# MANGO

**exploring Manycore Architectures for Next-GeneratiOn HPC systems**

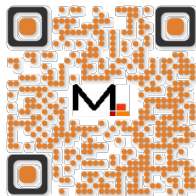# The MANGO Process for Designing and Programming Multi-Accelerator Multi-FPGA Systems

H2RC'18
Sunday, November 11th, 2018
Dallas, Texas, USA

Rafael Tornero, José Flich, José María Martínez, Tomás Picornell, Vincenzo Scotti
Email: ratorga@disca.upv.es

# MANGO Context

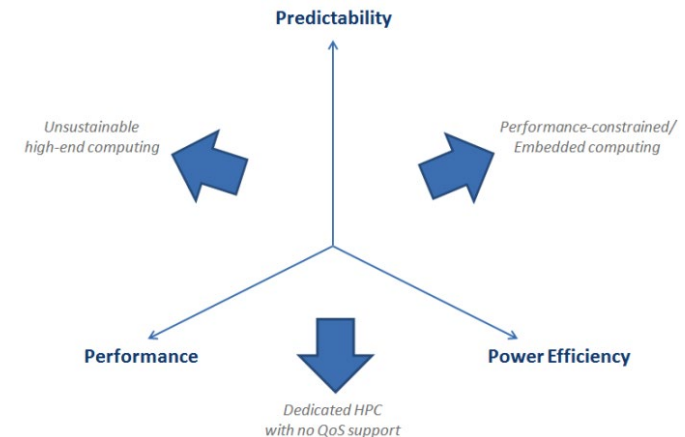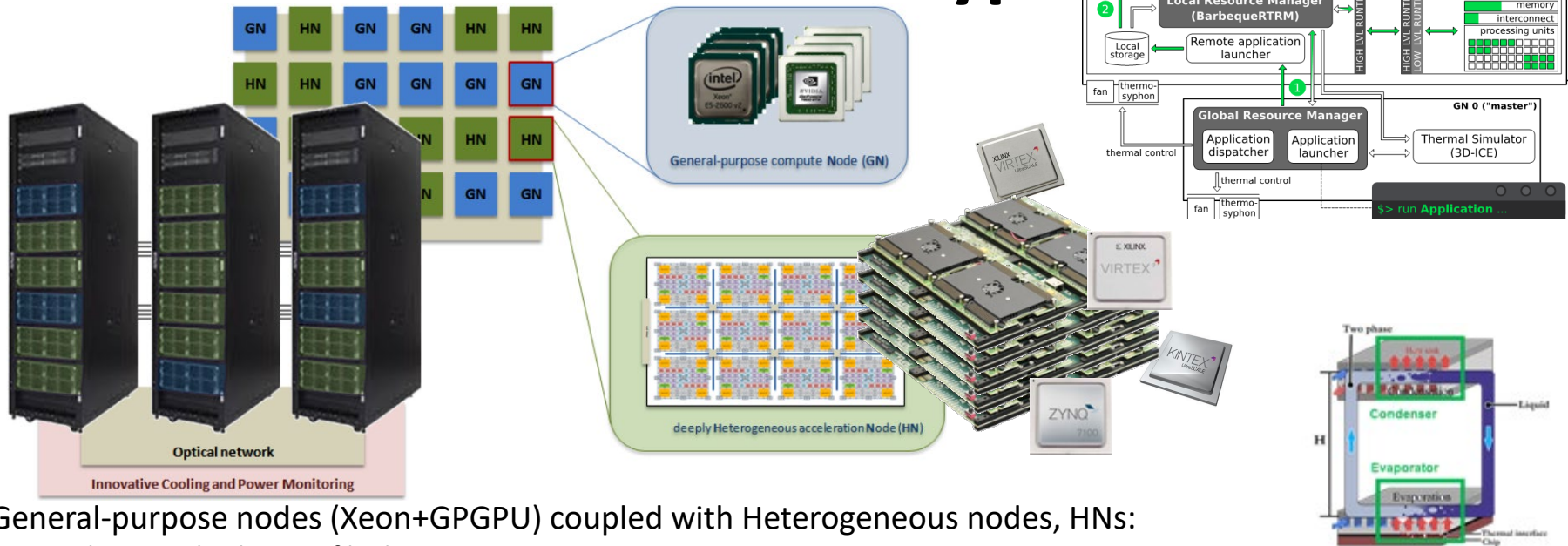o MANGO **FETHPC**-2014 project:
- is about *manycore architecture exploration in HPC*

o HPC quest for performance/power improvement
- Trend in using heterogeneous components
  - GPUs, manycores, and even FPGAs
  - Goal is to get closer to the Intrinsic Computational Efficiency (ICE)
- MANGO focuses on **heterogeneity**
  - How we combine heterogeneous components for the best achievement of computational efficiency
  - How to program/manage them for the best achievement of computational efficiency

o Emerging requirements on HPC systems:

- Predictability (QoS; time sensitivity)
  - Due to the merging of HPC with Big Data

- Capacity computing
  - Run as many application instances as possible

- MANGO addresses **predictability** and **capacity computing**
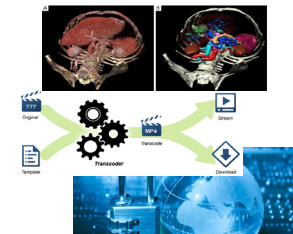  - **3P model** (Performance/Power/Predictability)

Future and Emerging Technologies

# MANGO Context

o MANGO builds a **prototyping system** for 3P space exploration
o Goals:
  - Hardware
    - Develop a flexible prototype for rapid exploration of architectures
    - Explore new deeply heterogeneous manycore architectures
    - Real-time support exploring the PPP design space
    - Provide a unified and simple (homogeneous) access to the system via a smart interconnect
  - Software
    - Adapt programming models and compiler support to the new architectures
    - Develop the right resource manager to deal with the system
  - Infrastructure
    - Provide new monitoring tools to the system
    - Provide new cooling techniques to the system
  - Applications
    - Analyze impact of on a set of real applications
    - Support of video transcoding, medical imaging, security and surveillance applications

# MANGO Prototype



- General-purpose nodes (Xeon+GPGPU) coupled with Heterogeneous nodes, HNs:
  - A large-scale cluster of high-capacity FPGAs
  - A robust, scalable interconnect for a **multi-FPGA manycore** system
  - Will enable FPGA acceleration *at scale*:
    - → a key ingredient for the EsD roadmap
  - A continuum from FPGA emulation to the final physical platform (might be an ASIC manycore, FPGA, mixed…)
    - → *under a <u>stable software environment</u>*
  - Native isolation and partitioning mechanisms for **QoS-aware capacity computing** HPC applications
- Two-phase passive **energy-efficient cooling**
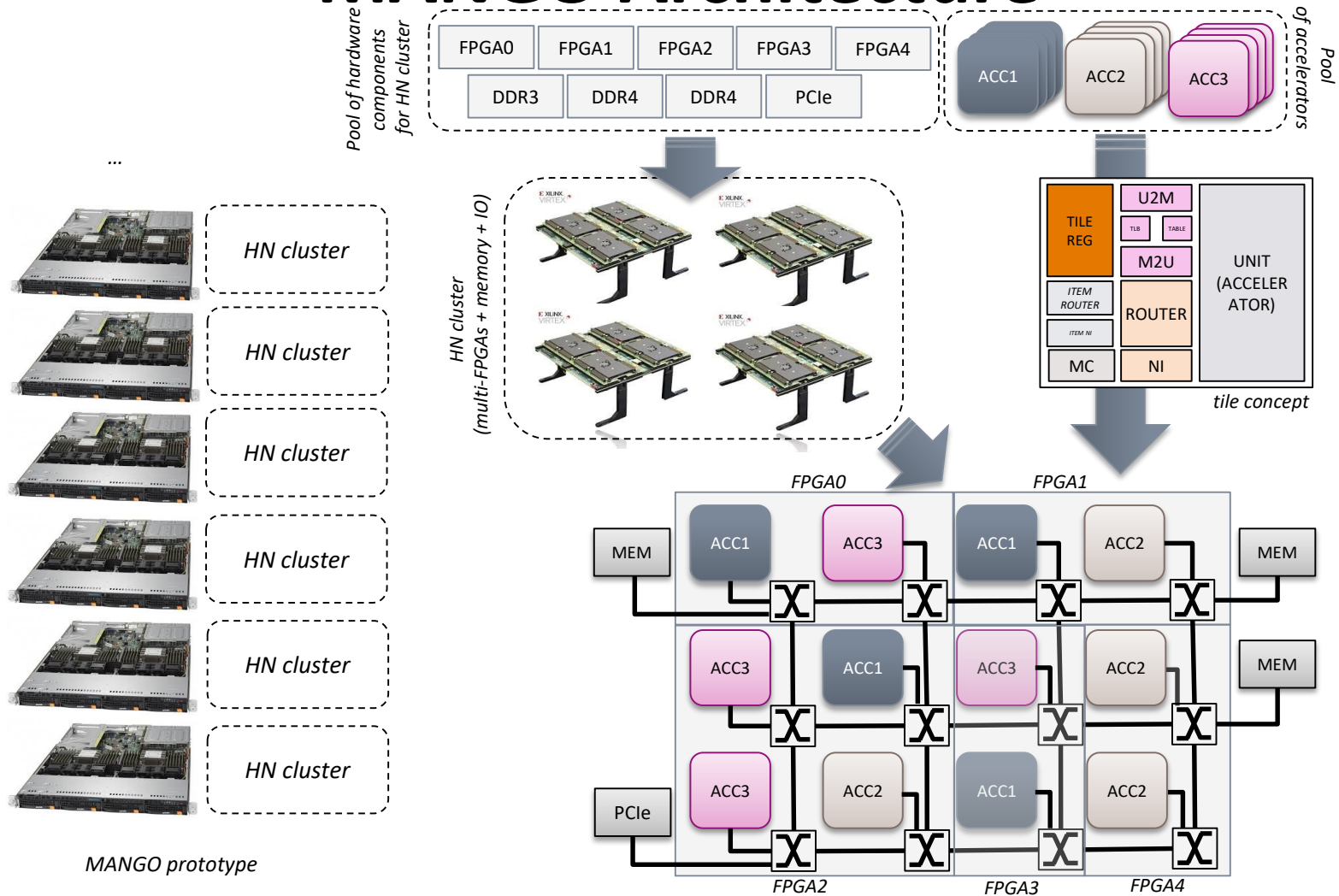- Demonstrated applications with stringent high-performance and QoS requirements

