# KALRAY

# MPPA® Manycore Processor: At the heart of Intelligent Systems

Benoît Dupont de Dinechin, CTO

June 2019

www.kalrayinc.com

# Kalray at Glance

**We design processors at the heart of new intelligent systems**

**3**
sites
Grenoble (France) ,
Los Altos (USA), Tokyo (Japan)

**~85**
Staff Members
~ 75 engineers & PhDs

**23**
patent families

**Financial and industrial shareholders**

cea investissement
AMORÇAGE TECHNOLOGIQUE

INOCAP Gestion

MBDA
MISSILE SYSTEMS

RENAULT NISSAN MITSUBISHI

ACE
PRIVATE EQUITY

Pengpai

SAFRAN

bpifrance

EURONEXT

**~ 48M€**
raised at IPO in June 2018

KALRAY

# Outline

**Intelligent Systems**

Manycore Processors

Kalray MPPA® Processors

Deep Learning Inference

Model-Based Design

Applications & Outlook

KALRAY

# Intelligent Systems

## Cyber-Physical Systems

- Information processing and physical processes are tightly integrated
- Time constraints associated with information manipulation
- Distributed systems
- Functional safety
- Cyber-security

## Intensive Computing

- Numerical computing
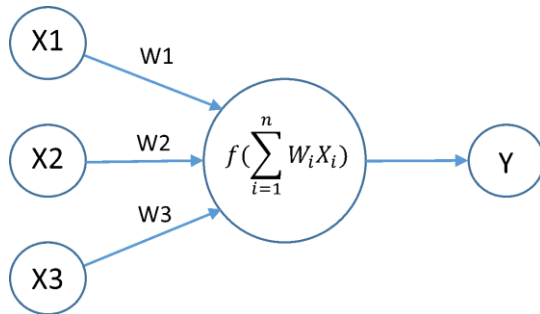- Signal processing
- Image processing
- Graph computing

## Artificial Intelligence

- The science and engineering of creating intelligent machines (J. McCarthy, 1956)
- Mostly represented by the Machine Learning field, in particular Deep Learning
- Association causation level: "the objective of curve-fitting is to maximize fit, while deep learning tries to minimize over-fit" (J. Pearl 2018)

KALRAY

# Machine Learning (ML)

Give computers the ability to learn without being explicitly programmed (Arthur Samuel, 1959)

- **Rule Extraction** Goal is to identify statistical relationships in data
- **Clustering** Group similar data together, while increasing the gap between the groups
- **Classification & Regression** Map a set of new input data to a set of discrete or continuous valued output, respectively
- **Artificial Neural Networks (ANN)** General implementation model for nonlinear classifiers, trained by using back-propagation algorithms
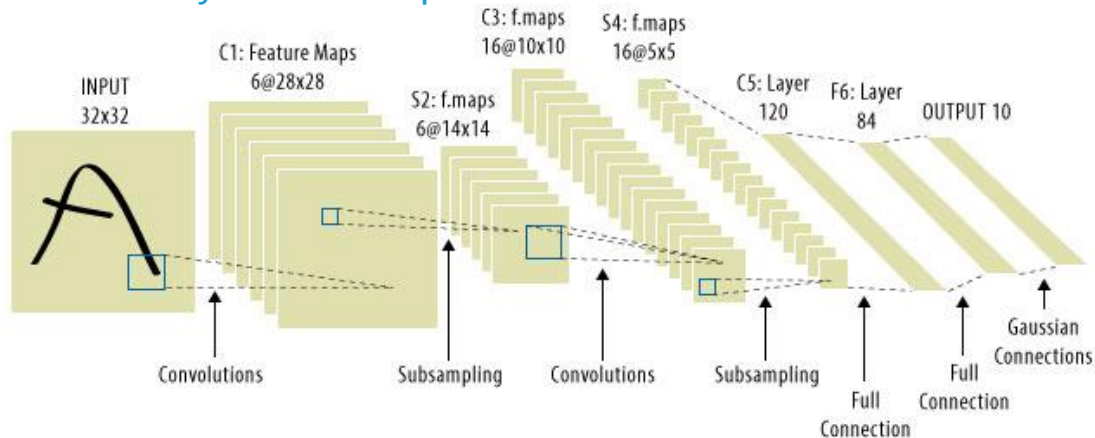


Weighted sum of inputs

$$Y = f\big((W_1 \times X_1) + (W_2 \times X_2) + (W_3 \times X_3)\big)$$

$f(\ )$ is the "activation function"
$f(x) = \max(0, x)$, "Rectified Linear Unit"
$f(x) = \tanh(x)$
…

KALRAY

# Deep Learning (DL)

Computational models composed of multiple processing layers to learn representations of data with multiple levels of abstraction (Yann Le Cun et al., 2015)

- **Convolutional Neural Networks (CNN)** Networks where most filtering operations performed by feature maps are discrete convolutions



- **Recurrent Neural Networks (RNN)** Networks with feedback loops
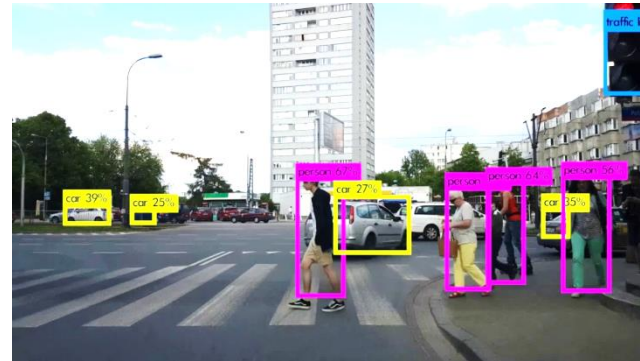
KALRAY

# Machine Learning Steps

## Training (datacenter)

- Learning part or Machine Learning
  - Supervised (classification & regression)
  - Unsupervised (clustering)
  - Reinforcement (decision-making)
- Off-line processing of large data sets
- Floating-point 32-bit arithmetic
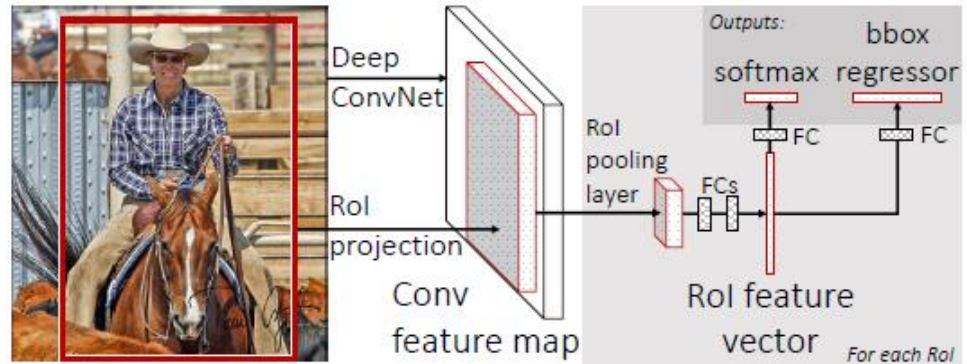
## Inference (intelligent system)

- Classification / Segmentation / Detection
- On-line / Real-time data stream processing
- Floating-point 16.32-bit or "bfloat16" arithmetic
- Integer 8.32 arithmetic (quantization) for CNN

KALRAY

# R-CNN, Fast & Faster R-CNN (Girshick & Ren, 2014-2016)

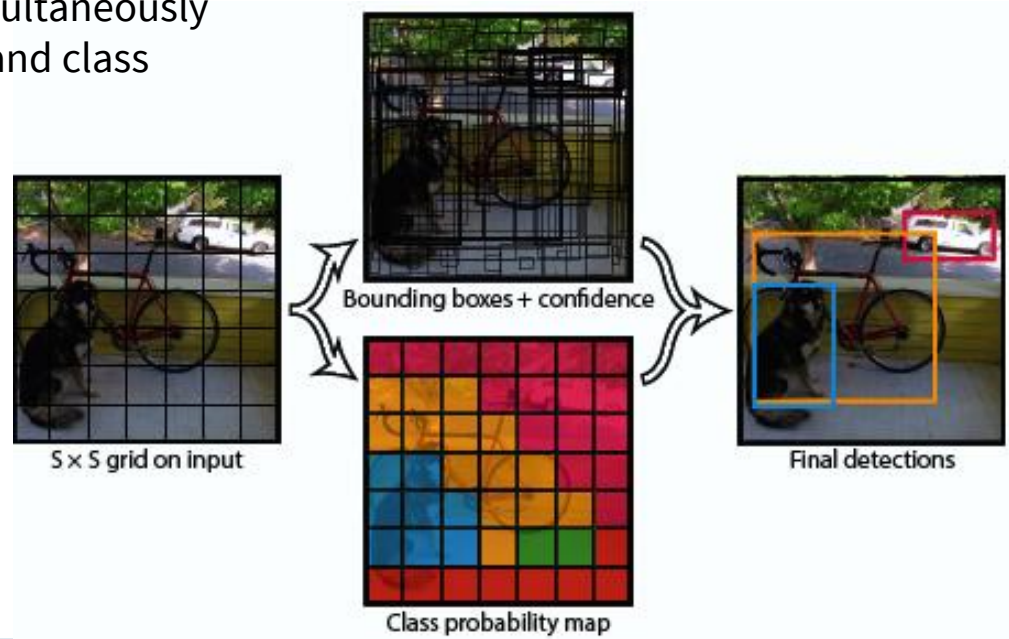Regional CNN and improvements use two steps for object detection

1) Proposal of candidate regions (initially by segmentation, then by neural computing)

2) Classification of candidate regions (neural computing and refinment steps)

# YOLO v1-3 « You Only Look Once » (Redmon 2016-2018)

## Single-step method (unlike « R-CNN » family)

- A single convolutional network simultaneously predicts multiple bounding boxes and class probabilities for those boxes



S × S grid on input

Bounding boxes + confidence

Class probability map

Final detections

KALRAY

# Cyber-Security Requirements

| | Defense | Avionics | Automotive |
|---|:---:|:---:|:---:|
| Hardware root of trust (HSM) | ✓ | ✓ | ✓ |
| Authenticated software | ✓ | ✓ | ✓ |
| Encrypted boot firmware | ✓ | ✓ | |
| Encrypted application code | ✓ | ✓ | ✓ |
| Event data record encryption | | ✓ | ✓ |
| Secured communication | ✓ | ✓ | ✓ |
| Physical attack protection | ✓ | | ✓ |

**KALRAY**

# Secure Boot for Trusted Software Deployment

## Measured boot

- Enables external agent to attest the platform state after the boot process
- Provides a secure measurement and reporting chain to external agent
- Detect modified boot code, settings and boot paths
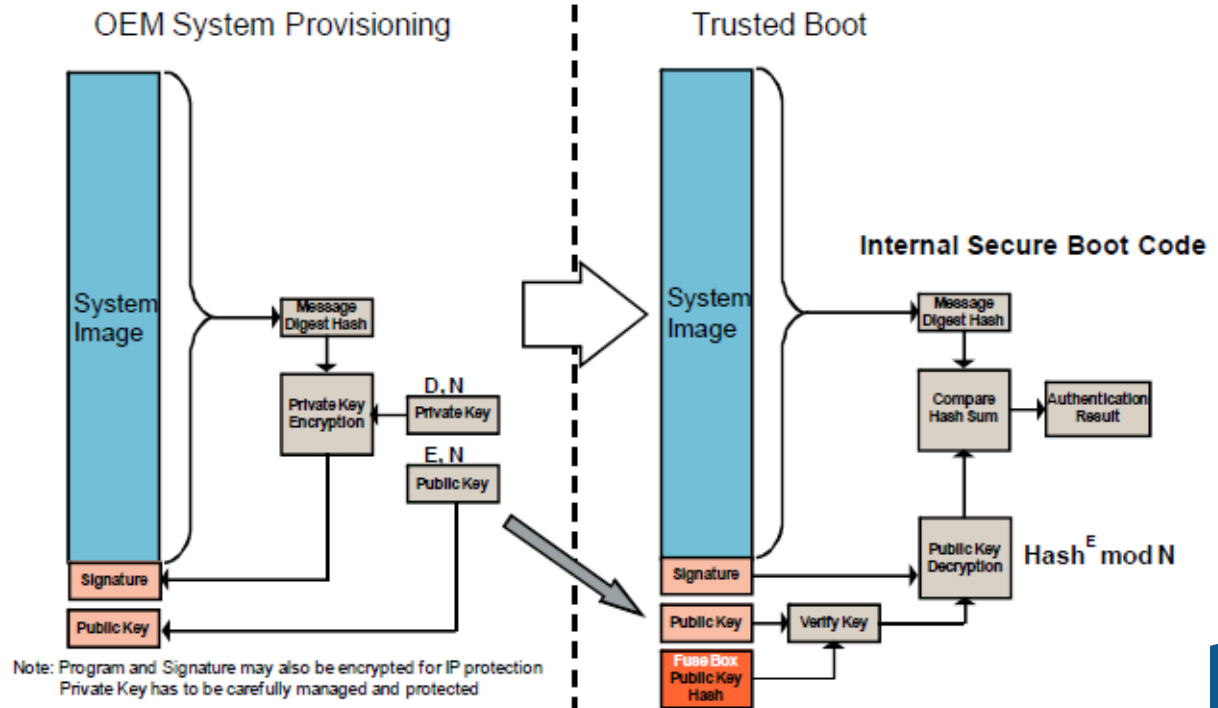- Typically used on servers, by associating UEFI and a TPM chip

