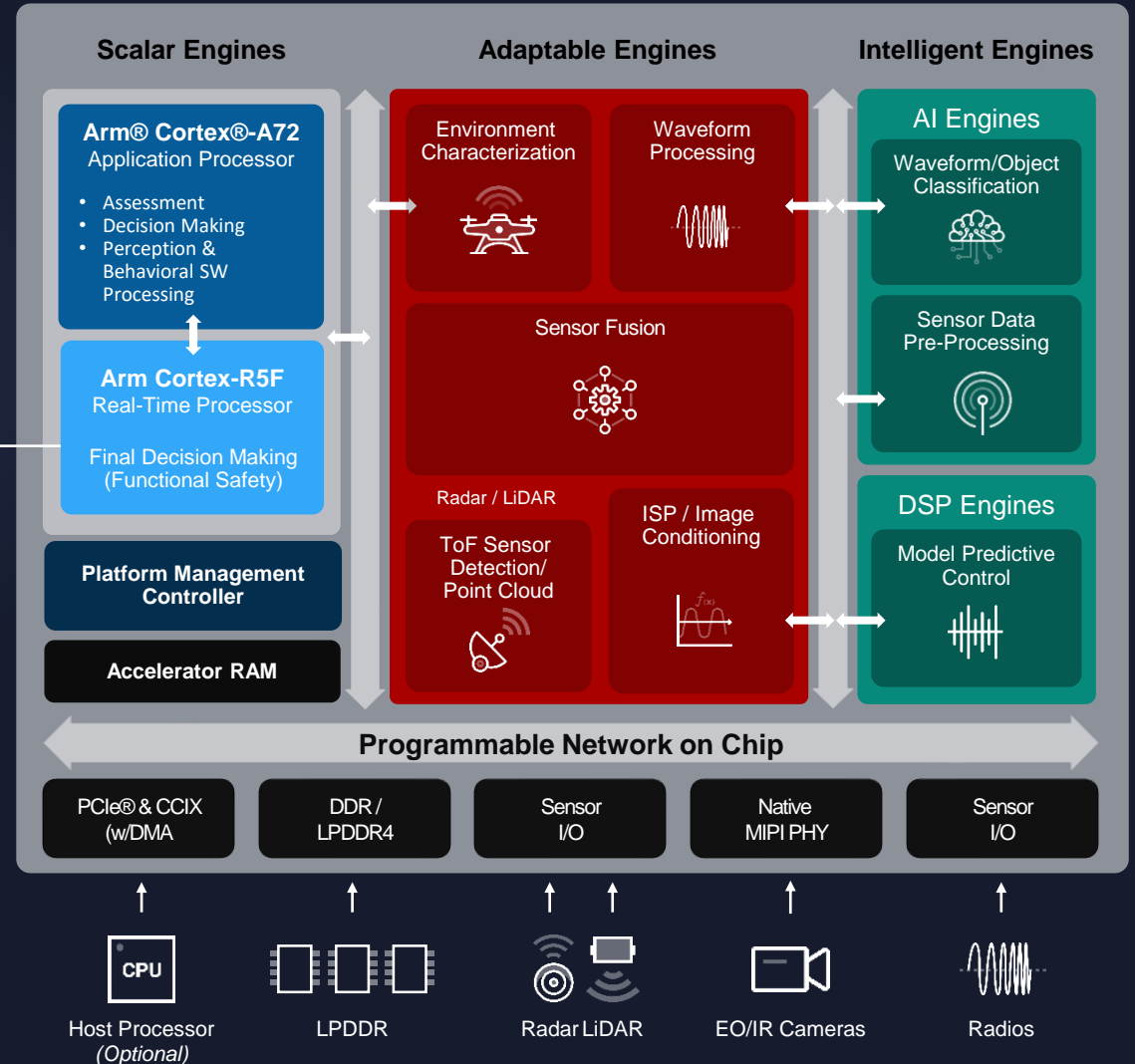
**Need AI Compute in Limited Size, Weight, and Power (SWaP) and Thermal Envelope**

XILINX

# Versal AI Edge for Unmanned Aerial Vehicles

## Multi-Mission and Situationally Aware UAVs with Low SWaP

▸ Adaptable Engines for sensor fusion and pre-processing

▸ Intelligent Engines for low power, low latency AI and signal conditioning

▸ Scalar Engines for command and control

▸ Ruggedized packaging and military-temp grade (XQ)