# Consortium



MANGO: exploring Manycore Architectures for Next-GeneratiOn HPC systems

# Agenda

o MANGO Architecture

- HN Hardware and assembly

- Heterogeneity

- Network

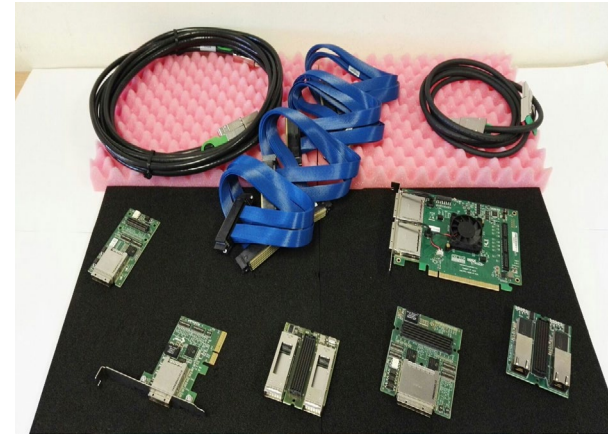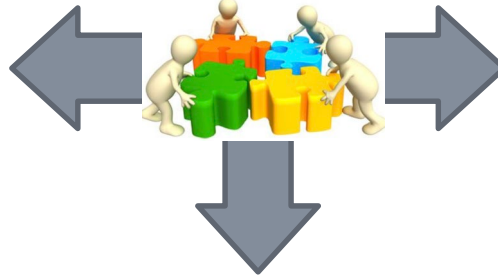- Accelerator Interface

- MANGO Design Flow

o FPGA resource utilization

o Conclusions

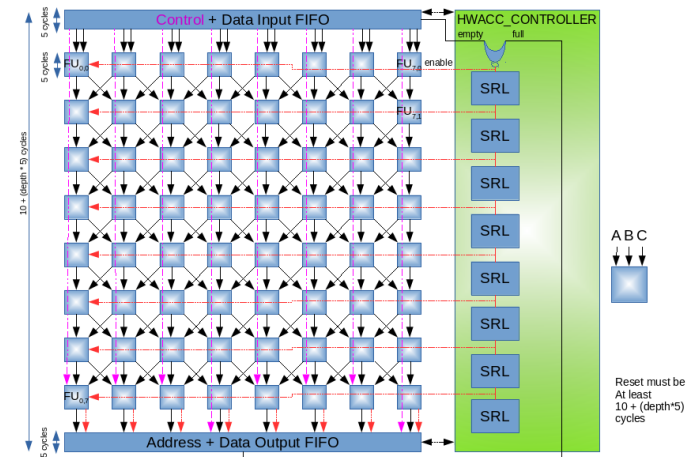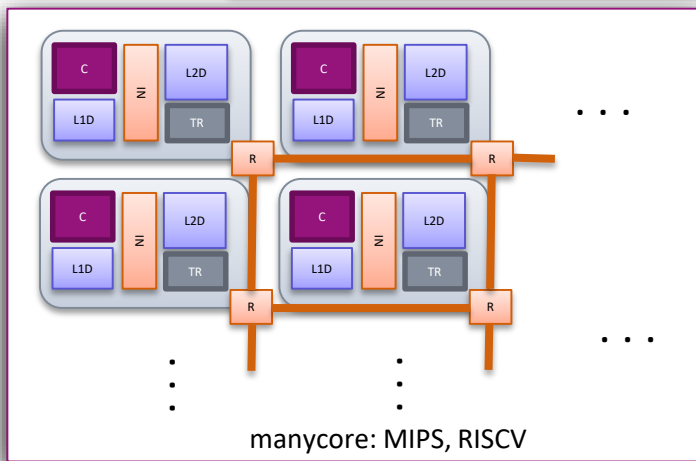# MANGO Architecture

