



A Model-driven development framework for highly Parallel and Energy-Efficient computation supporting multi-criteria optimisation

# D5.1

## Reference parallel heterogeneous hardware selection

### Version 1.3

#### Documentation Information

|                             |  |
|-----------------------------|--|
| <b>Contract Number</b>      | 871669   |
| <b>Project Website</b>      | <a href="http://www.ampere-euproject.eu">www.ampere-euproject.eu</a> |
| <b>Contractual Deadline</b> | 30.06.2020   |
| <b>Dissemination Level</b>  | PU   |
| <b>Nature</b>               | R  |
| <b>Author</b>               | Enkhtuvshin Janchivnyambuu and Jan Rollo (SYS)                       |
| <b>Contributors</b>         | Claudio Scordino (EVI) and Saoud Hadi (THALIT)                       |
| <b>Reviewer</b>             | Sara Royuela (BSC)   |
| <b>Keywords</b>             | Parallel heterogeneous, heterogeneous hardware, reference platform   |



The AMPERE project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 871669.

## Change Log

| Version | Description Change   |
|---------|--|
| V0.1    | Drafted first version by Enkhtuvshin Janchivnyambuu            |
| V0.2    | Review and contributions by Claudio Scordino and Sara Royuela. |
| V0.3    | Added Use-case requirements and NVIDIA Jetson AGX Xavier DK    |
| V0.4    | Contributions by Claudio Scordino and SAOUD Hadi               |
| V1.0    | Updated based on comments from Reviewers                       |
| V1.1    | Review and contributions by Pressler Michael                   |
| V1.2    | Review and contributions by SORRENTINO Viola                   |
| V1.3    | Review and contributions by Sara Royuela                       |

## Table of Contents

|   |    |
|---|----|
| 1. Executive Summary .....  | 3  |
| 1.1. Introduction .....   | 3  |
| 2. Heterogeneous Platform Selection Requirements.....   | 5  |
| 2.1. Requirements for Use-Cases.....  | 5  |
| 2.1.1. Requirements from the automotive use-case: Intelligent Predictive Cruise Control (PCC) ..... | 5  |
| 2.1.2. Requirements from the railway use-case: Obstacle Detection and Avoidance System (ODAS) ..... | 6  |
| 2.2. Requirements for Hypervisor and Operating Systems .....  | 8  |
| 2.3. Survey of existing COTS Platforms .....  | 9  |
| 2.4. Reference Parallel Heterogeneous Hardware .....  | 14 |
| 3. The Selected Board Description .....   | 15 |
| 3.1. Xilinx Zynq UltraScale+ MPSoC ZCU102 Evaluation Kit.....                                       | 15 |
| 3.2. NVIDIA Jetson AGX Xavier Developer Kit .....   | 18 |
| 4. Other Potential Board.....   | 20 |
| 4.1. KONIC200-FH Board .....  | 20 |
| 4.2. Evaluation of the MPPA-3 Coolidge processor .....  | 22 |
| 5. Acronyms and Abbreviations.....  | 24 |
| 6. References .....   | 25 |

# 1. Executive Summary

The European project AMPERE (A Model-driven development framework for highly Parallel and Energy-Efficient computation supporting multi-criteria optimization) [1] will implement an innovative software architecture that takes into account the non-functional requirements inherited from the cyber-physical interactions, such as time predictability, energy-efficiency, safety and security.

The aim is to provide the high-performance capabilities needed for the most advanced functionalities of cyber-physical systems (CPS). This novel technology will be employed in the automotive and railway domains.

Specifically, AMPERE will develop a novel system design and computing software ecosystem for the development and execution of CPS, targeting the most advanced energy-efficient and highly-parallel heterogeneous platforms, with the objective of fully exploiting the benefits of performance demanding emerging technologies, such as artificial intelligence or big data analytics. It will achieve that by the combination of model-driven engineering (MDE) and parallel execution, two important technical challenges at the system design and the computing software stack of CPS.

In accordance to the project's planning, Task 5.1. "Platform selection" of Work Package 5 is to select the most suitable and promising Commercial Off-The-Shelf (COTS) parallel heterogeneous processor architecture to better match the project requirements. Task 5.1 has been carried out successfully by MS1, and this deliverable includes its outcomes.

## 1.1. Introduction

This deliverable concerns the selection of the most-suitable commercial off-the-shelf (COTS) parallel heterogeneous hardware platform to be used as reference platform for the project. The consortium has already identified four potential candidates:

1. Xilinx Versal ACAP, featuring two ARM-based CPUs, a FPGA and two vector processors (GPU and DSP) [2].
2. Xilinx Zynq UltraScale+ MPSoC (CPU, GPU, FPGA and DSP) [3].
3. Kalray MPPA3 Coolidge, featuring 80 cores [4].
4. NVIDIA Jetson AGX Xavier, featuring an 8-core ARM-based CPU and a GPU composed of 512 NVIDIA cores [5], for the initial benchmarking.

Moreover, the open-source research platform HERO [8] developed at ETHZ is considered as one of the decision metrics. HERO, which will run on the FPGA board we select, is a soft system (Ariane RV64 host + PULP RV32 many-core accelerator) synthesized on the FPGA.

This selection is performed based on three key aspects: (1) previous partners' experience, (2) availability in the market, and (3) a collection of the following technical metrics:

- Cost
- Processing units
- Programmable logics

- Memory
- Compatibility with tools to be included within the AMPERE ecosystem
- Parallel programming model supported
- Debugging interface
- Efficient power management

In the scope of this project, the AMPERE ecosystem will be demonstrated in a relevant environment by developing and executing two real-world use cases (Intelligent Predictive Cruise Control in Automotive, and Obstacle Detection and Avoidance System in Railway [6]) very close to production, integrating the AMPERE ecosystem conforming to the already existing targeted environment of use-cases.

## 2. Heterogeneous Platform Selection Requirements

The characteristics of the platforms as well as the selection criteria is illustrated in this document, which will also serve to drive the selection of the further platforms to be supported during the exploitation and after the end of the project.

### 2.1. Requirements for Use-Cases

The Intelligent Predictive Cruise Control (PCC) and the Obstacle Detection and Avoidance System (ODAS) use cases have been carefully chosen to demonstrate the technologies which will be developed in the scope of AMPERE project. Both use cases are well-suited examples for the increasing needs of the heterogeneous platforms for the autonomous systems by virtue of the functional and non-functional requirements they impose (further details about the use-cases are described in the deliverable *D1.1 System models requirement and use case selection* [31]). Of particular interest are the artificial intelligence (AI) components of both use-cases, which are a core part of the corresponding implementations. Therefore, we considered the efficient implementation of Deep Neural Networks (DNN) on the computing platforms as an important factor for the platform selection.

#### 2.1.1. Requirements from the automotive use-case: Intelligent Predictive Cruise Control (PCC)

Automotive electrical/electronic (E/E) architectures are currently undergoing a radical shift in the way they are designed, implemented, and deployed. Especially, the computational power and communication bandwidth required for new functionalities, such as automated driving or connected vehicle functions (e.g. path planning, object recognition, predictive cruise control), exceed the capabilities of current compute nodes (mainly micro-controller SoCs); this is leading to a reorganization of automotive systems following the paradigm of so-called centralized E/E architectures that are based on a new class of computing nodes featuring more powerful micro-processors and accelerators such as GPUs.

One consequence of these centralized E/E architectures is that heterogeneous applications will be co-existing on the same HW platform, heterogeneous not only in their model of computation (ranging from classical periodic control over event-based planning to stream-based perception applications) but also in their criticality, in terms of real-time and safety requirements. In comparison to the previous practice to integrate mono-functional ECUs on the network level, the burden of integration is shifted from the network to the ECU level and in this regard typically from the vehicle manufacturer to the supplier of the control unit.

The PCC use case, illustrated in Figure 1, exposes the needs of the newest E/E architectures. As an illustration, shows the block diagram of the PCC architecture.