## Trusted boot

- Unbroken chain of trust across all stages of the OS boot:
  - Phase-0 (internal ROM) check platform, validate & launch Phase-1
  - Phase-1 (Flash) initialize peripherals, validate & launch Phase-2
  - Phase-2 (network or disk), validate & launch operating system
- Typically used on embedded systems

**KALRAY**

# FreeScale QorIQ Trusted Boot

## QorIQ processors target consumer, industrial, medical, networking

- OEM public keys and the intent to secure (ITS) bit in immutable storage (fuses)
- When the ITS bit is set, jump to internal boot ROM (IBR) for Phase-0
- Phase-1 firmware digitally signed using OEM private signature key
- Phase-0 verifies firmware signature using public key



Note: Program and Signature may also be encrypted for IP protection
Private Key has to be carefully managed and protected

KALRAY

# Outline

Intelligent Systems

**Manycore Processors**

Kalray MPPA® Processors

Deep Learning Inference

Model-Based Design

Applications & Outlook

KALRAY

# Homogeneous Multicore Processor

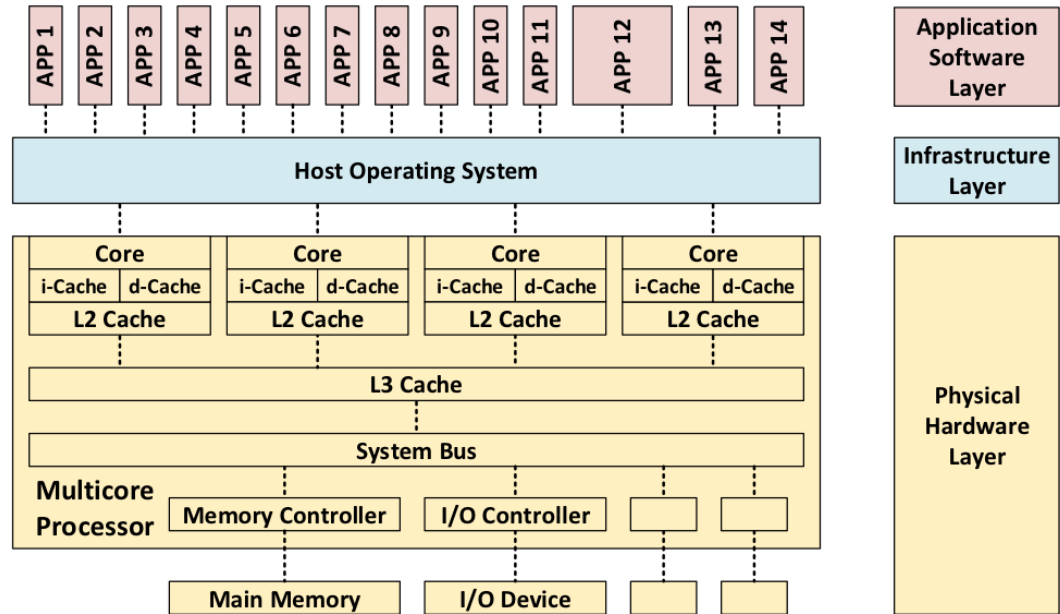## Multiple cores sharing a cache-coherent memory hierarchy

- Private L1 i-cache and d-cache
- Shared or clustered L2 cache
- Shared L3 cache

## Application programming

- C/C++, Python, Java
- Pthreads, std::thread, OpenMP
- Rich operating system (Linux)

## Application partitioning

- Virtual machine monitor



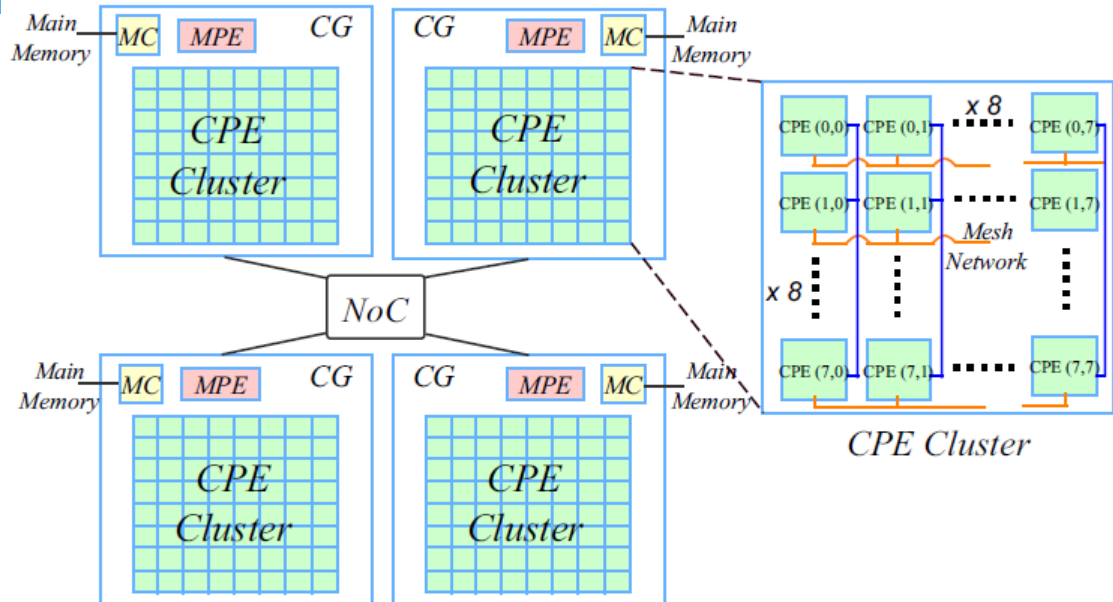https://insights.sei.cmu.edu/sei_blog/2017/08/multicore-processing.html

KALRAY

# Manycore Processors

## Multiple 'Compute Units' connected by a network-on-chip (NoC)

- Group of cores + DMA engine
- Scratch-pad memory (SPM)
- Software-managed caches
- Local cache coherency

## SW26010 Manycore processor

- Node of the Sunway TaihuLight supercomputer (#1 TOP 500 in 2016)
- 4 'core groups' with MPE core, CPE core cluster, collective DMA engine
- 64KB SPM per CPE core



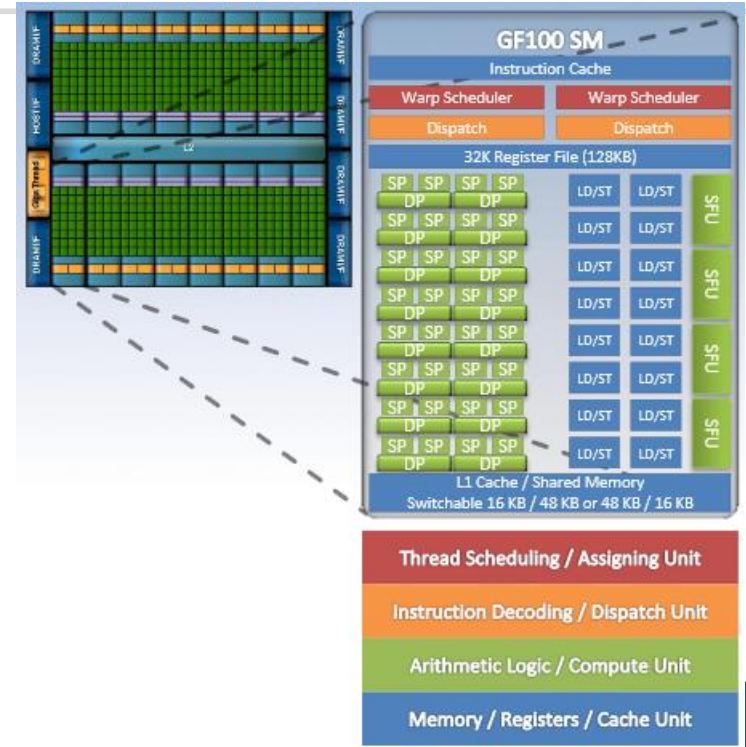Z. Xu, J. Lin, S. Matsuoka, «Benchmarking SW26010 Many-Core Processor» IPDPS 2017

**KALRAY**

# GPGPUs as Manycore Processors (NVIDIA)

## Classic GPGPU: NVidia Fermi architecture

- GPGPU 'compute units' are called Streaming Multiprocessors (SM)
- Each SM comprises 32 'streaming cores' or 'CUDA cores' that share a local memory, caches and a global memory hierarchy
- Threads are scheduled and executed atomically by 'warps', which execute the same instruction or are inactive at any given time
- Hardware multithreading enables warp execution switching on each cycle, helping cover memory access latencies

## GPGPU programming models (CUDA, OpenCL)

- Each SM executes 'thread blocks', whose threads may share data in the local memory and access a common memory hierarchy
- Synchronization inside a thread block by barriers, local memory accesses, atomic operations, or shuffle operations (NVIDIA)
- Synchronization between thread blocks through host program or global memory atomic operations in kernels

**KALRAY**

# GPGPU Tensor Cores for Deep Learning (NVIDIA)

## NVidia Volta architecture
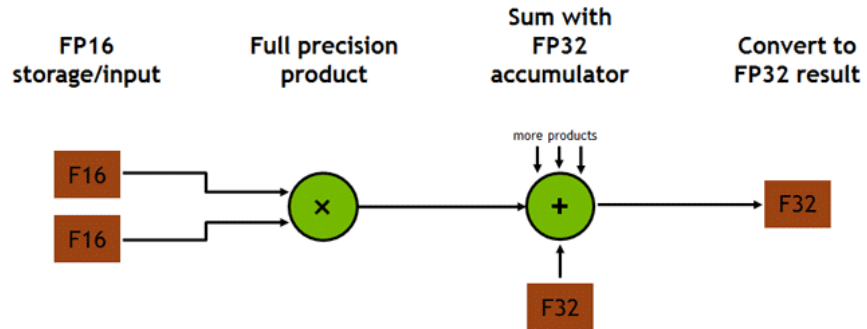
- 64x FP32 cores per SM
- 32x FP64 cores per SM
- 8x Tensor cores per SM

## Tensor core operations

- Tensor Core perform D = A x B + C, where A, B, C and D are matrices
- A and B are FP16 4x4 matrices
- D and C can be either FP16 or FP32 4x4 matrices
- Higher performance is achieved when A and B dimensions are multiples of 8
- Maximum of 64 floating-point mixed-precision FMA operations per clock

KALRAY

# Limitations of GPGPUs

## Restrictions of GPGPU programming

- CUDA is a proprietary programming environment
- Writing OpenCL programs implies writing host code and device code, then connecting them through a low-level API
- GPGPU kernel programming lacks standard features of C/C++, such as recursion or accessing a (virtual) file system
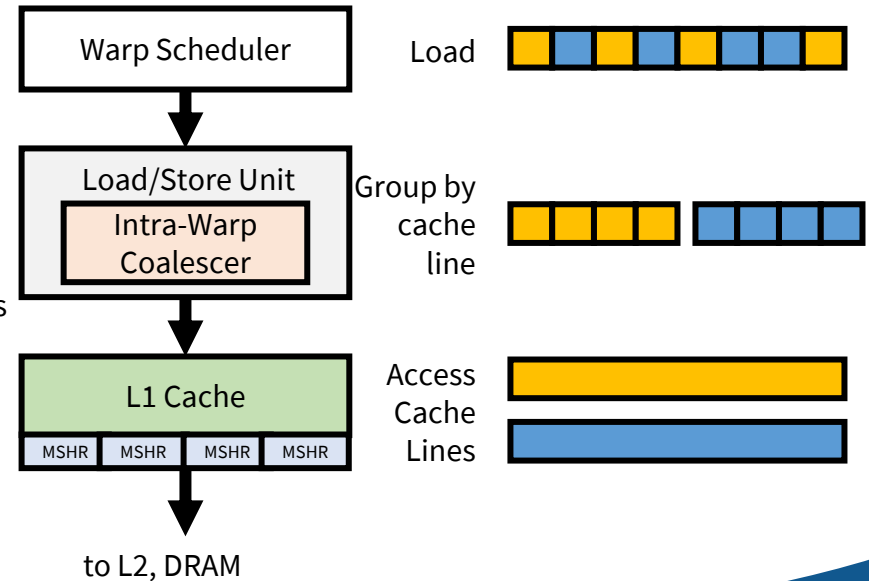
## Performance issues with 'thread divergence'