Actuator

**Scalar Engines**

**Arm® Cortex®-A72**
Application Processor

- Assessment
- Decision Making
- Perception & Behavioral SW Processing

**Arm Cortex-R5F**
Real-Time Processor

Final Decision Making
(Functional Safety)

**Platform Management Controller**

**Accelerator RAM**

**Adaptable Engines**

Environment Characterization

Waveform Processing

Sensor Fusion

Radar / LiDAR

ToF Sensor Detection/ Point Cloud

ISP / Image Conditioning

**Intelligent Engines**

AI Engines

Waveform/Object Classification

Sensor Data Pre-Processing

DSP Engines

Model Predictive Control

**Programmable Network on Chip**

| PCIe® & CCIX (w/DMA | DDR / LPDDR4 | Sensor I/O | Native MIPI PHY | Sensor I/O |

Host Processor *(Optional)*    LPDDR    Radar LiDAR    EO/IR Cameras    Radios

CPU

XILINX.

# Versal ACAP Development Experience for All Developers



HW Developer

SW Developer

Data Scientist

VIVADO™

XILINX VITIS™

C, C++, Python

XILINX VITIS™ | AI

TensorFlow    PyTorch    Caffe

Spark    FFmpeg    mxnet

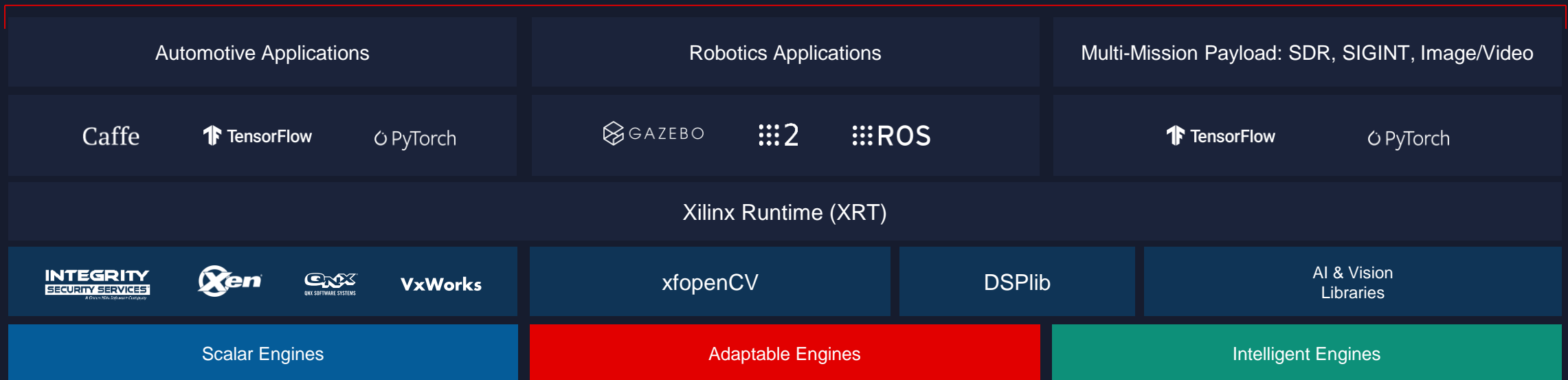| OS & Embedded Run-Time | Custom HW | HW IP & Accelerated Libraries | HW Accelerated Libraries |
|---|---|---|---|
| Scalar Engines | Adaptable Engines | | Intelligent Engines |

Versal™ AI Edge ACAP

XILINX

# Market-Specific Application Stacks
## Examples for Automotive, Robotics, and Multi-Mission Payload Applications

▸ One platform with market-specific libraries, frameworks, and ecosystem to enable all developers

▸ Following industry standards for developing safety critical software on silicon

| Automotive Applications | Robotics Applications | Multi-Mission Payload: SDR, SIGINT, Image/Video |
|---|---|---|
| Caffe · TensorFlow · PyTorch | GAZEBO · 2 · ROS | TensorFlow · PyTorch |