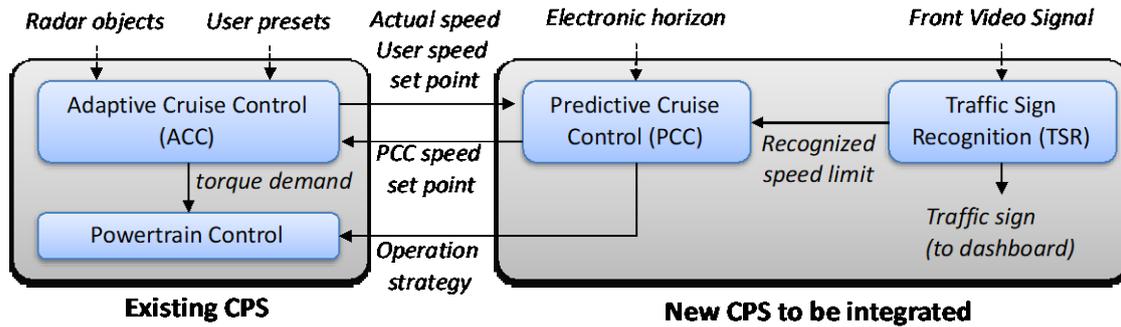


Figure 1. PCC block diagram.

As a result of the trends in E/E architectures, established methods for performance modelling and analysis that are currently used in industry are no longer effective or adequate. The following requirements for the AMPERE hardware can be derived to execute these systems and to support the analysis and control of the system performance for the PCC use case:

- The hardware platform should be supported by the Linux and the QNX operating systems.
- The hardware platform should be supporting programmable performance counters to analyse the performance of the system such as executed instructions, cycles and branch prediction misses. Highly important are counters related to the memory subsystem like cache hit/misses, write backs, and DRAM busy cycles to name a few.
- The hardware platform should have a GPU to run the trained algorithms of CNNs for a traffic sign recognition algorithm.
- The hardware platform should provide Ethernet and USB interfaces to stream sensor data from camera and sensors (e.g. radar).
- The hardware platform should include virtualization support for hypervisors.
- The hardware platform should be based on in-order cores with at least 4 cores in the system.

### 2.1.2. Requirements from the railway use-case: Obstacle Detection and Avoidance System (ODAS)

The ODAS system is intended to be a helper for a tram driver to prevent the collision against obstacles along the tracks, taking advantage of the fusion between the selected RADAR, LiDAR and optical cameras. If a potential collision is detected according to the outputs of the sensor fusion algorithm, the driver will be warned with a visual and sound alert.

The block diagram of the ODAS architecture for data processing is illustrated in Figure 2.

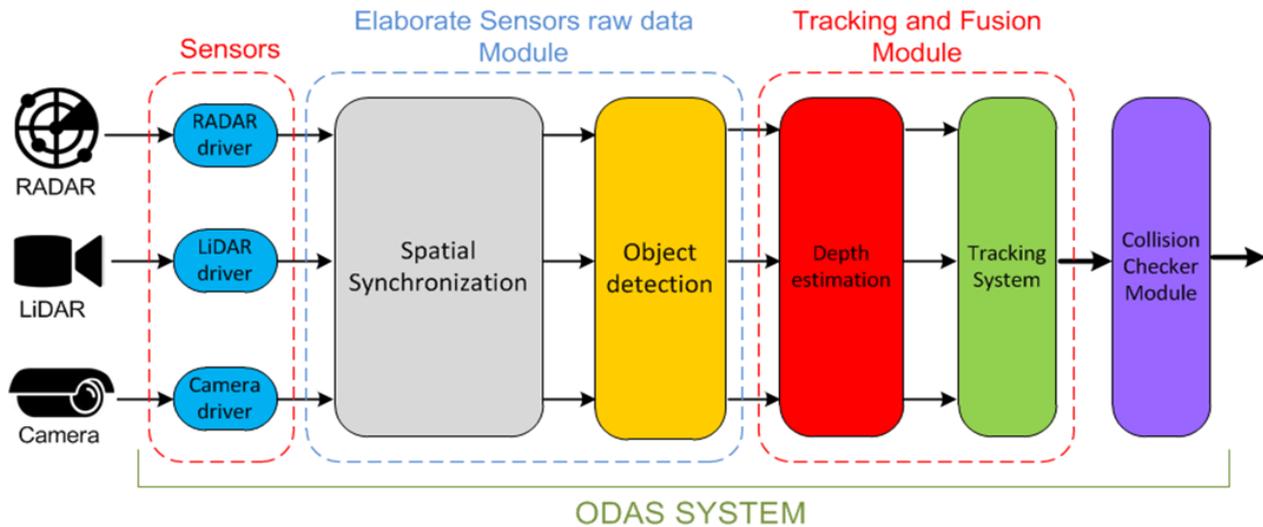


Figure 2. ODAS block diagram.

The main functionalities of the ODAS are:

- Detect objects in the Field of View (FoV) of selected sensors.
- Associate and track detected objects in the FoV.
- Warn the driver about potential collision with obstacles.

The platform should provide a Linux-based system with Robot Operating System (ROS) [32]. ROS describes a *computation graph model* where *nodes* represent processes and *topics* are named buses used by the nodes to send and receive messages. The system defines a tree-like structure that is split into levels, and a hierarchical naming structure with a namespace encapsulating each level, and defined by a forward slash followed by a name (topic names must be unique within a namespace). To send messages to a topic, a node must publish to that topic, while to receive messages it must subscribe.

In the ODAS, ROS nodes are the RADAR, the LiDAR, the CAMERAs and other entities mandatory for the correct functioning of the system, and messages can contain sensors data, sensors status, commands, etc.

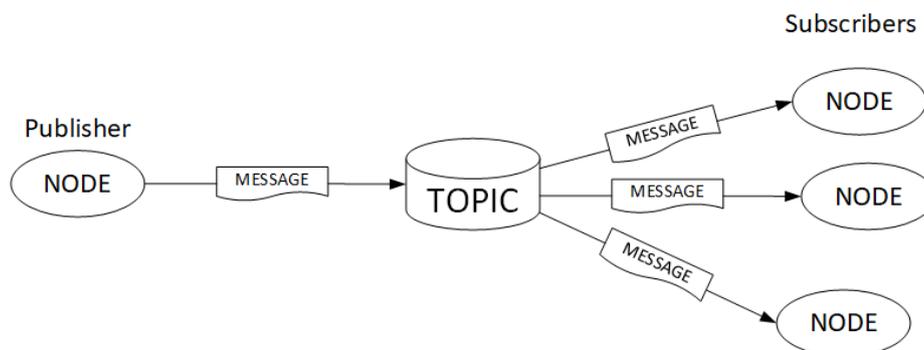


Figure 3. ROS Publisher/Subscriber paradigm adopted in the ODAS.

Due to the ODAS architectural design, the hardware requirements of the ODAS system are:

- The Ethernet and USB interfaces are required for connection with camera, RADAR and LiDAR sensors.
- The Sensor Fusion Algorithm of the ODAS shall be executed in the CPU.
- The HW platform should have a GPU to run the trained algorithms of CNN for object detection on video frames.
- HW platform's processing unit should have capability to provide the high-performance parallel execution for Object Detection, Unscented Kalman Filter and Sensor fusion algorithms.
- The HW/SW platform shall include the following libraries: ROS, OpenCV, GStreamer, Qt, Python3, PyTorch/Tensorflow, CUDA and Spdlog.
- The HW power consumption shall be less than 200 W.
- The HW platform should be able to perform dynamic power management both in the active and passive states of the system.
- The ODAS system shall use the 24 V power supply of the tramway vehicle.

## 2.2. Requirements for Hypervisor and Operating Systems

For the hypervisor and real-time operating system, the processor architecture, memory system, processor modes, hardware virtualization support and some characteristics of the microarchitecture (e.g., pipeline order, pipeline depth, branch prediction, and cache hierarchies which could prevent predictability if not properly handled) have been considered as important factors (a detailed survey on these and other architectural aspects of the different boards is provided in Table 4).