- Branch divergence: if...then...else construct will force all threads in a warp to execute both the "then" and the "else" path
- Memory divergence: when hardware cannot coalesce the set of warp global memory accesses into one or two L1 cache blocks
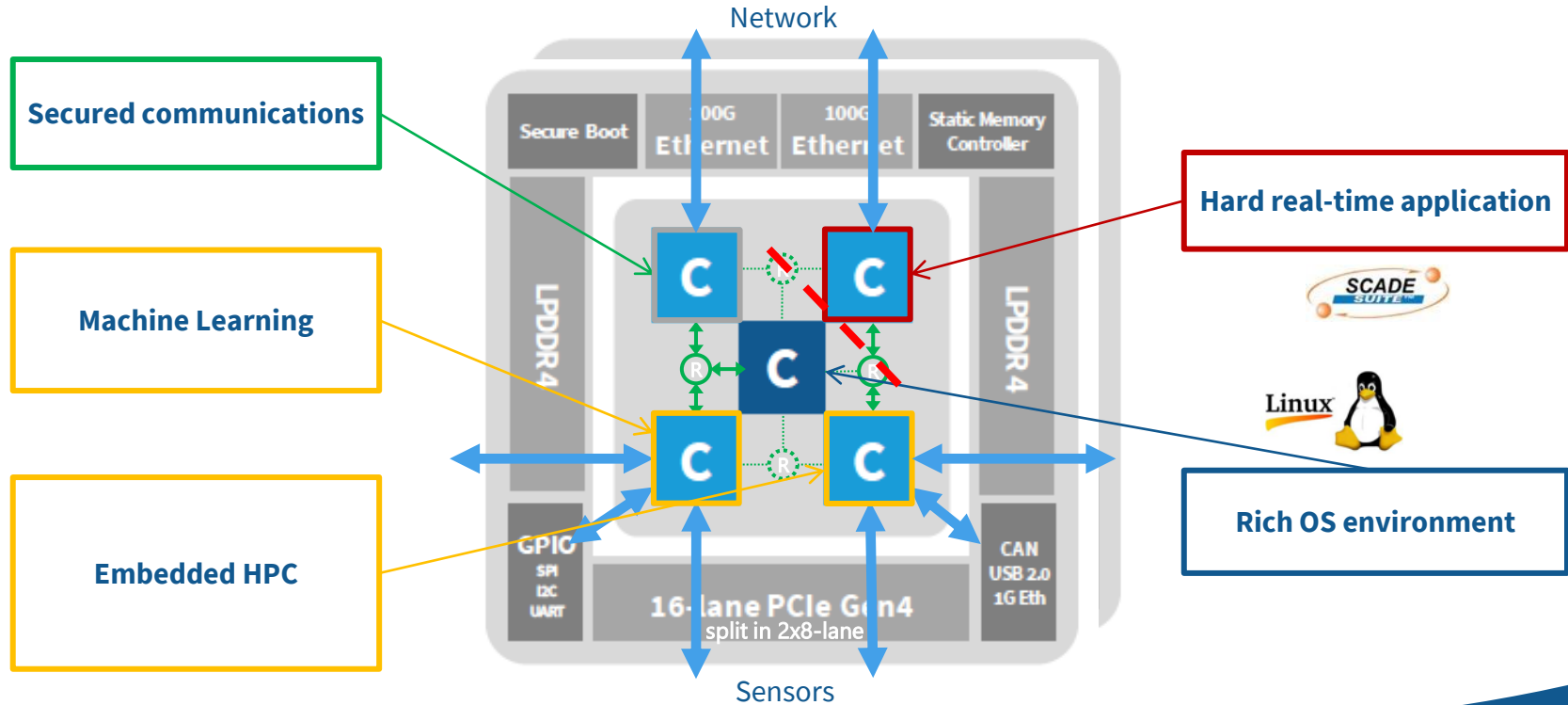
## Time-predictability issues

- Dynamic allocation of thread blocks to SMs
- Dynamic warp scheduling and out-of-order execution on a SM

## Memory access coalescing (Kloosterman et al.)



Warp Scheduler

Load/Store Unit
Intra-Warp Coalescer

L1 Cache
MSHR  MSHR  MSHR  MSHR

to L2, DRAM

Load

Group by cache line

Access Cache Lines

**KALRAY**

# Mapping Intelligent System Functions to Compute Units



©2019 – Kalray SA All Rights Reserved

# Outline

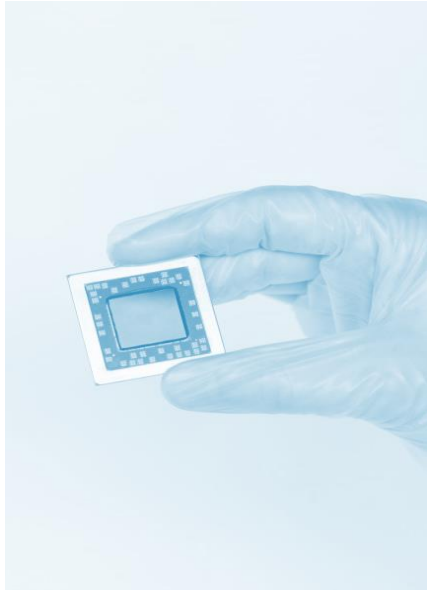Intelligent Systems

Manycore Processors

**Kalray MPPA® Processors**

Deep Learning Inference

Model-Based Design

Applications & Outlook

KALRAY

# Kalray's MPPA® Manycore Architecture

## MPPA® (Massively Parallel Processor Array) Platform

**Hardware**

Manycore CPU architecture

*Compute clusters of 16 high-performance CPU cores with local memory*

DSP-like timing predictability

*'Fully timing compositional' cores for accurate static timing analysis*
*Service guarantees of local memory system and network-on-chip*

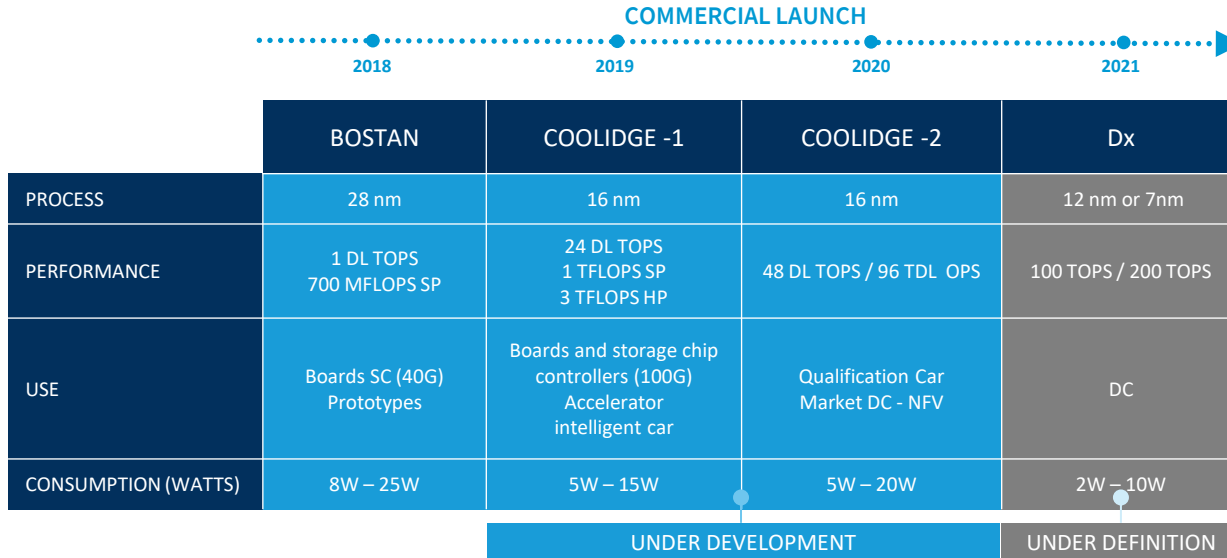FPGA-like I/O capabilities

**Software**

CPU programming

*Standard C/C++/OpenMP/OpenCL, OpenVX*
*Library code generators (MetaLibm, KaNN)*
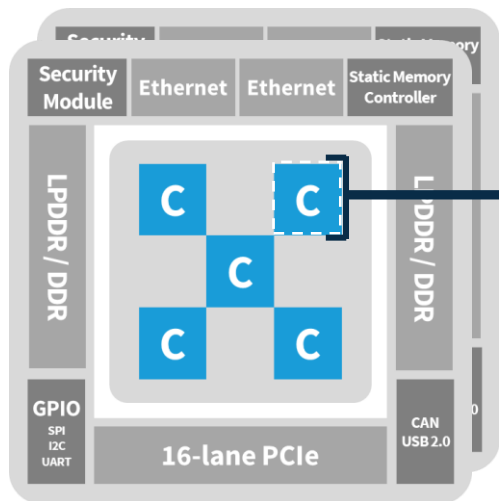*Model-based (SCADE Suite®, Simulink®)*

**KALRAY**

# MPPA® Processor Family and Roadmap

**COMMERCIAL LAUNCH**

2018     2019     2020     2021

|  | BOSTAN | COOLIDGE -1 | COOLIDGE -2 | Dx |
|---|---|---|---|---|
| PROCESS | 28 nm | 16 nm | 16 nm | 12 nm or 7nm |
| PERFORMANCE | 1 DL TOPS<br>700 MFLOPS SP | 24 DL TOPS<br>1 TFLOPS SP<br>3 TFLOPS HP | 48 DL TOPS / 96 TDL OPS | 100 TOPS / 200 TOPS |
| USE | Boards SC (40G)<br>Prototypes | Boards and storage chip controllers (100G)<br>Accelerator<br>intelligent car | Qualification Car<br>Market DC - NFV | DC |
| CONSUMPTION (WATTS) | 8W – 25W | 5W – 15W | 5W – 20W | 2W – 10W |
|  |  | UNDER DEVELOPMENT | | UNDER DEFINITION |

**MANYCORE TECHNOLOGY** THAT ENABLES PROCESSOR OPTIMIZATION
BASED ON EVOLVING MARKET REQUIREMENTS
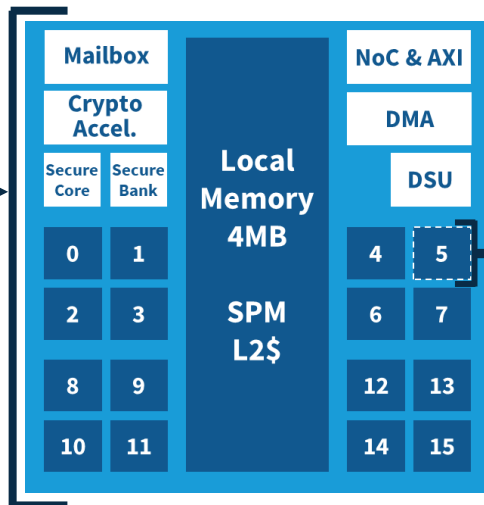
KALRAY

# MPPA3-80 Processor (TSMC 16FFC, 1.2GHz)
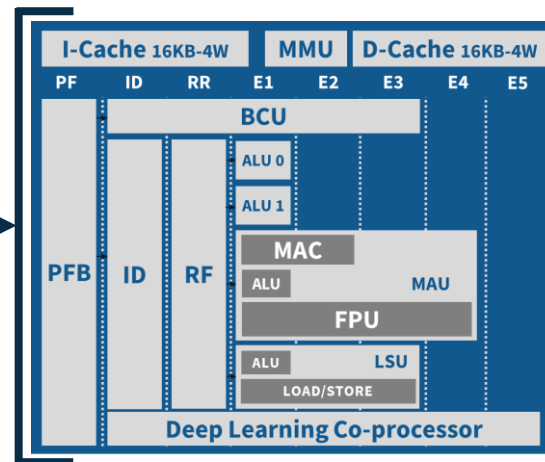## 1TFLOPS FP32, 3TFLOPS FP16.32, 24 DLTOPS INT8.32