# HN Hardware and Assembly



Lego-like excercise

# Heterogeneity: Pool of Accelerators



NU+ MANYCORE

GPU-like manycore

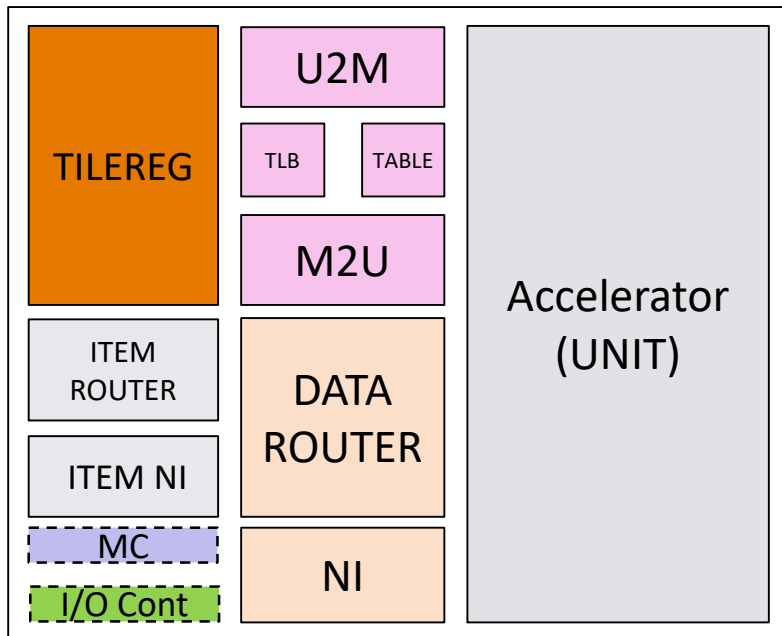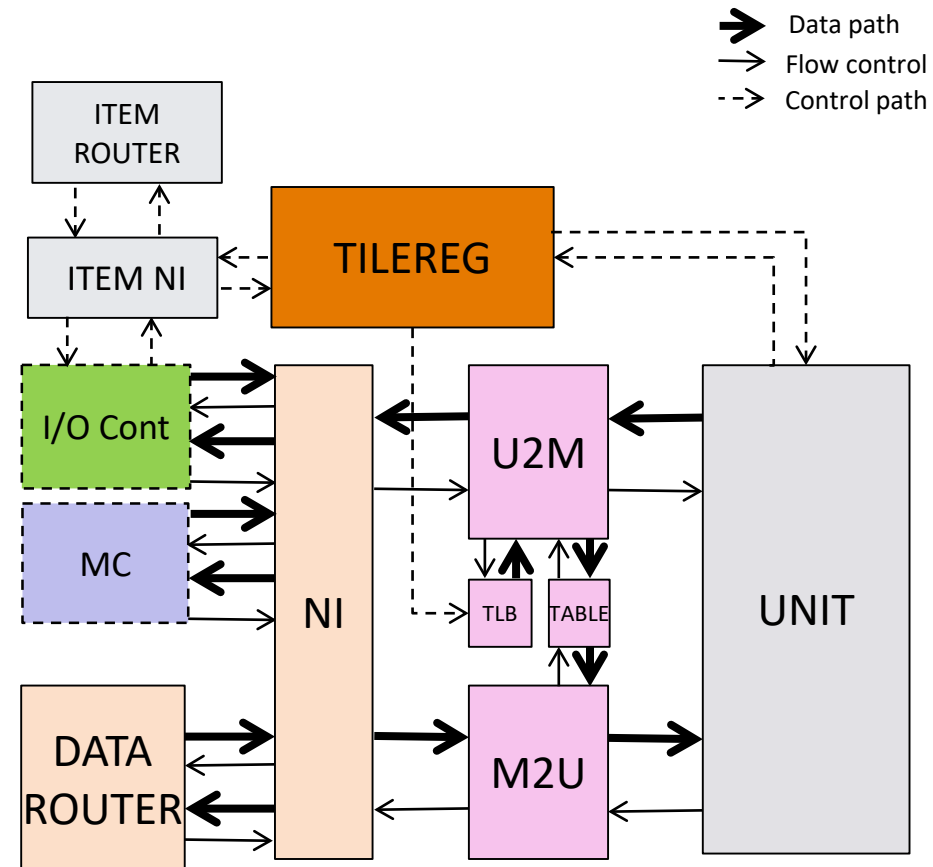manycore: MIPS, RISCV