For the PikeOS hypervisor [9], the following are the minimum requirements on the hardware platform to be considered for the selection:

- CPU with a memory management unit (MMU) which is capable of restricting accesses (e.g. destinations of load and store CPU instructions) of non-privileged executables to certain memory regions and only be configurable from a privileged CPU mode.
- Support for at least two different CPU privilege modes ("user" and "supervisor" mode). Only the PikeOS Separation Kernel itself and privileged executables may run in the "supervisor" mode. Non-privileged executables always run in "user mode". In "user mode", only a limited set of instructions is available; in "supervisor mode", all instructions are available.
- The hardware (CPU or CPUs) shall provide instructions to switch between privilege modes and to use the memory management to set up different segments of memory.
- The timer facilities provided by the hardware shall be sufficient for the timing requirements (e.g., timer resolution).
- The PikeOS supports Symmetric Multiprocessing (SMP) for many cores. For implementation of the SMP, the processing unit must have cache coherence.

- For hardware-virtualized guest, the hardware platform must provide the hardware-assisted virtualization.

## 2.3. Survey of existing COTS Platforms

Based on the requirements described for the two use cases, as well as the hypervisor and operating system, this section presents a survey of existing COTS platforms.

The section is organized as follows: Tables 1, 2 and 3 summarize the most important characteristics of the COTS multi-core heterogeneous platforms currently available on the market from the HERO soft-core, hypervisor and use-case point of view; and Table 4 shows some features of CPU to be considered as key factor for portability of Hypervisor and operating system.

Table 1. List of boards proposed by partners.

| Manufacturer   | Xilinx, USA                               | Xilinx, USA                               | Xilinx, USA                               | AVNET, USA                                |
|--|---|---|---|---|
| Board name   | Xilinx Zynq ZCU102 [7]                    | Xilinx Zynq ZCU104 [10]                   | Xilinx Zynq ZCU106 [11]                   | UltraZed-EV SOM [12]                      |
| Board SoC / Module                                   | Zynq UltraScale+ XCZU9EG-2FFVB1156E MPSoC | Zynq UltraScale+ XCZU7EV-2FFVC1156 MPSoC  | Zynq UltraScale+ XCZU7EV-2FFVC1156 MPSoC  | Zynq UltraScale+ XCZU7EV-1FBVB900 MPSoC   |
| Cost   | \$2,495                                   | \$1,295                                   | \$2,495                                   | \$999                                     |
| Application Processor (CPU)                          | Quad-core ARM Cortex-A53                  | Quad-core ARM Cortex-A53                  | Quad-core ARM Cortex-A53                  | Quad-core ARM Cortex-A53                  |
| Graphics Processor (GPU)                             | ARM Mali-400 MP2                          | ARM Mali-400 MP2                          | ARM Mali-400 MP2                          | ARM Mali-400 MP2                          |
| FPGA   | System Logic Cells (K)                    | 600                                       | 504                                       | 504                                       |
|  | CLB Flip-Flops (K)                        | 548                                       | 461                                       | 461                                       |
|  | CLB LUT (K)                               | 274                                       | 230                                       | 230                                       |
|  | Block RAM (Mb)                            | 32.1                                      | 11  | 11  |
|  | Ultra RAM (Mb)                            | -   | 27  | 27  |
|  | DSP Slices                                | 2,520                                     | 1728                                      | 1728                                      |
| Memory   | 4GB DDR4                                  | 2GB DDR4                                  | 4GB DDR4                                  | 4GB DDR4                                  |
| Interfaces   | Ethernet, USB, CAN, PCIe                  |
| JTAG   | Yes                                       | Yes                                       | Yes                                       | Yes                                       |
| Parallel programming model (H: host; A: accelerator) | H: OpenMP/OmpSs/OpenCL<br>A: OmpSs/OpenCL | H: OpenMP/OmpSs/OpenCL<br>A: OmpSs/OpenCL | H: OpenMP/OmpSs/OpenCL<br>A: OmpSs/OpenCL | H: OpenMP/OmpSs/OpenCL<br>A: OmpSs/OpenCL |
| Hardware cache coherency                             | Within cluster + CCI                      |
| Suitable for HERO soft-core                          | Yes                                       | No  | No  | No  |
| Suitable for Hypervisor/OS                           | Yes                                       | Yes                                       | Yes                                       | Yes                                       |
| Suitable for PCC use-case                            | Yes                                       | Yes                                       | Yes                                       | Yes                                       |
| Suitable for ODAS use-case                           | No  | No  | No  | No  |

Table 2. List of boards proposed by partners (continuation).

| Manufacturer  |                        | AVNET, USA                          | NVIDIA, USA                                      | Xilinx, USA                         | Xilinx, USA                         |
|---|------------------------|-------------------------------------|--|-------------------------------------|-------------------------------------|
| Board name  |                        | Ultra96-V2 [13]                     | Jetson AGX Xavier DK [14]                        | Versal Prime VMK180 EK [15]         | Versal AI Core VCK190 EK [16]       |
| Board SoC / Module                                      |                        | Zynq UltraScale+ MPSoC ZU3EG A484   | Jetson AGX Xavier                                | Versal VM1802 ACAP                  | Versal VC1902 ACAP                  |
| Cost  |                        | \$249                               | \$699  | Not available                       | Not available                       |
| Application Processor (CPU)                             |                        | Quad-core ARM Cortex-A53            | 8-Core Carmel ARM v8.2                           | Dual-Core ARM Cortex-A72            | Dual-Core ARM Cortex-A72            |
| Graphics processor (GPU)                                |                        | ARM Mali-400 MP2                    | 512-Core Volta w/Tensor Core                     | -                                   | -                                   |
| FPGA  | System Logic Cells (K) | 154                                 | -  | 1,968                               | 1968                                |
|   | CLB Flip-Flops (K)     | 141                                 | -  | 1,450                               | -                                   |
|   | CLB LUT (K)            | 71                                  | -  | 899.8                               | 899.8                               |
|   | Block RAM (Mb)         | 7.6                                 | -  | 34                                  | 34                                  |
|   | Ultra RAM (Mb)         | -                                   | -  | 130                                 | 130                                 |
|   | DSP Slices             | 360                                 | -  | 1,968                               | 1968 (3.2 TFLOPS/1.9 TOPS)          |
| Memory  |                        | 2GB LPDDR4                          | 32 GB LPDDR4                                     | 8GB DDR4                            | 8GB DDR4                            |
| Interfaces  |                        | Ethernet, USB, CAN, PCIe            | Ethernet, USB A/C, HDMI, PCIe                    | Ethernet, CAN, USB                  | Ethernet, CAN, USB                  |
| JTAG  |                        | No                                  | -  | Yes                                 | Yes                                 |
| Parallel programming model<br>(H: host; A: accelerator) |                        | H: OpenMP/OmpSs/OpenCL<br>A: OpenCL | H: OpenMP/OmpSs/OpenCL<br>A: CUDA/OpenCL/OpenACC | H: OpenMP/OmpSs/OpenCL<br>A: OpenCL | H: OpenMP/OmpSs/OpenCL<br>A: OpenCL |
| Hardware cache coherency                                |                        | Within cluster + CCI                | Within cluster + I/O coherency                   | Within cluster + CCI                | Within cluster + CCI                |
| Suitable for HERO soft-core                             |                        | No                                  | No   | Yes                                 | Yes                                 |
| Suitable for Hypervisor/OS                              |                        | Yes                                 | Yes  | -                                   | -                                   |
| Suitable for PCC use-case                               |                        | Yes                                 | Yes  | No                                  | No                                  |
| Suitable for ODAS use-case                              |                        | No                                  | Yes  | No                                  | No                                  |

Table 3. List of boards proposed by partners (continuation).