**COOLIDGE PROCESSOR**

5 compute clusters at 1200 MHz
2x 100Gbps Ethernet, 16x PCIe Gen4

**COMPUTE CLUSTER**

16+1 cores, 4 MB local memory
NoC and AXI global interconnects

**6-ISSUE VLIW CORE**

64x 64-bit register file
128MAC/c tensor coprocessor

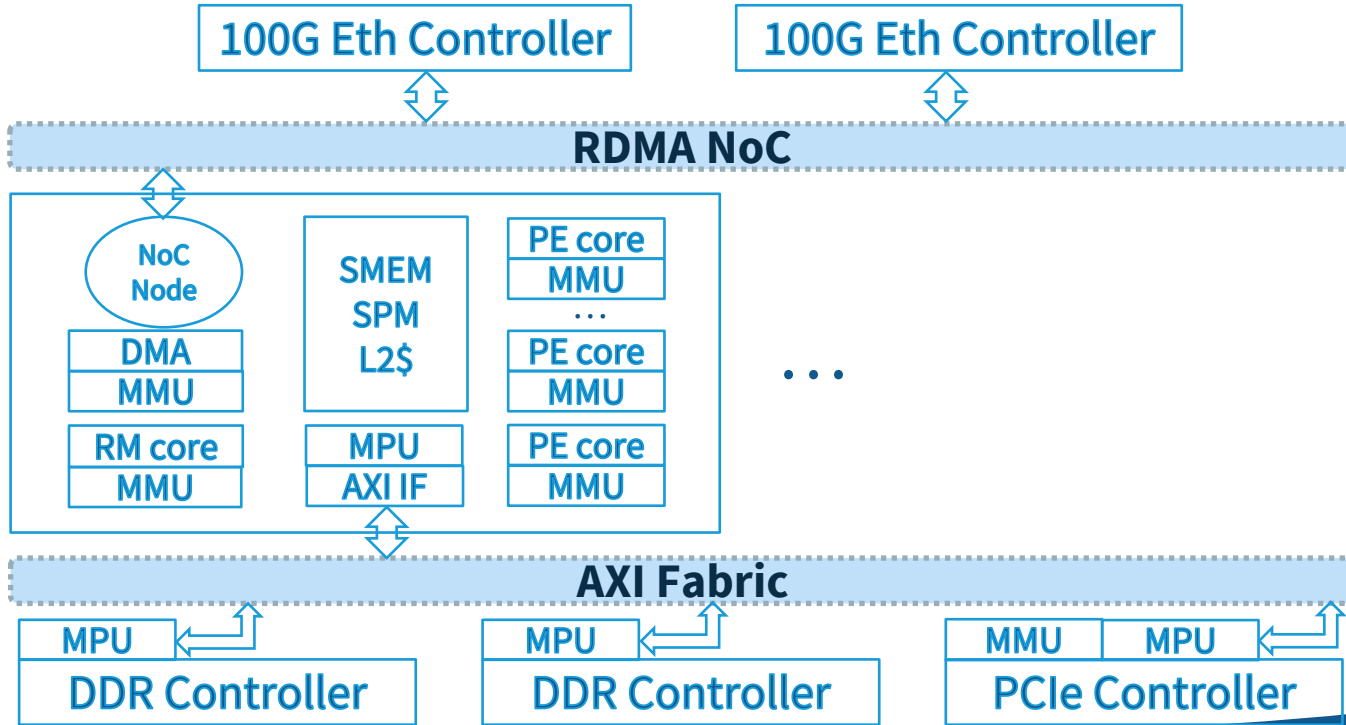KALRAY

# Network-on-Chip for Global Interconnects

## NoC as generalization of busses

- Connectionless
- Address-based transactions
- Flit-level flow control
- Implicit packet routing
- Inside a coherence domain
- Reliable communication
- Coherency protocol messages
- Coordinate with DDR memory controller front-end (Ex. Arteris FlexMem Memory Scheduler)

## NoC as integrated macro-network

- Connection-oriented
- Stream-based transactions
- [End-to-end flow control]
- Explicit packet routing
- Across address spaces (RDMA)
- [Packet loss or packet reordering]
- Traffic shaping for QoS (application of DNC)
- Terminate macro-network (Ethernet, InfiniBand)
- Support of multicasting

KALRAY

# MPPA3 Global Interconnects
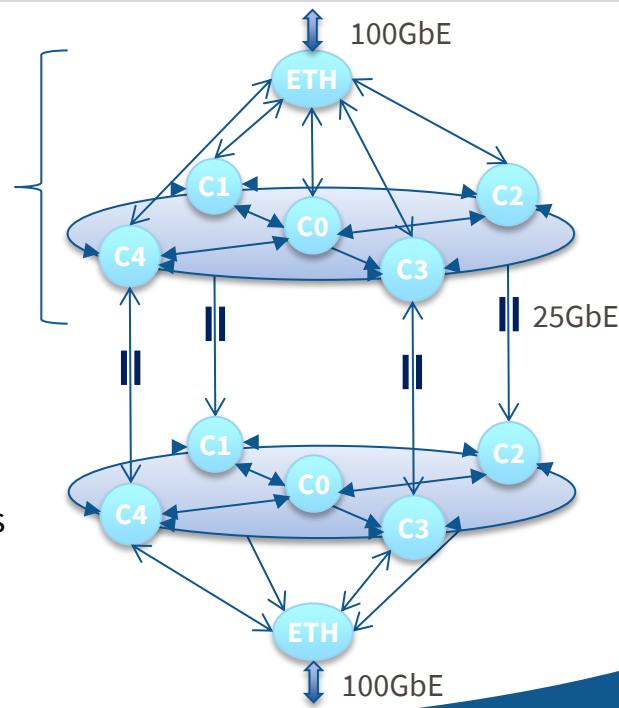
# MPPA3 RDMA NoC

## MPPA3 NoC architecture

- Wormhole switching with source routing
- 2 virtual channels, 4x TX DMA channels
- RDMA, remote queues, remote atomics
- 128-bit flits, up to 17 flits/packet (256B payload)

## 4x 25Gbps Ethernet lanes reused for NoC extension

- NoC packet encapsulation into IEEE 802.1Q standard for VLAN
- Designed for direct connections between 2 to 4 chips (using FEC)
- VCs map to IEEE 802.1Qbb Priority-based Flow Control (PFC) classes

MPPA3-80 Processor

100GbE

25GbE

100GbE

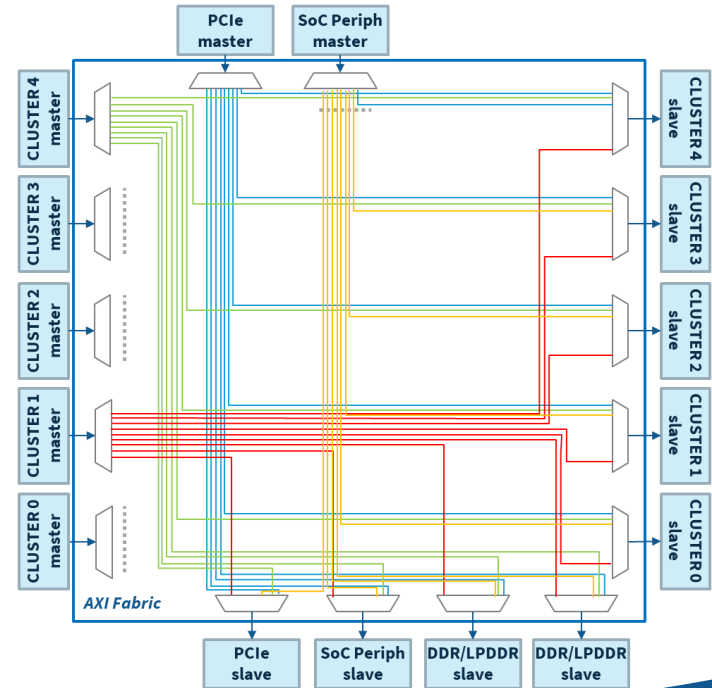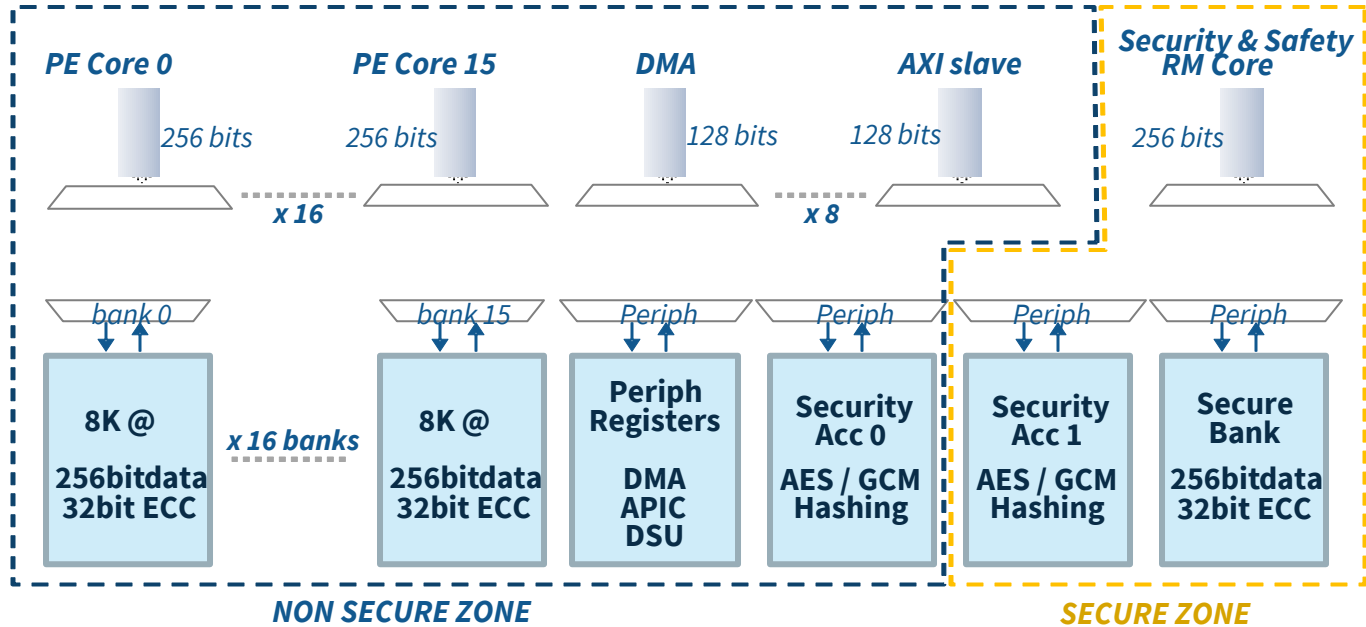| MAC dst 6 bytes | MAC src 6 bytes | VLAN etype 0x8100 2 bytes | VLAN TCI PFC (3 bits) / CFI (1 bit) NoC pkt nb (12 bits) 2 bytes | NoCX etype 0xB000 2 bytes | NoC pkt0 | NoC pkt1 | FCS 4 bytes |
|---|---|---|---|---|---|---|---|

KALRAY

# MPPA3 AXI Fabric

## Deficit Round-Robin (DRR) Arbitration

- Assing a 'quantum' of flits $Q_1 \dots Q_n$ to each input
- Associate a 'deficit counter' in flits $DC_1 \dots DC_n$ to each input
- Iterate on the non-empty inputs; for each input $i$:

1. $DC_i$ += $Q_i$
2. Transfer packets to output while cumulative flit count $\leq DC_i$
3. $DC_i$ -= transferred cumulative flit count
4. $DC_i$ := 0 if input is empty
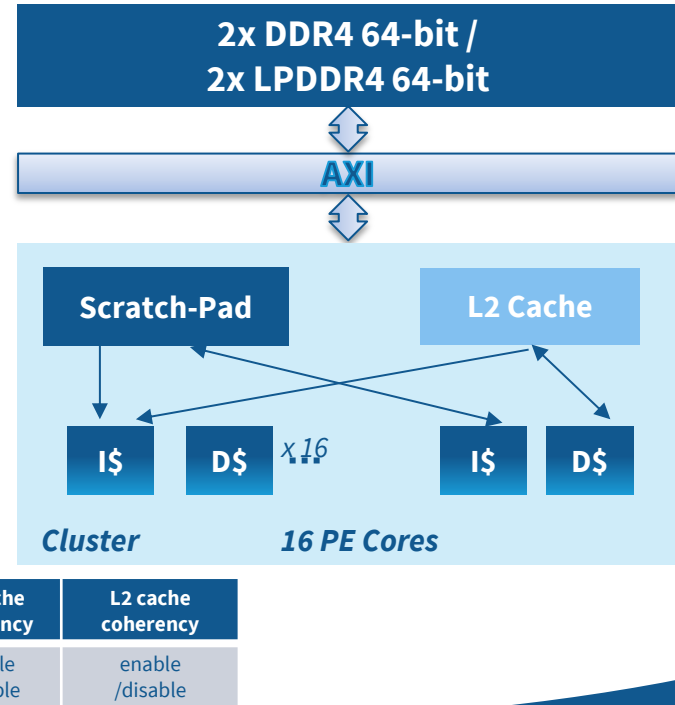
KALRAY