**Xilinx Runtime (XRT)**

| INTEGRITY SECURITY SERVICES · Xen · QNX SOFTWARE SYSTEMS · VxWorks | xfopenCV | DSPlib | AI & Vision Libraries |
|---|---|---|---|

| Scalar Engines | Adaptable Engines | Intelligent Engines |
|---|---|---|

**Versal™ AI Edge ACAP**

XILINX®

# World's Most Adaptable and Scalable Edge Platform

XILINX

# Adaptability: From Domain Specific Architectures (DSAs) to Dynamic Function Exchange

## DSAs for Diverse Platform Requirements
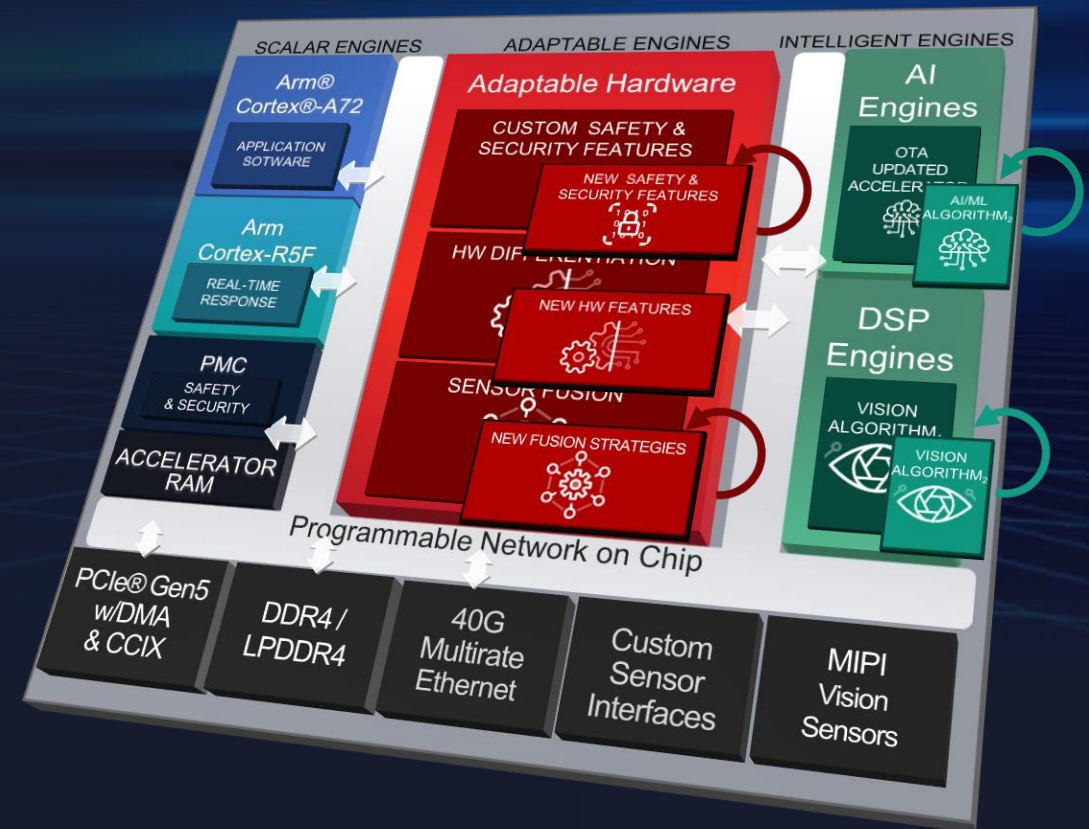
- Implement custom AI, vision, sensor strategies
- Design for different safety and security targets
- One platform for diverse end-customers' requirements

## Hardware/Software Over-the-Air Updates

- Update your AI accelerator or fusion algorithms
- Future proof for emerging security threats
- Avoid recalls or costly re-deployment

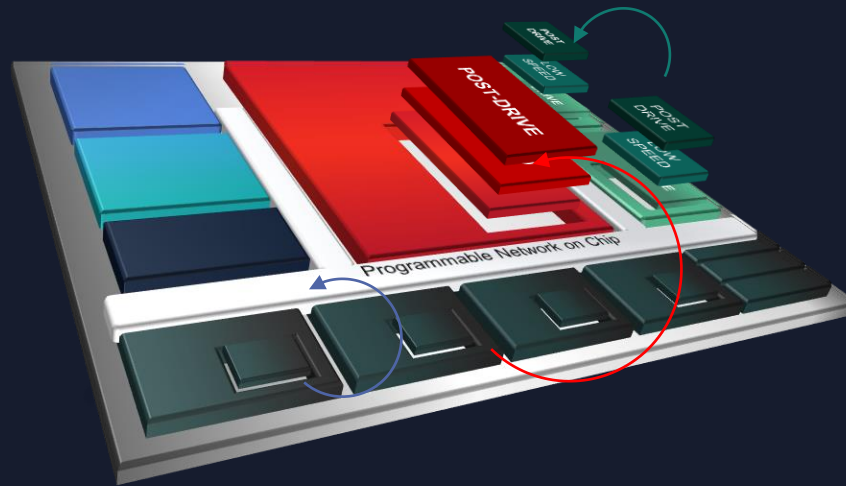## Dynamic Function Exchange (DFx)