| Manufacturer  |                        | NXP, Netherlands                    | Toradex, Switzerland                | Kalray, France                  | Kalray, France                  |
|---|------------------------|-------------------------------------|-------------------------------------|---------------------------------|---------------------------------|
| Board name  |                        | NXP iMX 8QuadMax MEK [23]           | Apalis iMX8 CoM [24]                | KONIC board [25]                | TurboCard [25]                  |
| Board SoC / Module                                      |                        | NXP iMX 8QuadMax                    | Apalis iMX8 – i.MX8QuadMax          | KONIC200 – FH/LP/HP             | TurboCard K200                  |
| Cost  |                        | \$1498 (\$999 + \$499)              | \$252.9                             | -                               | -                               |
| Application Processor (CPU)                             |                        | 4x Cortex-A53 + 2x Cortex-A72       | 4x Cortex-A53 + 2x Cortex-A72       | 1-2x MPPA3 Coolidge             | 1x MPPA3 Coolidge               |
| Graphics processor (GPU)                                |                        | 2x GC7000Lite/XSVX                  | 2x GC7000Lite/XSVX                  | -                               | -                               |
| FPGA  | System Logic Cells (K) | -                                   | -                                   | -                               | -                               |
|   | CLB Flip-Flops (K)     | -                                   | -                                   | -                               | -                               |
|   | CLB LUT (K)            | -                                   | -                                   | -                               | -                               |
|   | Block RAM (Mb)         | -                                   | -                                   | -                               | -                               |
|   | Ultra RAM (Mb)         | -                                   | -                                   | -                               | -                               |
|   | DSP Slices             | HiFi 4 DSP                          | HiFi 4 DSP                          | -                               | -                               |
| Memory  |                        | LPDDR4                              | 2-4GB LPDDR4 RAM                    | 4-64GB DDR4                     | 2x 4GB DDR4                     |
| Interfaces  |                        | Ethernet, USB                       | Ethernet, PCIe, USB, CAN            | PCIe                            | PCIe                            |
| JTAG  |                        | Yes                                 | Yes                                 | Yes                             | Yes                             |
| Parallel programming model<br>(H: host; A: accelerator) |                        | H: OpenMP/OmpSs/OpenCL<br>A: OpenCL | H: OpenMP/OmpSs/OpenCL<br>A: OpenCL | OpenMP/OpenCL                   | OpenMP/OpenCL                   |
| Hardware cache coherency                                |                        | Within cluster + CCI                | Within cluster + CCI                | Within cluster                  | Within cluster                  |
| Suitable for HERO soft-core                             |                        | No                                  | No                                  | No                              | No                              |
| Suitable for Hypervisor/OS                              |                        | Yes                                 | Yes                                 | No, due to no details about ISA | No, due to no details about ISA |
| Suitable for PCC use-case                               |                        | No                                  | No                                  | No                              | No                              |
| Suitable for ODAS use-case                              |                        | No                                  | No                                  | No                              | No                              |

Table 4. Key characteristics of the board processing units.

| CPU                               | ISA       | CPU Frequency | Memory system | Pipeline order | Pipeline depth       | Cache hierarchies  | Branch prediction | Hardware virtualization | Suitable for hypervisor |
|-----------------------------------|-----------|---------------|---------------|----------------|----------------------|--|-------------------|-------------------------|-------------------------|
| <b>ARM Cortex-A72</b> [17]        | ARMv8-A   | Up to 1.6GHz  | MMU           | Out of order   | 15-stage superscalar | L1: Instruction/data<br>L2: Single shared                              | Yes               | Yes                     | Yes                     |
| <b>ARM Cortex-A53</b> [18]        | ARMv8-A   | Up to 1.5GHz  | MMU           | In order       | 8-stage superscalar  | L1: Instruction/data<br>L2: Single shared                              | Yes               | Yes                     | Yes                     |
| <b>ARM Cortex-R5</b> [19]         | ARMv7-R   | Up to 600MHz  | MPU           | In order       | 8-stage dual issue   | L1: Instruction/data<br>TCM  | Yes               | No                      | No                      |
| <b>Nvidia Carmel ARMv8.2</b> [20] | ARMv8.2-A | Up to 2.26GHz | MMU           | -              | -                    | L1: Instruction/data<br>L2: Shared per duplex<br>L3: Shared by cluster | Yes               | Yes                     | Yes                     |
| <b>ARM Cortex-M4F</b> [26]        | ARMv7-M   | 120 MHz       | MPU           | -              | 3-stage              | No   | No                | No                      | No                      |
| <b>Kalray MPPA3 Coolidge</b> [4]  | 32/64-bit | 1.2GHz        | MMU           | -              | -                    | L1: Instruction/data<br>L2: Shared by cluster                          | -                 | Yes                     | No                      |

## 2.4. Reference Parallel Heterogeneous Hardware

After evaluating the wide range of boards proposed by partners in Tables 1, 2 and 3, and taking into account the technical factors, the previous experience of the partners and the requirements listed above, the following platforms have been selected as reference hardware for further development:

- Xilinx Zynq UltraScale+ MPSoC ZCU102 Evaluation Kit [7], and
- NVIDIA Jetson AGX Xavier Development Kit [14]

Majority of partners voted to select the Xilinx Zynq UltraScale+ MPSoC ZCU102 Evaluation Kit due to PCC use-case, HERO soft-core and Hypervisor/OS requirement matching as well as their previous experience in this platform. Moreover, it contains a programmable FPGA, allowing research on soft-core RISC-V technology by ETHZ as well as predictable acceleration by SSSA.

Selection of the NVIDIA Jetson AGX Xavier Development Kit was driven by ODAS use case requirement which would use the CUDA for its development.

One of the important factors we considered is the availability of this board in the current market and all partners can purchase it directly for the further development without depending on its market availability. Due to this consideration as well as no support of the PikeOS hypervisor, Xilinx Versal Prime, Xilinx Versal AI Core, Kalray KONIC and Kalray TurboCard boards were not chosen (*please note that at the moment selecting reference hardware platforms, those were not available in the market*).

To make a decision between a different Xilinx Zynq evaluation boards with same MPSoC, a bigger size of Block RAM has been considered due to the open-source research platform HERO requirement as the programmable logic (PL) memory. Generally, a RAM of the Xilinx Zynq UltraScale+ evaluation board is split up between Block RAM and UltraRAM, which are different modules that have qualitative and quantitative differences. UltraRAMs are denser but slower while Block RAM is faster.

As shown in Tables 1-3, only Zynq UltraScale+ XCZU9EG-2FFVB1156E MPSoC-based platform, that is the Xilinx Zynq UltraScale+ MPSoC ZCU102 Evaluation Kit, satisfies the HERO soft-core requirement for a bigger size of Block RAM.

## 3. The Selected Board Description

This section describes the architectural details of the selected boards.

### 3.1. Xilinx Zynq UltraScale+ MPSoC ZCU102 Evaluation Kit

The Xilinx Zynq ZCU102 is a general purpose evaluation board for rapid prototyping of automotive, industrial, video, and communications applications. It is based on the Zynq UltraScale+ XCZU9EG-2FFVB1156E MPSoC, which combines a powerful processing system (PS) and user-programmable logic (PL) into the same device. The overview of the ZCU102 evaluation board is shown in **Error! Reference source not found.**

The evaluation board comes with the following key features:

- Optimized for quick application prototyping with Zynq UltraScale+ MPSoC.
- DDR4 SODIMM – 4GB 64-bit w/ ECC attached to processing system (PS).
- DDR4 Component – 512MB 16-bit attached to programmable logic (PL).
- PCIe® Root Port Gen2 x4, USB3, Display Port, and SATA.
- 4x SFP+ interfaces for Ethernet.
- 2x FPGA Mezzanine Card (FMC) interfaces for I/O expansion, including 16 16.3Gb/s GTH transceivers and 64 user-defined differential I/O signals.
- Configuration from QSPI.
- Configuration from SD card.
- Configuration over JTAG with PC4 header.
- Configuration over JTAG with Arm 20-pin header.
- Configuration over USB-to-JTAG bridge.