# MPPA3 Compute Cluster

# MPPA3 Memory Hierarchy

## VLIW Core L1 Caches

- 16KB / 4-way LRU instruction cache per core
- 16KB / 4-way LRU data cache per core
- 64B cache line size
- Write-through, write no-allocate (write around)
- Coherency configurable across all L1 data caches

## Cluster L2 Cache & Scratch-Pad Memory

- **Scratch-pad from 2MB to 4MB**
  - **16 independent banks, full crossbar**
  - **Interleaved or banked address mapping**
- L2 cache from 0MB to 2MB
  - 16-way Set Associative
  - 256B cache line size
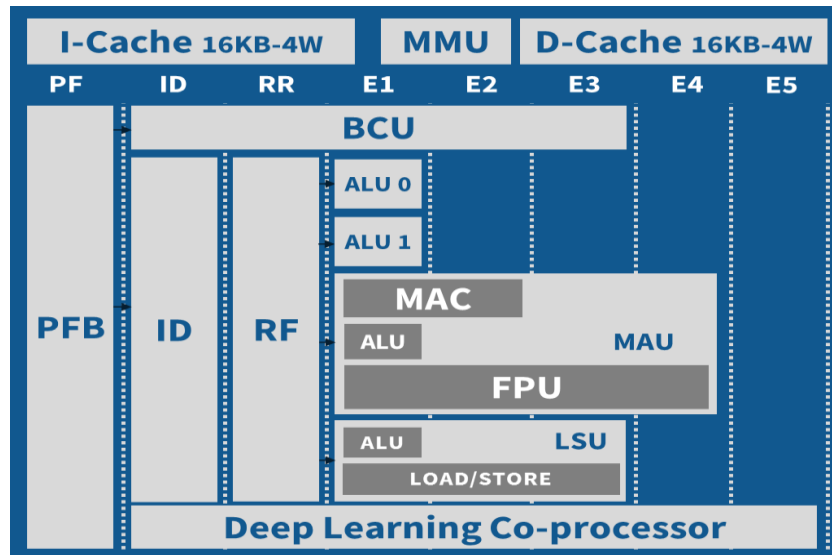  - Write-back, write allocate



| 2x DDR4 64-bit / 2x LPDDR4 64-bit |
| --- |

**AXI**

Scratch-Pad    L2 Cache

I\$    D\$    *x16*    I\$    D\$

*Cluster*    *16 PE Cores*

| L1 cache coherency | L2 cache coherency |
| --- | --- |
| enable /disable | enable /disable |

**KALRAY**

# MPPA3 64-Bit VLIW Core

## Vector-scalar ISA

- 64x 64-bit general-purpose registers
- Operands can be single registers, register pairs (128-bit) or register quadruples (256-bit)
- Immediate operands up to 64-bit, including F.P.
- 128-bit SIMD instructions by dual-issuing 64-bit on the two ALUS or by using the FPU datapath

## FPU capabilities

- 64-bit x 64-bit + 128-bit → 128-bit
- 128-bit op 128-bit → 128-bit
- FP16x4 SIMD 16 x 16 + 32 → 32
- FP32x2 FMA, FP32x4 FADD, FP32 FMUL Complex
- FP32 Matrix Multiply 2x2 Accumulate



K1C VLIW CORE PIPELINE

KALRAY

# MPPA3 Tensor Coprocessor

## Extend VLIW core ISA with extra issue lanes

- Separate 48x 256-bit wide vector register file
- Matrix-oriented arithmetic operations (CNN, CV …)
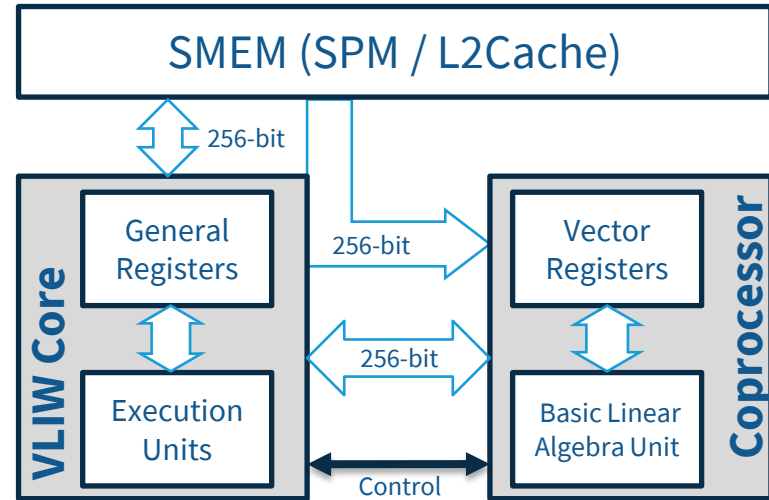
## Full integration into core instruction pipeline

- Move instructions supporting matrix-transpose
- Proper dependency / cancel management

## Leverage MPPA memory hierarchy

- SMEM directly accessible from coprocessor
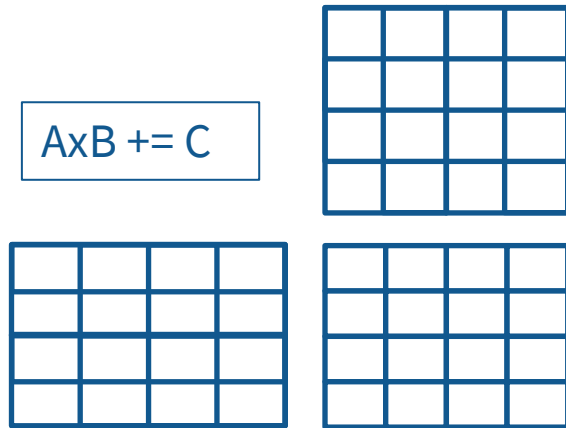- Memory load stream aligment operations

## Arithmetic performances

- 128x INT8→INT32 MAC/cycle
- 64x INT16→INT64 MAC/cycle
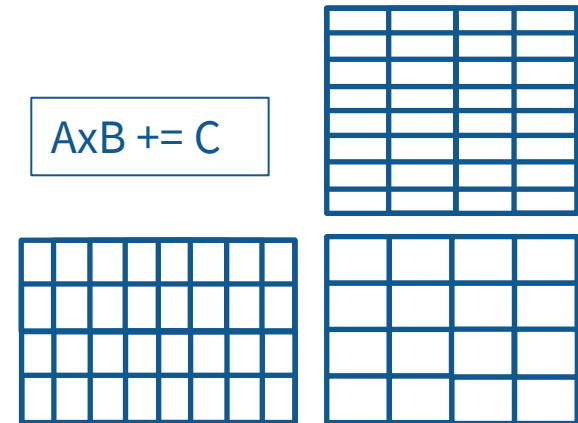- 16x FP16→FP32 FMA/cycle

# MPPA3 Coprocessor Matrix Operations

- INT16 to INT64 convolutions:

$$(4x4)_{int16} \cdot (4x4)_{int16} \mathrel{+}= (4x4)_{int64}$$

- INT8 to INT32 convolutions

$$(4x8)_{int8} \cdot (8x4)_{int8} \mathrel{+}= (4x4)_{int32}$$

AxB += C

AxB += C

KALRAY

# Outline

Intelligent Systems
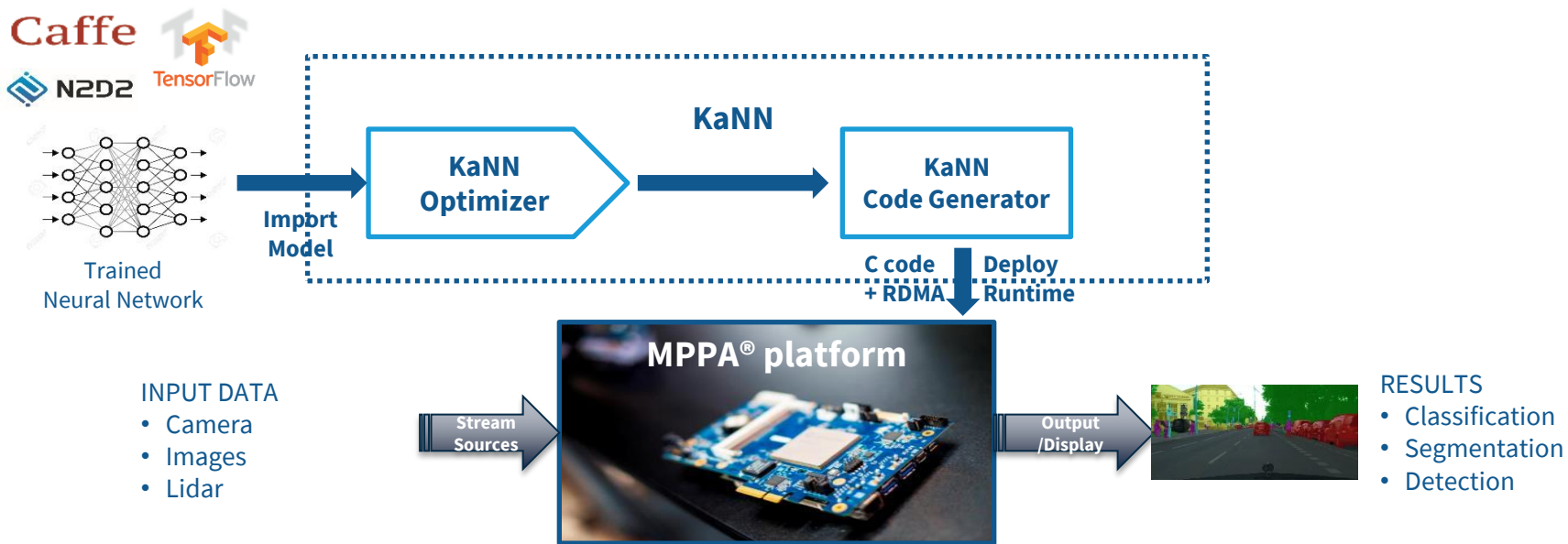
Manycore Processors

Kalray MPPA® Processors

**Deep Learning Inference**

Model-Based Design

Applications & Outlook

KALRAY

# KaNN (Kalray Neural Network) Inference Code Generator



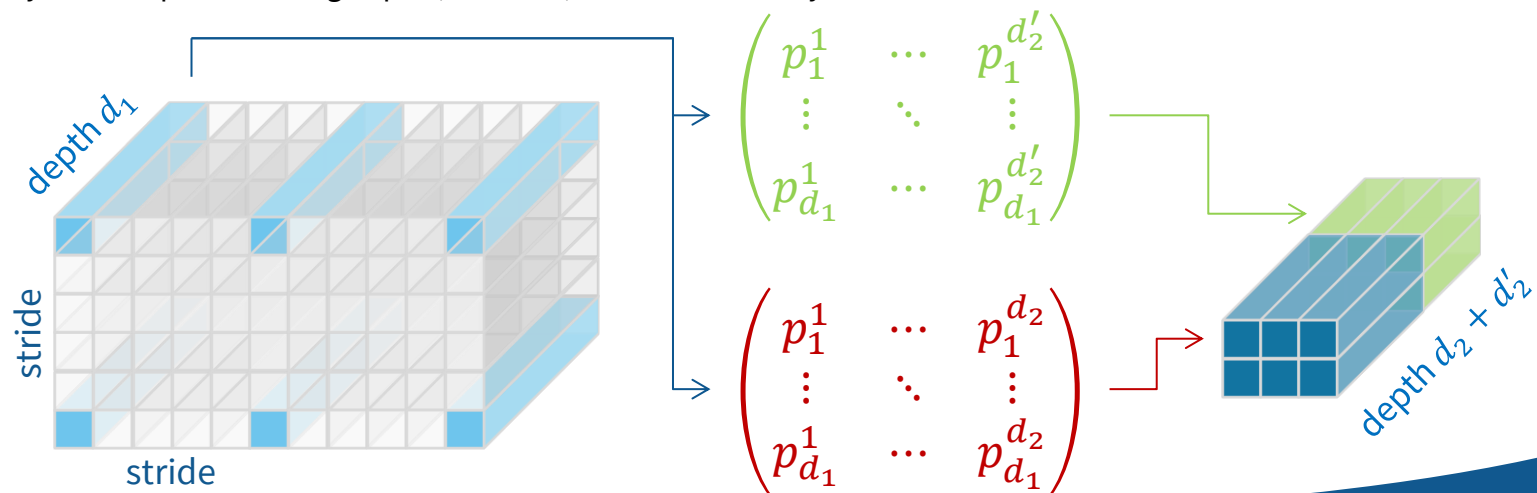©2018 – Kalray SA All Rights Reserved

# CNN Inference on a MPPA Processor (1)

Compute one DNN layer at a time in topological sort order of the network

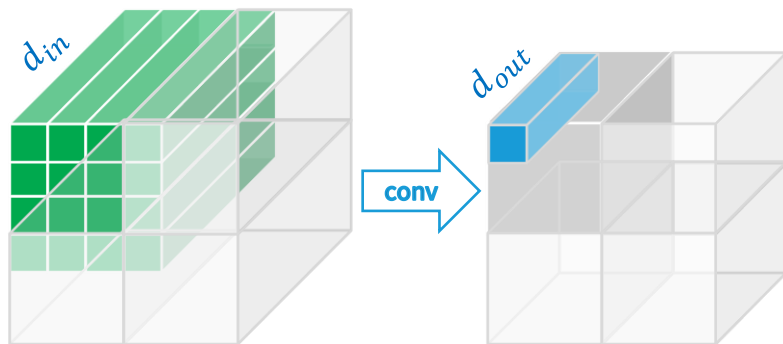Decompose NxN convolutions as accumulations of $N^2$ 1x1 convolutions

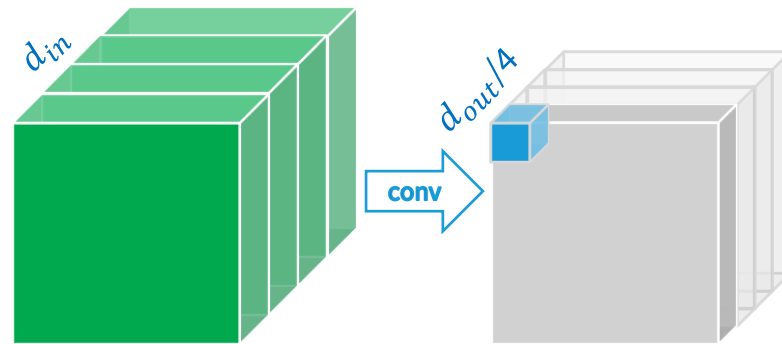- Pixels layout is sequential along depth (channels) for dense memory accesses

**KALRAY**

# CNN Inference on a MPPA Processor (2)

## Distribute activations across clusters SPMs, splitting along spatial and/or depth dimensions

- Spatial dimension splitting requires that the full set of parameters be loaded from external memory
- Channel dimension splitting requires access to the whole input image and a subset of the parameters
- Leverage NoC multicasting of parameters from external memory in case of spatial dimension splitting
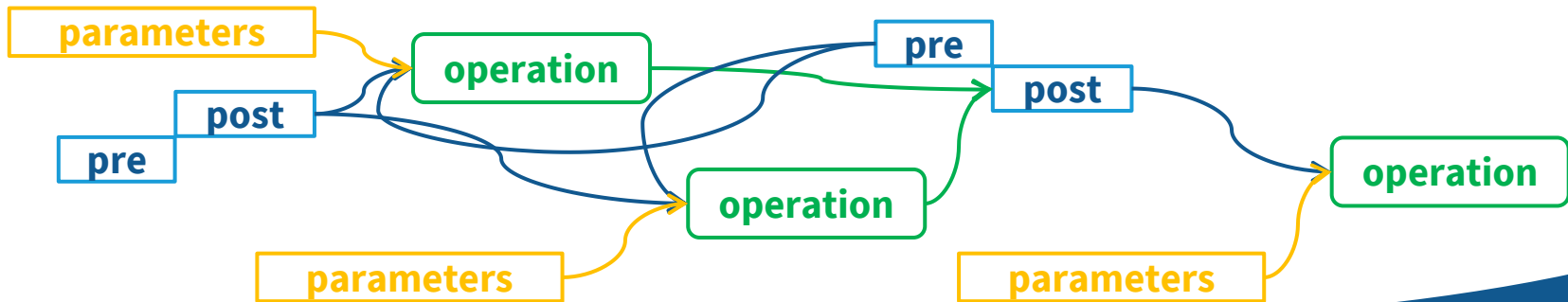


$$[3][3][d_{in}][d_{out}]$$

$$[3][3][d_{in}][d_{out}/4]$$

KALRAY

# CNN Inference on a MPPA Processor (3)

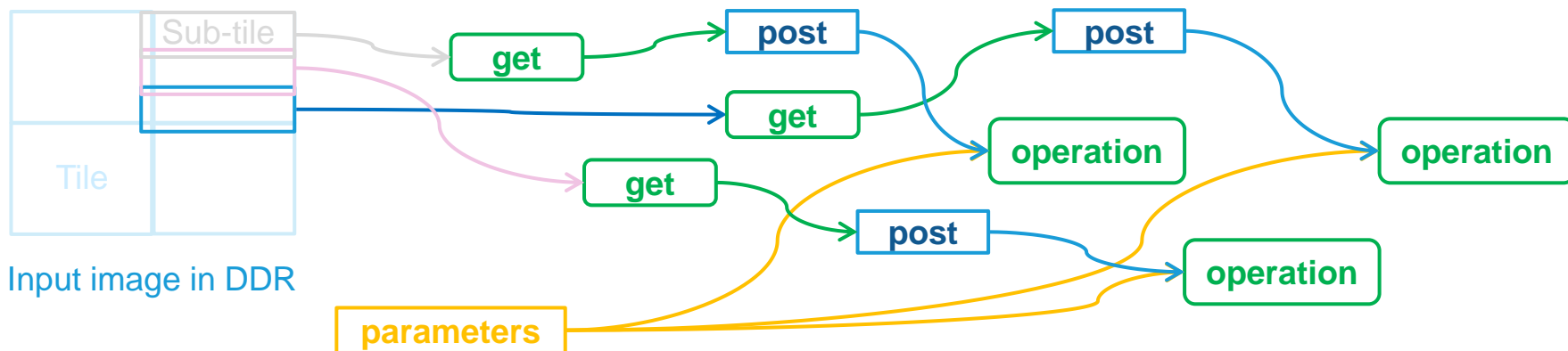## SPMD (Single Program Multiple Data) execution that leverages NoC multicasting of parameters

- Build a local memory buffer allocation and task execution schedule in each cluster
- Overlap parameter transfers from external memory with computations on local memory
- Allocation and scheduling are performed on the CNN network
  - an image corresponds to pre and post tasks,
  - layer compute operations corresponds to a malleable task
  - pre tasks load biases from external memory into the local memory buffer

KALRAY

# CNN Inference on a MPPA Processor (4)
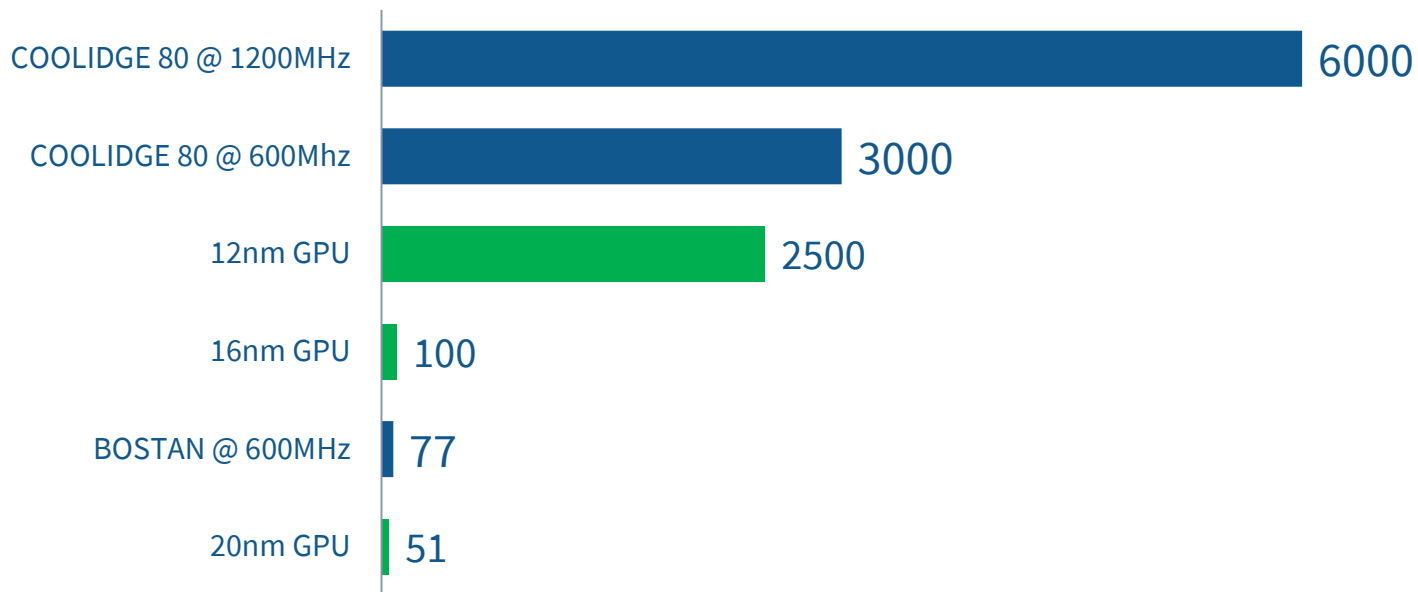
## For layers where images do not fit on-chip, stream sub-tiles from DDR memory

- All clusters remote write their tile of output image to DDR memory, then enter a synchronization barrier
- After clusters leave the barrier, they pipeline the remote read from DDR / operate / put to DDR of sub-tiles
- Larger sub-tiles factor more control overhead but reduce the amount of pipelining

# Deep Learning Inference on Caffe GoogLeNet

## Batch 1 performances in Frames per Second (FPS)



| | FPS |
|---|---|
| COOLIDGE 80 @ 1200MHz | 6000 |
| COOLIDGE 80 @ 600Mhz | 3000 |
| 12nm GPU | 2500 |
| 16nm GPU | 100 |
| BOSTAN @ 600MHz | 77 |
| 20nm GPU | 51 |