- Swap functionality in milliseconds
- Available in Adaptable Engines, DSP, AI Engines
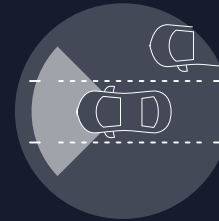- Fewer system components → reduce power and cost

XILINX.

# Dynamic Function Exchange (DFx) in Automotive
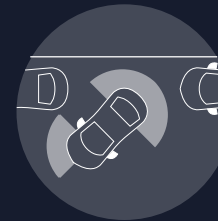
Swap Functionality in
Milliseconds

Programmable Network on Chip

Dynamic Regions
(Engines, Integrated Cores, I/O)

**Drive Mode**
(Lane Departure Warning)

**Low Speed Mode**
(Parking Assist)

**Post-Drive Mode**
(Dog Left Behind)

Fewer Devices to Reduce System-Wide Power and Cost

XILINX

# Scale from Edge Sensor to CPU Accelerator

Accelerator

Edge Aggregation &
Autonomous Systems

Intelligent Edge Sensor
& End Point

|  |  | VE2002 | VE2102 | VE2202 | VE2302 | VE2602 | VE1752 | VE2802 |
|---|---|---|---|---|---|---|---|---|
| **Engines** | Total AI Compute (INT4) | 14 TOPS | 22 TOPS | 47 TOPS | 67 TOPS | 256 TOPS | 166 TOPS | 479 TOPS |
|  | Total AI Compute (INT8) | 7 TOPS | 10 TOPS | 21 TOPS | 31 TOPS | 120 TOPS | 124 TOPS | 228 TOPS |
|  | AIE / AIE-ML[1] | 8 | 12 | 24 | 34 | 152 | 304 | 304 |
|  | Adaptable Engines | 20K LUTs | 37K LUTs | 105K LUTs | 150K LUTs | 375K LUTs | 449K LUTs | 521K LUTs |
|  | Processing Subsystem | Dual-Core Arm® Cortex®-A72 Application Processing Unit  /  Dual-Core Arm Cortex-R5F Real-Time Processing Unit | | | | | | |
| **RAM** | Accelerator RAM (4MB) | ✓ | ✓ | ✓ | ✓ | - | - | - |
|  | Total Memory | 95Mb | 103Mb | 156Mb | 172Mb | 554Mb | 253Mb | 575Mb |
|  | 32G Transceivers | - | - | 8 | 8 | 32 | 44 | 32 |
|  | PCIe® | - | - | ✓ | ✓ | - | - | - |
|  | PCIe + CCIX | - | - | - | - | ✓ | ✓ | ✓ |
|  | Estimated Power | 6–9W | 7–10W | 15–20W | ~20W | 50–60W | 50–60W | 75W |

1: VE2xx2  based on AIE-ML, VE1752 device base based on AIE

# The Only Edge AI Platform that Scales from Sensor to Accelerator on a Single Architecture[1,2]

1–100 Watts

| | 0–10TOPS | 10–25TOPS | 25–75TOPS | 100+ TOPS |
|---|---|---|---|---|
| XILINX VERSAL AI EDGE | ● | ● | ● | ● |
| NVIDIA (Jetson) | ● | ● | ● | |
| NVIDIA (T4) | | | | ● |
| NVIDIA (ORIN) | | | ● | ● |
| mobileye An Intel Company | ● | ● | ● | |
| TEXAS INSTRUMENTS | ● | ● | ● | |
| Qualcomm snapdragon | ● | ● | | |
| Qualcomm® Cloud AI 100 | | | ● | ● |
| NXP (iMX8) | ● | | | |
| RENESAS | ● | | ● | |