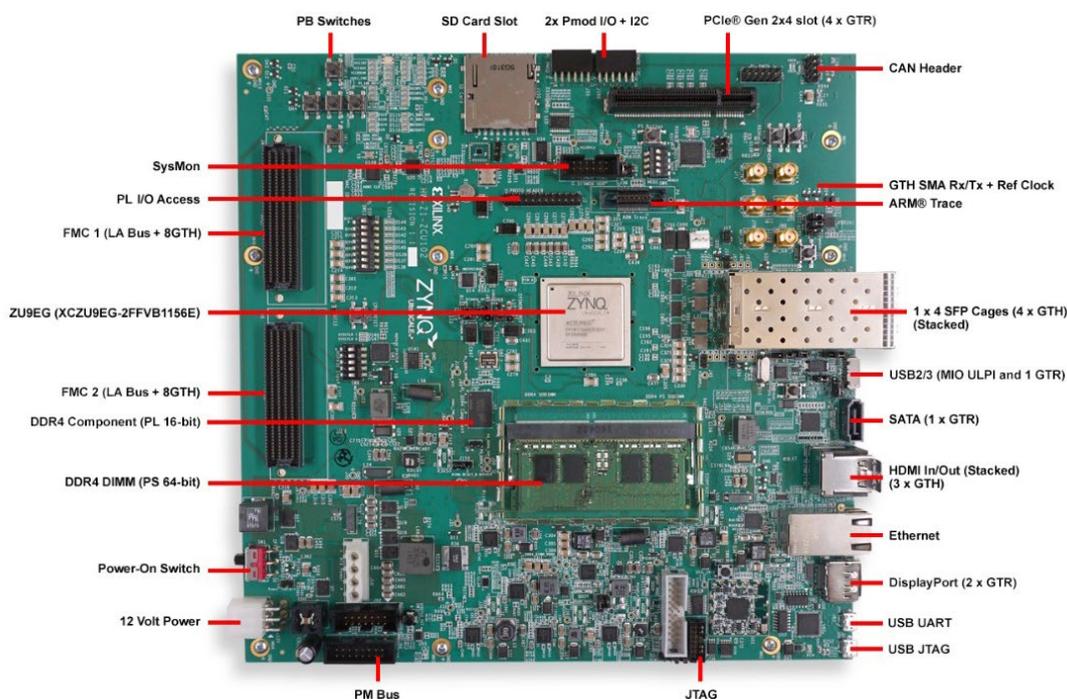


Figure 4. Features of ZCU102 Evaluation Board.

For more details, please refer to the ZCU102 board user manual [21].

The Zynq UltraScale+ MPSoC Top-Level Block Diagram is shown in Figure 5 and it has three major processing units:

- Cortex-A53 application processing unit (APU)-Arm v8 architecture-based 64-bit quad-core multiprocessing CPU with 8-stage pipelined processor with 2-way superscalar, in-order execution pipeline and with memory management unit (MMU).
- Cortex-R5 real-time processing unit (RPU)-Arm v7 architecture-based 32-bit dual real-time processing unit with dedicated tightly coupled memory (TCM) and memory protection unit (MPU).
- Mali-400 graphics processing unit (GPU)-graphics processing unit with pixel and geometry processor and 64 KB L2 cache.

Also, the Zynq UltraScale+ MPSoC PS has four high-speed serial I/O (HSSIO) interfaces supporting the following protocols:

- Integrated block for PCI Express interface-PCIe™ base specification version 2.1 compliant.
- SATA 3.1 specification compliant interface.
- DisplayPort interface-implements a DisplayPort source-only interface with video resolution up to 4K x 2K-30 (300 MHz pixel rate).
- USB 3.0 interface-compliant to USB 3.0 specification implementing a 5 Gb/s line rate.
- Serial GMII interface-supports a 1 Gb/s SGMII interface.

The ARM Cortex-A53 application processing unit of the Zynq UltraScale+ MPSoC is a high efficiency processor and has the following features:

- 8-stage pipelined processor with 2-way superscalar, in-order execution pipeline.
- DSP and NEON SIMD extensions are mandatory per core.
- VFPv4 Floating Point Unit onboard (per core).
- Hardware virtualization support.
- TrustZone security extensions.
- 64-byte cache lines.
- 10-entry L1 TLB, and 512-entry L2 TLB.
- 4 KiB conditional branch predictor, 256-entry indirect branch predictor.

For additional information on Zynq UltraScale+ MPSoC devices, please see the Zynq UltraScale+ MPSoC Data Sheet [3]. For more information about Zynq UltraScale+ MPSoC configuration options, please refer to the Zynq UltraScale+ MPSoC Technical Reference Manual [22].

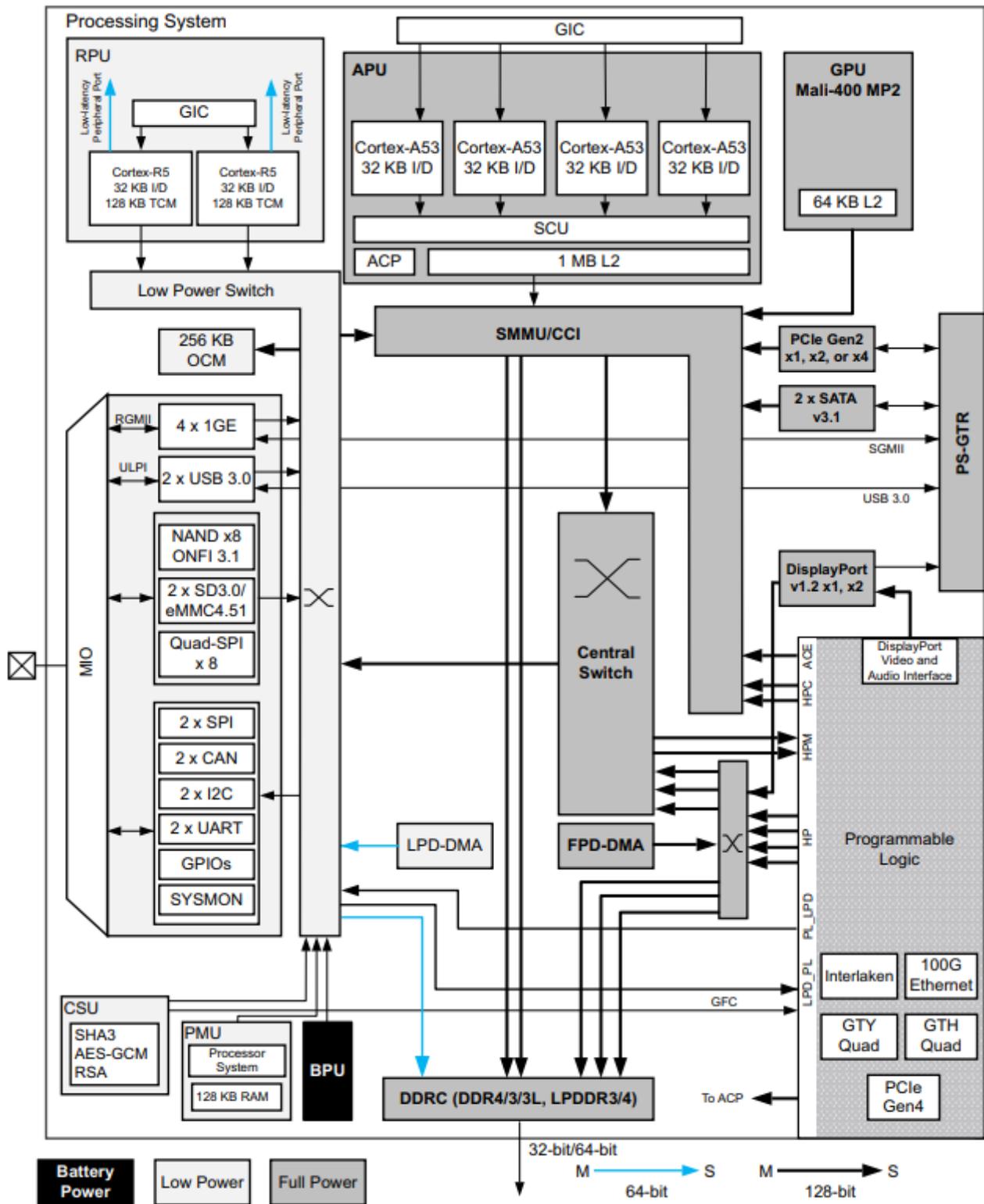


Figure 5. Zynq UltraScale+ MPSoC Top-Level Block Diagram.