KALRAY

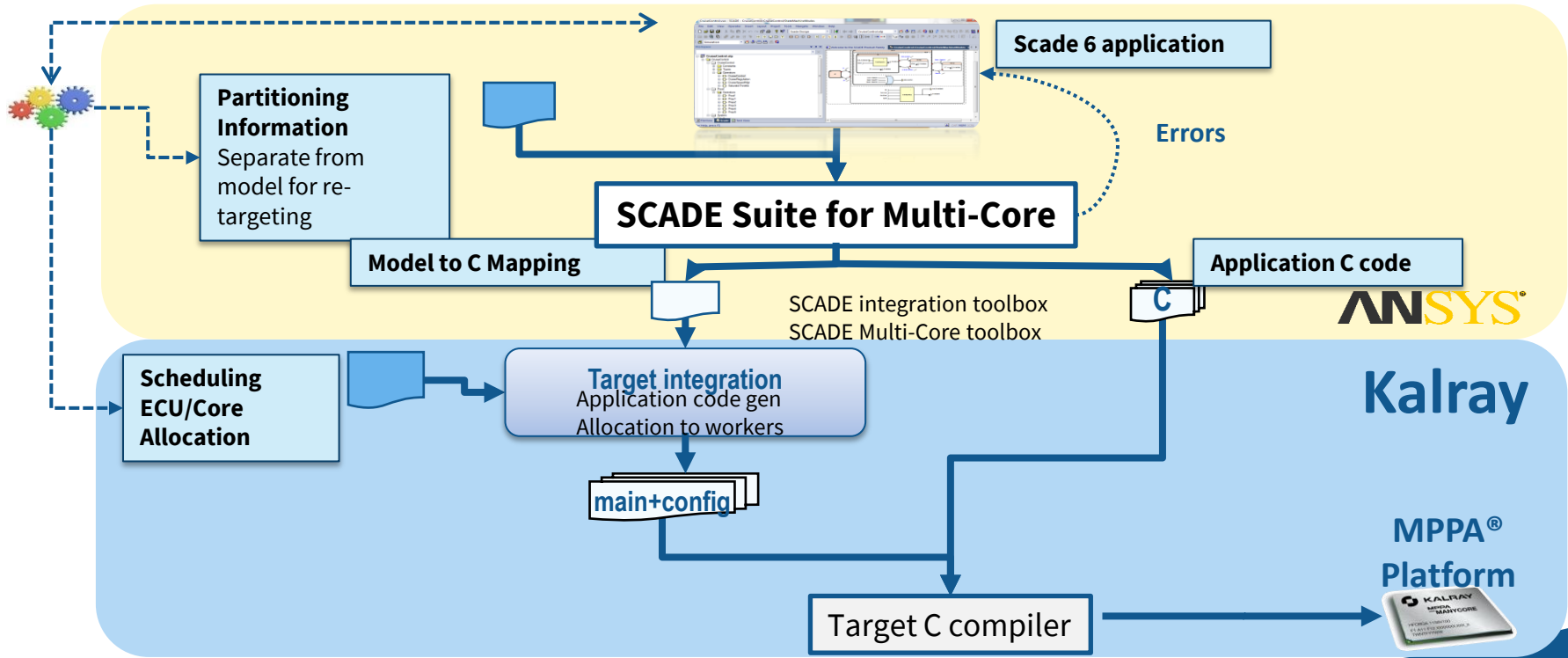# Outline

Intelligent Systems

Manycore Processors

Kalray MPPA® Processors

Deep Learning Inference

**Model-Based Design**

Applications & Outlook
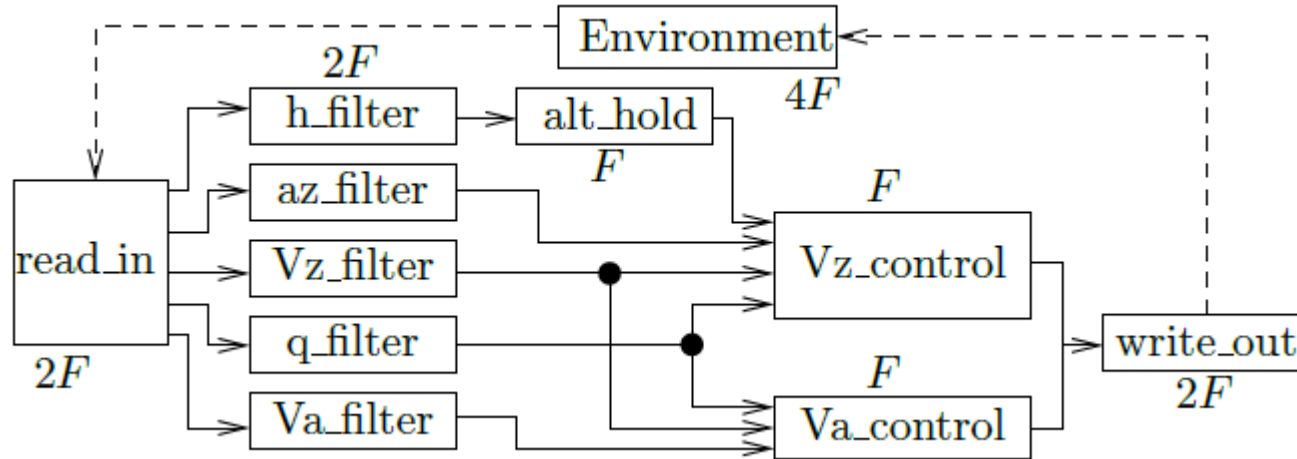
KALRAY

# SCADE Suite Multi-Core Code Generation Flow



**Partitioning Information**
Separate from model for re-targeting

**Scade 6 application**

**Errors**

**SCADE Suite for Multi-Core**

**Model to C Mapping**

**Application C code**

SCADE integration toolbox
SCADE Multi-Core toolbox

**ANSYS**®

**Scheduling ECU/Core Allocation**

**Target integration**
Application code gen
Allocation to workers

**Kalray**

**main+config**

**MPPA® Platform**

Target C compiler

**KALRAY**

# ROSACE Demonstration Application

- Simplified controller for the longitudinal motion of a medium-range civil aircraft in en-route phase: cruise and change of cruise level sub-phases
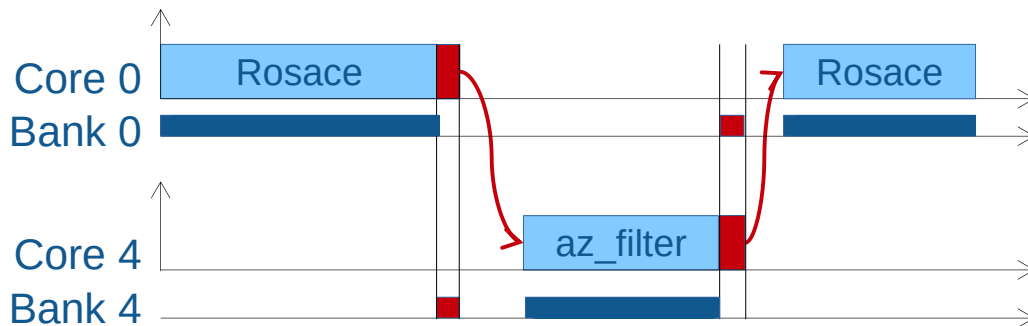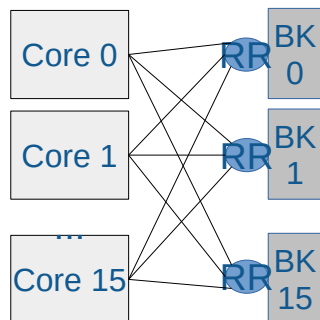


- Original application has 3 harmonic periods: F, 2F, 4F

KALRAY

# SCADE Suite MCG Code Generation (1)

- MCG generates a set of tasks communicating trough one-to-one channels:
    - The root task executes the root operator of the input model
    - One task for each operator instance annotated in the input model
    - Each task receives data on an input channel, calls the operator and then sends the result on an output channel
    - Channels are single-producer, single-consumer FIFOs of size one

- The platform provider (Kalray) integrates MCG generated code by:
    - Providing *workers*, each able to execute sequentially a set of tasks
    - Implementing communication *channels* with their *send/recv* methods
    - Applying the prescribed scheduling and mapping of tasks to workers

**KALRAY**

# SCADE Suite MCG Code Generation (2)

- ## Exploit the MPPA cluster configuration for 'high-integrity' execution
  - ### Enable the cluster local memory mapping of one bank per core



- ## Precisely compute the task WCETs (Worst-Case Execution Times)
  - ### Static analysis or measurement for the WCET of tasks in isolation
  - ### Refine the WCET with interferences using fixed-point [Rihani RTNS'16]
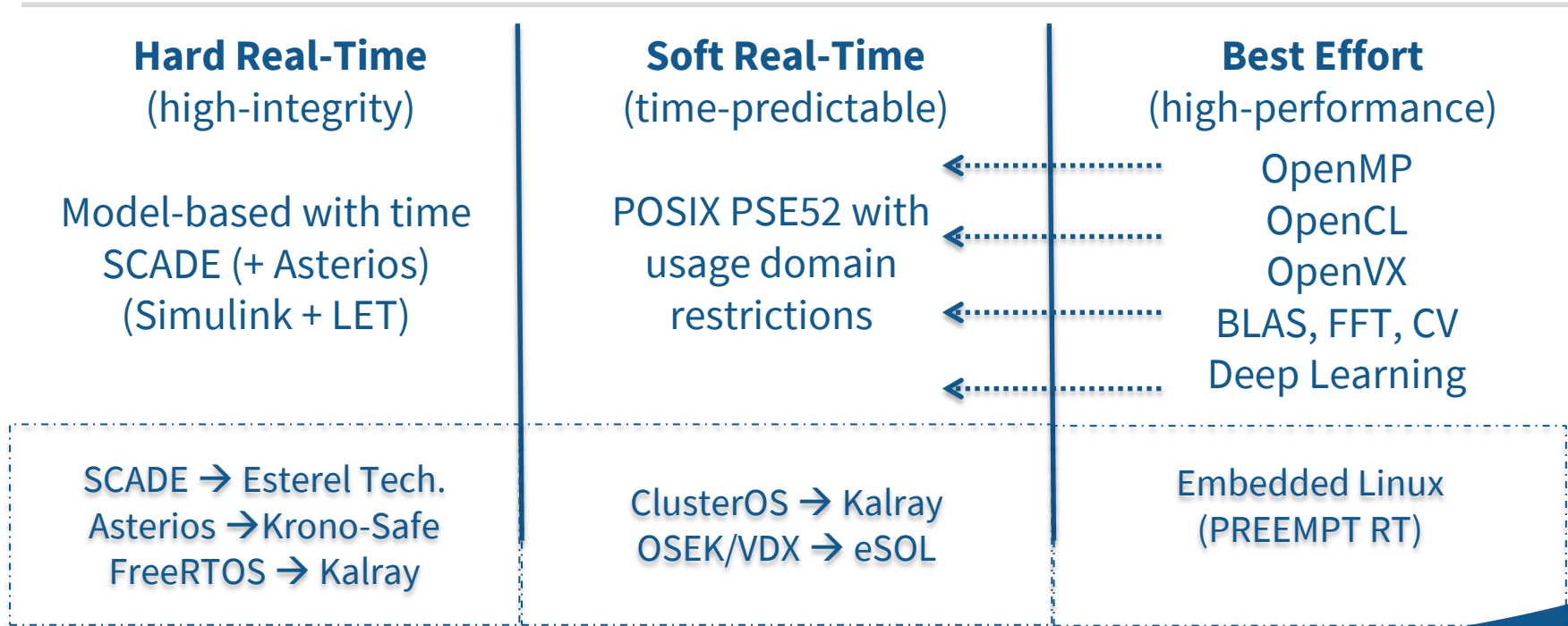
# Outline