1: Shown in INT8 TOPS
2: Based on published sources

XILINX®

# Scalable for Different Requirements and Product Features

## Scale for Varying Levels of Compute and Safety

▶ Scale number of sensors, AI compute, vision and video processing

▶ e.g., Scale from Level-3 ADAS to Level-5 automated drive on a single platform

*Level 2    Level 3    Level 4    Level 5*

## Scale a Low-End to High-End End-Product Portfolio

▶ Design once, scale with same tools, SW, ecosystem, safety certification

▶ Scale for different price points and capabilities

## Explore Distributed vs. Centralized Architectures

▶ "Load Balance" across the system

▶ Shift compute from edge sensor to central compute across a single system
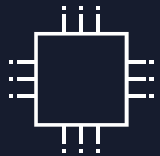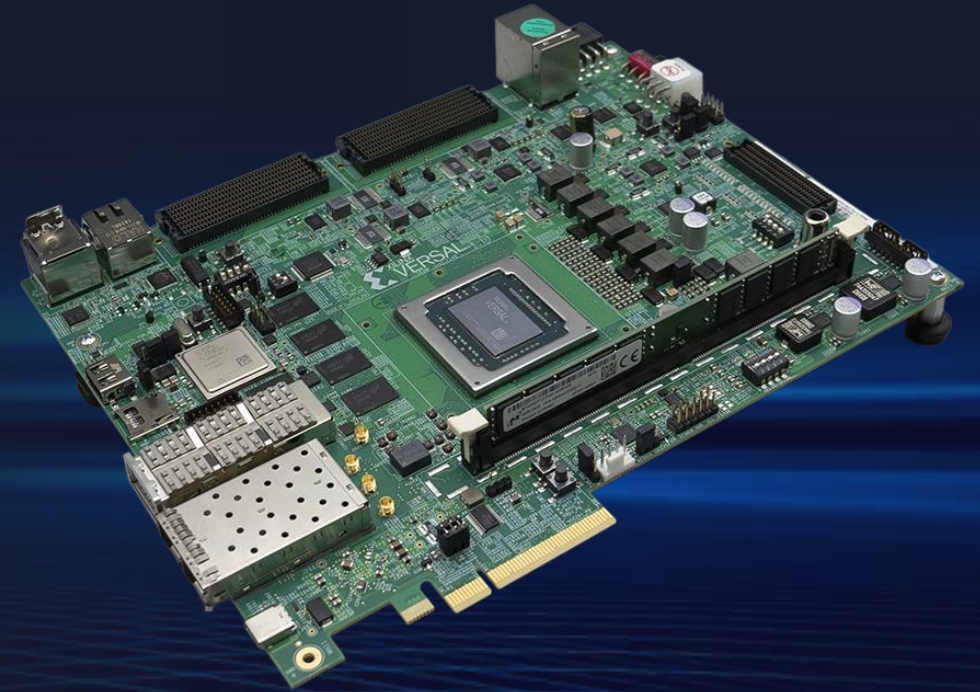
XILINX

# How Customers Can Get Started

# Availability

▶ Documentation Available Now

▶ Tools Available in 2nd Half of '21

▶ ES & Production Silicon in 1st Half '22

▶ Versal™ AI Edge ACAP Eval Kit in 2nd Half '22

XILINX.

# Start Prototyping Now

**Start Now with Versal AI Core ACAP
VCK190 Evaluation Kit**
Migrate Later to Versal AI Edge Device

| | | | | |
|---|---|---|---|---|
| Evaluate Key Blocks in Versal™ AI Edge | Leverage Vitis™ Accelerated Libraries | Breadth of Interfaces for System Testing | System-Design Methodology Guides | Guided Flows in Vitis and Vivado® Tools |

www.xilinx.com/vck190

XILINX

# Versal AI Edge ACAP: Intelligence Unleashed
*From Sensor to AI to Real-Time Control*

▶ 4X AI Performance/Watt vs. GPUs[1] with Innovations in AI Engines and Memory Hierarchy

▶ 10X Compute Density[2] with Highest Levels of Safety and Security

▶ World's Most Scalable and Adaptable Platform for Edge Systems

Silicon Sampling in 1st Half 2022

1: vs. Jetson AGX Xavier (MAX N-Mode), ResNet50 224x224, batch=1, https://developer.nvidia.com/embedded/jetson-agx-xavier-dl-inference-benchmarks
2: Compared to Zynq® UltraScale+™ MPSoCs

XILINX

**Thank You**