## 3.2. NVIDIA Jetson AGX Xavier Developer Kit

NVIDIA Jetson AGX Xavier is an embedded system-on-module (SoM) from the NVIDIA AGX Systems family, including an integrated Volta GPU with Tensor Cores, dual Deep Learning Accelerators (DLAs), octal-core NVIDIA Carmel ARMv8.2 CPU, 32GB 256-bit LPDDR4x with 137GB/s of memory bandwidth, and 650Gbps of high-speed I/O including PCIe Gen 4 and 16 camera lanes of MIPI CSI-2. The NVIDIA Jetson AGX Xavier module delivers up to 32 TOPS of accelerated computing capability in a compact form factor consuming under 30 Watts. This advanced system-on-module is designed specifically for autonomous machines. Heterogeneous accelerated computing architecture delivers advanced edge capabilities. Plus, it comes with integrated memory, storage, power management, and an innovative thermal design to enable faster time to market.

Jetson AGX Xavier is supported by NVIDIA JetPack software development kit (SDK), which includes a board support package (BSP), an Ubuntu Linux OS, NVIDIA CUDA, cuDNN, and TensorRT software libraries for deep learning, computer vision, GPU computing, multimedia processing, and much more. It's also supported by the NVIDIA DeepStream SDK, which delivers a complete toolkit for real-time situational awareness through intelligent video analytics (IVA). This helps you boost performance and accelerate software development, while reducing development cost and effort.

Below is a list of the module's key features as depicted in Figure 6:

### Processing Components:

- Octal-core NVIDIA Carmel ARMv8.2 CPU @ 2.26GHz
- 512-core Volta GPU @ with 64 Tensor Cores
- Dual Deep Learning Accelerator (DLA) engines
- 32GB 256-bit LPDDR4x @ 2133MHz (137GB/s)
- 32GB eMMC 5.1
- Vision Accelerator engine
- (4x) 4Kp60 H.264/H.265 video encoder
- (2x) 8Kp30 / (6x) 4Kp60 H.265 video decoder

### I/O Interfaces & Ports:

- (16x) MIPI CSI-2 lanes, (8x) SLVS-EC lanes
- up to 6 active sensor streams and 36 virtual camera channels
- (5x) PCIe Gen 4 controllers | 1x8, 1x4, 1x2, 2x1
- (3x) Root Port & Endpoint
- (2x) Root Port
- (3x) USB 3.1 + (4x) USB 2.0
- (3x) eDP 1.4 / DP 1.2 / HDMI 2.0 @ 4Kp60
- 10/100/1000 BASE-T Ethernet + MAC + RGMII PHY
- Dual CAN bus controller
- UART, SPI, I2C, I2C, GPIOs

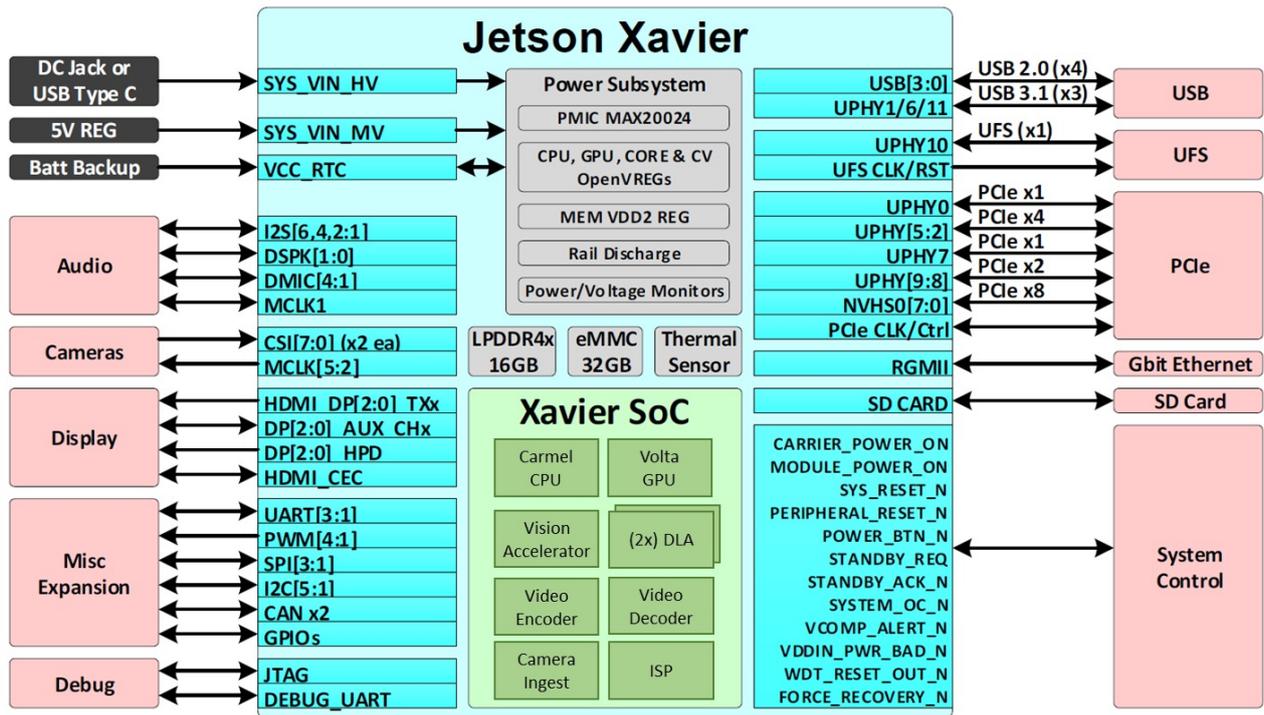


Figure 6. NVIDIA Jetson AGX Xavier SoM Block Diagram.

For more detailed information on NVIDIA Jetson AGX Xavier SoM, please see the NVIDIA Jetson Xavier System-on-Module [27].

## 4. Other Potential Board

As explained in section 2.4, the selected reference boards for the AMPERE ecosystem do not include a board with an MPPA-3 processor due to the difficulty to purchase such a board during the first months of the project and no support of the hypervisor. Nevertheless, Thales R&T had the opportunity to evaluate the KONIC200-FH with an MPPA-3 processor and to develop several benchmarks, providing performance measurements that can later be compared to the selected platforms. Section 4.1 introduces the KONIC200-FH board, and Section 4.2 discusses the evaluation of the MPPA-3 Coolidge processor.

### 4.1. KONIC200-FH Board

The KONIC200-FH [28] is an accelerator board from Kalray that embeds an MPPA-3 Coolidge processor [4]. It is one of the few boards with a Coolidge that is currently available (June 2020). The KONIC200-FH (K200) takes the form of a Full-Height, Half-Length (FHHL) PCIe card with one MPPA-3 Coolidge processor that includes 80 64-bit CPU cores and 80 AI coprocessors. The cores are organized in 5 clusters; each of them includes 16 Kv3 cores and 16 coprocessors. The processor is clocked at 1.2 GHz. The K200 board has a power consumption of 45W and it can deliver up to 25 TOPS with INT8 data, up to 1.15 TFLOPS in single precision and up to 284 GFLOPS in double precision.