Intelligent Systems

Manycore Processors

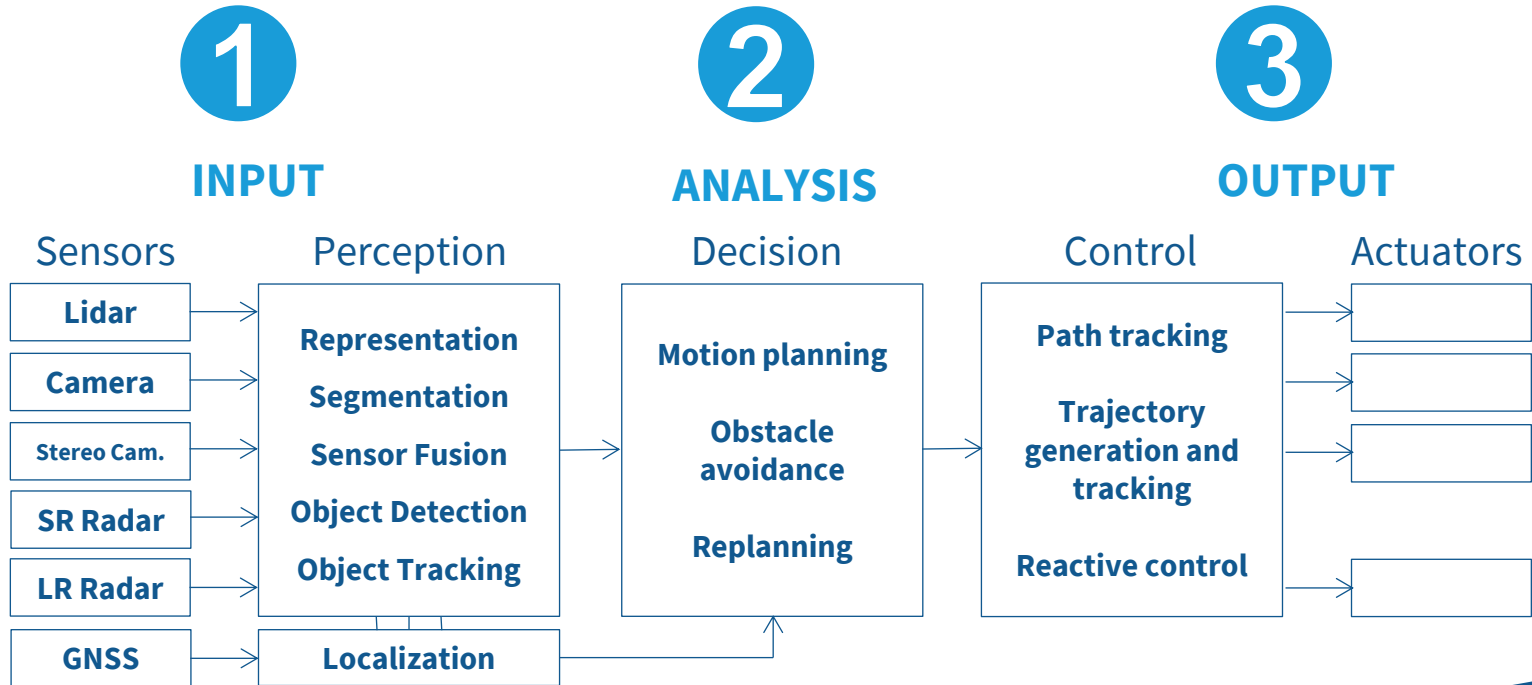Kalray MPPA® Processors

Deep Learning Inference
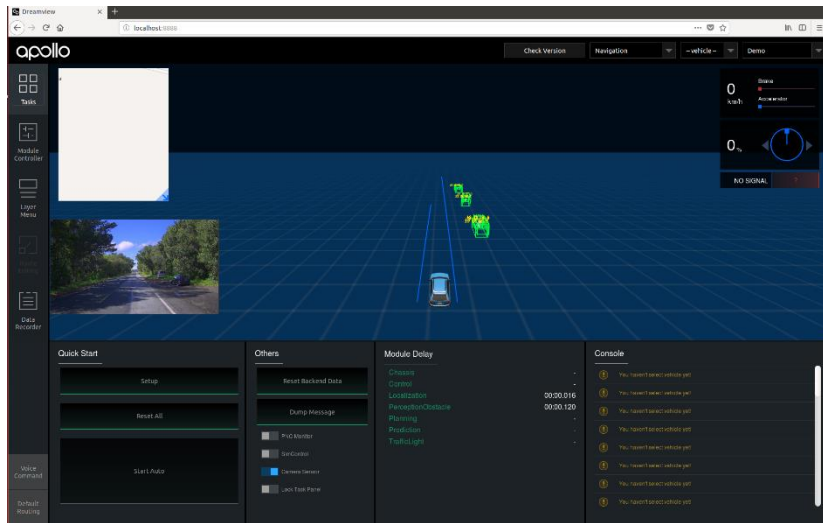
Model-Based Design

**Applications & Outlook**

KALRAY

# MPPA® Embedded Platform

| **Hard Real-Time** (high-integrity) | **Soft Real-Time** (time-predictable) | **Best Effort** (high-performance) |
|---|---|---|
| Model-based with time SCADE (+ Asterios) (Simulink + LET) | POSIX PSE52 with usage domain restrictions | OpenMP OpenCL OpenVX BLAS, FFT, CV Deep Learning |
| SCADE → Esterel Tech. Asterios → Krono-Safe FreeRTOS → Kalray | ClusterOS → Kalray OSEK/VDX → eSOL | Embedded Linux (PREEMPT RT) |

KALRAY

# Autonomous Driving System

| Sensors | Perception | Decision | Control | Actuators |
|---------|------------|----------|---------|-----------|
| Lidar | **Representation** | **Motion planning** | **Path tracking** | |
| Camera | **Segmentation** | **Obstacle avoidance** | **Trajectory generation and tracking** | |
| Stereo Cam. | **Sensor Fusion** | | | |
| SR Radar | **Object Detection** | **Replanning** | **Reactive control** | |
| LR Radar | **Object Tracking** | | | |
| GNSS | **Localization** | | | |

**KALRAY**

# KaNN Integration into
# 3rd Party Autonomous Software Platforms



MPPA2 Processing of
BAIDU Apollo (Perception)

MPPA2 Processing of
Autoware (Perception)

KALRAY

# Kalray News from CES 2019 (EETimes)



The Dutch semiconductor company revealed at the Consumer Electronics Show here that it has chosen a French startup called Kalray to fill in a void created by Qualcomm when it walked away last summer from a $44 billion deal to buy NXP.

**Under their new partnership, Kalray and NXP are developing a central computing platform that combines Kalray's MPPA processors with NXP's S32 processors.**
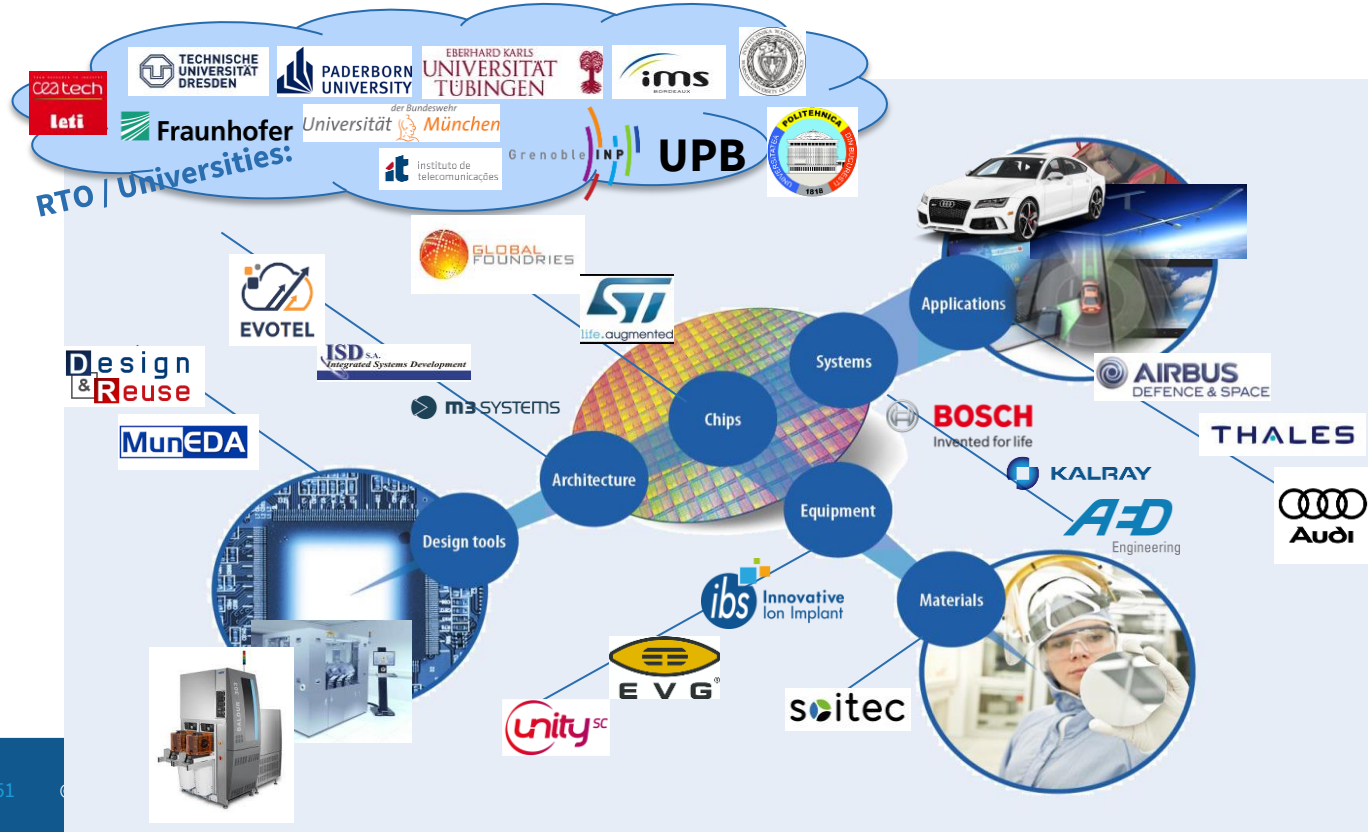
At CES, the companies demonstrate Kalray's MPPA and NXP BlueBox running together on Baidu's Apollo open automotive software.
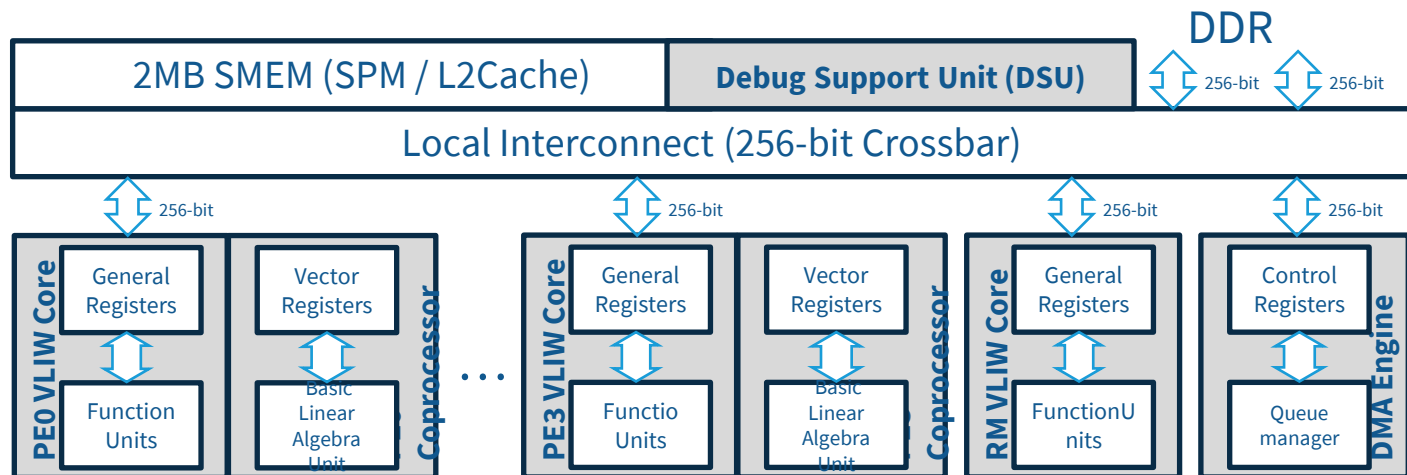
# Mont-Blanc 2020 and EPI Projects

# OCEAN 12 ECSEL Project
# Opportunity to Carry European Autonomous driviNg further with FDSOI technology up to 12nm node

# OCEAN12 Work Package 3 (IP Factory)

## Task 3.1: High-Performance Computing & Vision Signal Processing (Kalray, CEA, ISD, M3S)

- MPPA Cluster tile IP designed and running on FPGA emulation (Altera Stratix-10)
- Running deep learning inference (KaNN) and computer vision (dense Optical Flow)

# Conclusions

## CPU-based manycore accelerators

- C/C++/POSIX/OpenMP/OpenCL programmability
- Energy efficiency & time predictability

## Kalray manycore accelerators

- Learn from GPGPU and computer vision processors
- High compute intensity comes from 2D operations
- Leverage local memories with RDMA engines

## Programming environments

- Optimized application library generators
- Deep learning and graph-based frameworks
- High-performance using OpenCL and OpenMP
- Model-based programming for safety-critical

## European Projects Mont-Blanc 2020, EPI and OCEAN12

**KALRAY**

# Conclusions

## SAFETY

- Hardware partitioning
- Software partitioning
- Hypervisor support
- ISO26262 ASIL B/C

## SECURITY

- Hardware root of trust
- Secure boot
- Authenticated debug
- Trusted execution environment
- Encrypted application code

## DETERMINISM

- Fully timing compositional cores
- Banked on-chip memory
- Interference-free local interconnect
- Network-on-Chip (NoC) service guarantees

## PERFORMANCE

- High-end floating-point and bit-level processing
- DSP-style energy efficiency
- Scalability by replicating clusters

## STANDARDS

- Standard programming environments (C/C++, OpenMP, POSIX, OpenCL, OpenVX)
- Standard development tools (Eclipse, GCC, GDB, LLVM, Linux)

## SCALABLE

- Adaptability to E/E architecture
- Low range to high range car lines
- Allow distribution of functions

**KALRAY**

# THANK YOU

**KALRAY S.A. - GRENOBLE - FRANCE**
180 avenue de l'Europe,
38 330 Montbonnot - France
Tel: +33 (0)4 76 18 09 18
email: info@kalray.eu

**KALRAY INC. - LOS ALTOS - USA**
4962 El Camino Real
Los Altos, CA - USA
Tel: +1 (650) 469 3729
email: info@kalrayinc.com

**KALRAY**