The K200 board is included in a “MPPA Developer platform” (or MPPA-Dev) to form a powerful evaluation and development platform with a full set of debug and monitoring tools. One can either use the K200 in a standalone mode or in an accelerator mode:

- In the standalone mode, the K200 acts as a standalone board with the ability to run a Linux OS or an RTOS on one cluster and execute application(s) on the four others while the MPPA-Dev acts as a debug and monitoring probe (binary files and data can be loaded using the ethernet or the JTAG).
- In the accelerator mode, the MPPA-Dev becomes a host machine and the K200 acts as an accelerator card with all the cores available for the user applications. This mode is parallel to a GPU running an OpenCL program: the binary files and data are loaded from the MPPA-dev host through the PCIe interface.

The Coolidge SDK supports GCC and LLVM toolchains. It includes optimized mathematical libraries (BLAS, LAPACK,...), a support of OpenCL (1.2 Embedded profile), an AI framework (KaNN) and an optimized backend for OpenCV.

Figure 7 is an overview of the processor architecture with zooms on a single cluster and on a single core. The list below is a summary of the key characteristics and features of the Coolidge processor.

MPPA-3 Coolidge characteristics summary:

- Core
  - 64-bit/32-bit architecture
  - 6-issue VLIW
  - 16KB instruction cache / 16KB data cache with MMU
  - IEEE 754-2008 Floating Point Unit (FPU)
  - Square root and reciprocal operations in floating single precision

- 64-bit integer multiplication (Asymmetric cryptography)
- Up to 4 execution rings
- Up to 256-bits per cycle Load/Store
- Co-processor (one per Core)
  - INT8, INT16 or FP16 accuracy
  - Up to 128 MAC per cycle
  - 16 x 64-bit Cores + 1 dedicated safety/security Core
  - 4 MB L2 Cache – 600GB/s
  - Configurable cluster/chip cache coherency modes
  - Low Power (600MHz) / Standard (900MHz)/ High performance (1.2GHz ) modes
- System-on-Chip
  - 5 clusters (total of 80 Cores + 5 security Cores)
  - 40GB/s High Speed/Low Latency Network-on-Chip
  - Support of chip-to-chip connection to scale performance

KONIC200-FH characteristics summary:

- 1x MPPA-3 Coolidge
- 1x PCIe Gen3 x16
- 2x QSFP28 (2x 100G)
- 4GB DDR4
- JTAG over USB and trace

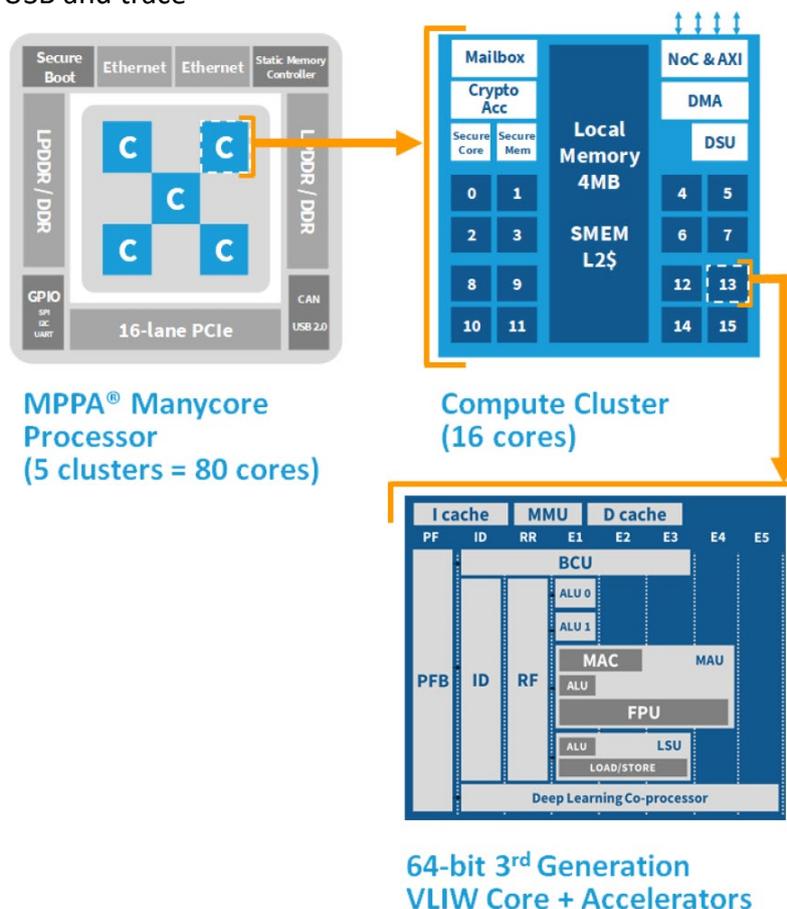


Figure 7. MPPA-3 architecture overview.

## 4.2. Evaluation of the MPPA-3 Coolidge processor

The applications used for the evaluation of the Coolidge are:

- CoreMark benchmark [29], executed on a single core of the Coolidge,
- YOLOv3 convolutional neural network [30]: the application has been implemented using the Kalray AI framework “KaNN” with FP16 values and executed on images from a 3840x2160 video stream. The KaNN framework uses the whole chip, i.e. the 5 clusters with both all the processor cores and all the coprocessors.
- A set of custom OpenCL applications that are representative of signal processing applications and more specifically Radar applications. These applications include Fast Fourier Transform (FFT), Pulse Compression (PC), Moving Target Indicator (MTI), Doppler operations and so on.

The applications have been executed on both an MPPA-Dev and an x86 computer with an nVidia GeForce GTX 980 Ti graphic card to get reference results. The nVIDIA GeForce® GTX 980Ti main characteristics are:

- 1x PCIe Gen3 x16
- Core Clocks 1076 MHz / 1000 MHz
- Power consumption of 250W

All the tests have been done using the latest Kalray’s AccessCore SDK (version 4.1).

Table 5. Coolidge benchmark results.

| Application | Coolidge Results |         |
|-------------|------------------|---------|
| CoreMark    | Score:           | 3.21    |
| YOLOv3      | Time/frame:      | 66.5 ms |
|             | Median power:    | 14.97W  |
|             | Peak power:      | 47.54W  |

Table 6. Coolidge & GeForceGTX 980Ti benchmark results.

|                  |                      | # OPs/cycles | # OPs/cycles /core | Performance efficiency | Effective GOPs/s | Effective GOPs/s/W | Effective GOPs/s/core |
|------------------|----------------------|--------------|--------------------|------------------------|------------------|--------------------|-----------------------|
| MPPA             | Custom OpenCL app. A | 50,33        | 0,63               | 0,10                   | 60,39            | 1,34               | 0,75                  |
|                  | Custom OpenCL app. B | 35,16        | 0,44               | 0,07                   | 42,19            | 0,94               | 0,53                  |
|                  | Custom OpenCL app. C | 47,21        | 0,59               | 0,10                   | 56,65            | 1,26               | 0,71                  |
| GeForceGTX 980Ti | Custom OpenCL app. A | 1443,48      | 0,51               | 0,51                   | 1443,48          | 5,77               | 0,51                  |
|                  | Custom OpenCL app. B | 1190,38      | 0,42               | 0,42                   | 1190,38          | 4,76               | 0,42                  |
|                  | Custom OpenCL app. C | 1428,77      | 0,51               | 0,51                   | 1428,77          | 5,72               | 0,51                  |

Table 5 shows the MPPA-3 score on CpuMark and the YOLOv3 execution characteristics. Table 6 includes the results of executing the custom OpenCL applications on the MPPA-3 and on the GeForceGTX 980Ti.

From the Table 6, one observes that the reference GPU computes more operation per second than the MPPA-3, which is mainly due to the higher number of core on the GPU (2816 CUDA cores vs. 80 Kv3 cores). Thus, a better metric to compare such opposite platforms would be to normalize the OPs/cycle results and compute the number of operation per cycle per core. From this metric, one observes that the MPPA-3 computes a slightly higher number of operations per core per cycle than its equivalent on the NVIDIA board.

The previous metric does not reflect the energy efficiency, which can be of prime concern when addressing embedded applications as the automotive or the avionic ones. So, we also computed the GOPs/s/W and GOPs/s/core, making it easier to compare the MPPA-3 with the NVIDIA GPU or with any other hardware platform with OpenCL support. From observing the GOPs/s/W column, one deduces that the energy efficiency of the Coolidge is up to four times less than the 980Ti.

From this benchmarks, it appears that the MPPA-3 has a great potential with a high value for GOPs/s/core, but have two main drawbacks:

- The delivered GOPs/s per W is not as high as on the GeForceGTX 980Ti , which means that the latter has a better power efficiency
- Even if the GOPs/s per MPPA-3 core is better, the total number of core and thus the GOPs/s is less than on the GeForceGTX 980Ti, making the Coolidge more suitable for applications with moderate computational power requirements.

## 5. Acronyms and Abbreviations

- AI – Artificial Intelligence
- AMPERE – A Model-driven development framework for highly Parallel and Energy-Efficient computation supporting multi-criteria optimization
- APU – Application Processing Unit
- COTS – Commercial Off-The-Shelf
- CNN – Convolutional Neural Network
- CPS – Cyber Physical System
- CPU – Central Processing Unit
- DK – Development Kit
- DMA – Direct Memory Access
- DNN – Deep Neural Networks
- DSP – Digital Signal Processing
- ECU – Electronic Control Unit
- E/E – Electrical/Electronic
- EK – Evaluation Kit
- FoV – Field of View
- FPD – Full Power Domain
- FPGA – Field Programmable Gate Array
- GIC – Generic Interrupt Controller
- GPU – Graphics Processing Unit
- HW – Hardware
- LPD – Low Power Domain
- MDE – Model-Driven Engineering
- MIO – Multiplexed I/O
- MMU – Memory Management Unit
- MPSoC – Multiprocessor SoC
- MPU – Memory Protection Unit
- ODAS – Obstacle Detection and Avoidance System
- OS – Operating System
- PL – Programmable Logic
- PMU – Platform Management Unit
- ROS – Robot Operating System
- RPU – Real-time Processing Unit
- SCU – Snoop-Control Unit
- SDK – Software Development Kit
- SW – Software
- SoC – System-on-Chip
- SoM – System-on-Module

## 6. References

- [1] "AMPERE Project," [Online]. Available: [www.ampere-euproject.eu](http://www.ampere-euproject.eu).
- [2] "Xilinx Versal ACAP," [Online]. Available: <https://www.xilinx.com/products/silicon-devices/acap/versal.html>.
- [3] "Xilinx Zynq UltraScale+ MPSoC Data Sheet," Xilinx, [Online]. Available: [https://www.xilinx.com/support/documentation/data\\_sheets/ds891-zynq-ultrascale-plus-overview.pdf](https://www.xilinx.com/support/documentation/data_sheets/ds891-zynq-ultrascale-plus-overview.pdf).
- [4] "MPPA3-80 Coolidge," [Online]. Available: <https://www.kalrayinc.com/portfolio/processors/>.
- [5] "Nvidia Jetson AGX Xavier," [Online]. Available: <https://developer.nvidia.com/embedded/jetson-agx-xavier>.
- [6] AMPERE Deliverable D1.6, "Use case evaluation", expected in December 2022.
- [7] "Xilinx Zynq UltraScale+ MPSoC ZCU102 Evaluation Kit," Xilinx, [Online]. Available: <https://www.xilinx.com/products/boards-and-kits/ek-u1-zcu102-g.html>.
- [8] "HERO Open Heterogeneous Research Platform," [Online]. Available: <https://pulp-platform.org/hero.html>.
- [9] "PikeOS Certified Hypervisor," [Online]. Available: <https://www.sysgo.com/products/pikeos-hypervisor/>.
- [10] "Xilinx Zynq UltraScale+ MPSoC ZCU104 Evaluation Kit," [Online]. Available: <https://www.xilinx.com/products/boards-and-kits/zcu104.html>.
- [11] "Xilinx Zynq UltraScale+ MPSoC ZCU106 Evaluation Kit," [Online]. Available: <https://www.xilinx.com/products/boards-and-kits/zcu106.html>.
- [12] "UltraZed-EV Starter Kit," [Online]. Available: <http://zedboard.org/product/ultrazed-ev>.
- [13] "Ultra96-V2 Development Board," [Online]. Available: <http://zedboard.org/product/ultra96-v2-development-board>.
- [14] "JETSON AGX XAVIER Development Kit," [Online]. Available: <https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-agx-xavier/>.
- [15] "Versal Prime Series VMK180 Evaluation Kit," [Online]. Available: <https://www.xilinx.com/products/boards-and-kits/vmk180.html>.
- [16] "Versal AI Core Series VCK190 Evaluation Kit," [Online]. Available: <https://www.xilinx.com/products/boards-and-kits/vck190.html>.
- [17] "ARM Cortex-A72," [Online]. Available: <https://developer.arm.com/ip-products/processors/cortex-a/cortex-a72>.
- [18] "ARM Cortex-A53," [Online]. Available: <https://developer.arm.com/ip-products/processors/cortex-a/cortex-a53>.

products/processors/cortex-a/cortex-a53.

- [19] "ARM Cortex-R5," [Online]. Available: <https://developer.arm.com/ip-products/processors/cortex-r/cortex-r5>.
- [20] "Nvidia Carmel ARM v8.2," [Online]. Available: <https://en.wikichip.org/wiki/nvidia/microarchitectures/carmel>.
- [21] "Xilinx Zynq UltraScale+ MPSoC ZCU102 Evaluation Board User Guide," Xilinx, [Online]. Available: [https://www.xilinx.com/support/documentation/boards\\_and\\_kits/zcu102/ug1182-zcu102-eval-bd.pdf](https://www.xilinx.com/support/documentation/boards_and_kits/zcu102/ug1182-zcu102-eval-bd.pdf).
- [22] "Xilinx Zynq UltraScale+ MPSoC Technical Reference Manual," [Online]. Available: [https://www.xilinx.com/support/documentation/user\\_guides/ug1085-zynq-ultrascale-trm.pdf](https://www.xilinx.com/support/documentation/user_guides/ug1085-zynq-ultrascale-trm.pdf).
- [23] "NXP i-mx8QuadMax Multisensory Enablement Kit," NXP, [Online]. Available: <https://www.nxp.com/design/development-boards/i-mx-evaluation-and-development-boards/i-mx-8quadmax-multisensory-enablement-kit-mek:MCIMX8QM-CPU>.
- [24] "Apalis iMX8," Toradex, [Online]. Available: <https://www.toradex.com/computer-on-modules/apalis-arm-family/nxp-imx-8>.
- [25] "Kalray Boards," Kalray, [Online]. Available: <https://www.kalrayinc.com/portfolio/boards/>.
- [26] "ARM Cortex-M4F," ARM, [Online]. Available: <https://developer.arm.com/ip-products/processors/cortex-m/cortex-m4>.
- [27] "NVIDIA Jetson Xavier System-on-Module Datasheet," [Online]. Available: <https://developer.nvidia.com/embedded/downloads#?search=Data%20Sheet>.
- [28] "KONIC200, Programmable Accelerator Cards for Data Centers," [Online]. Available: <https://www.kalrayinc.com/download/konic-80-200/>.
- [29] "CoreMark, an EEMBC benchmark", [Online]. Available: <https://www.eembc.org/coremark/>.
- [30] Joseph Redmon and Ali Farhadi, "YOLOv3: An Incremental Improvement," 2018.
- [31] Thales Italy, "Deliverable 1.1 - System models requirement and use case selection".
- [32] "Robotic Operating System," ROS, [Online]. Available: <https://docs.ros.org/en/foxy/index.html>.