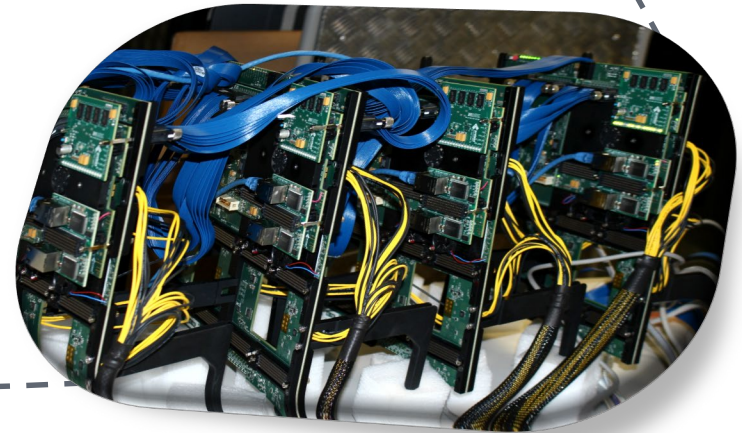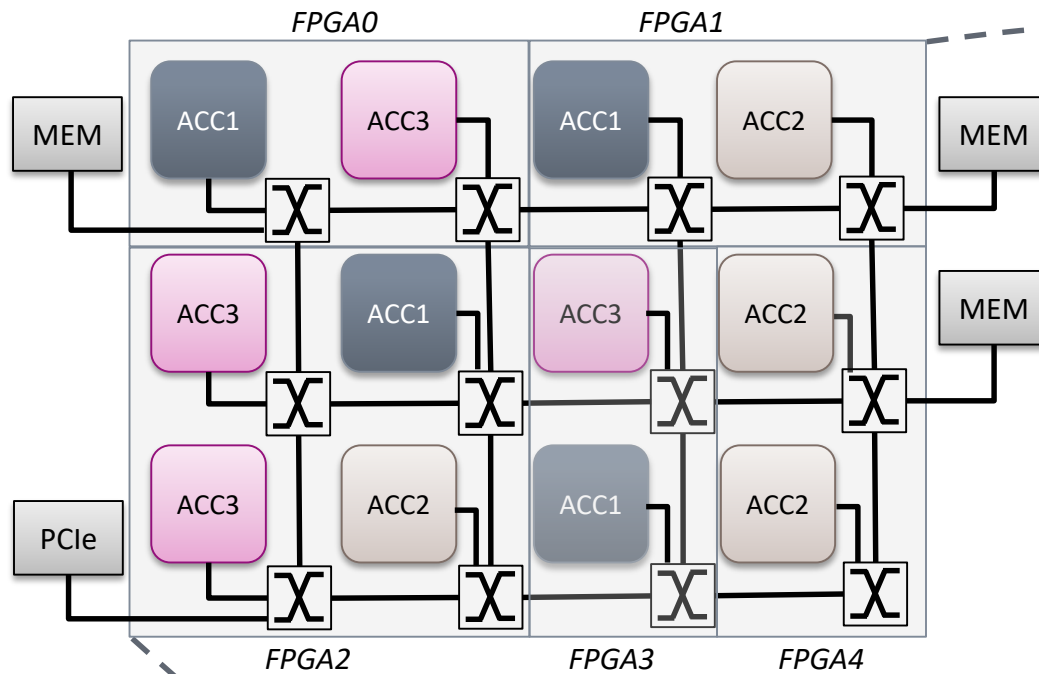General purpose                                    Custom made
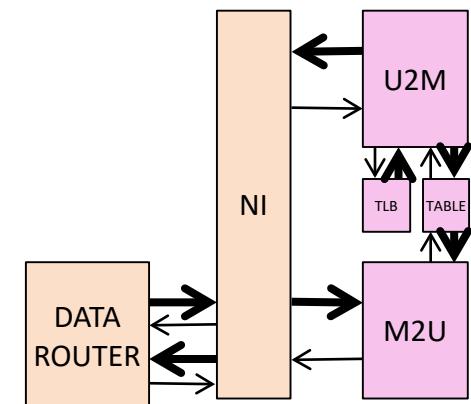
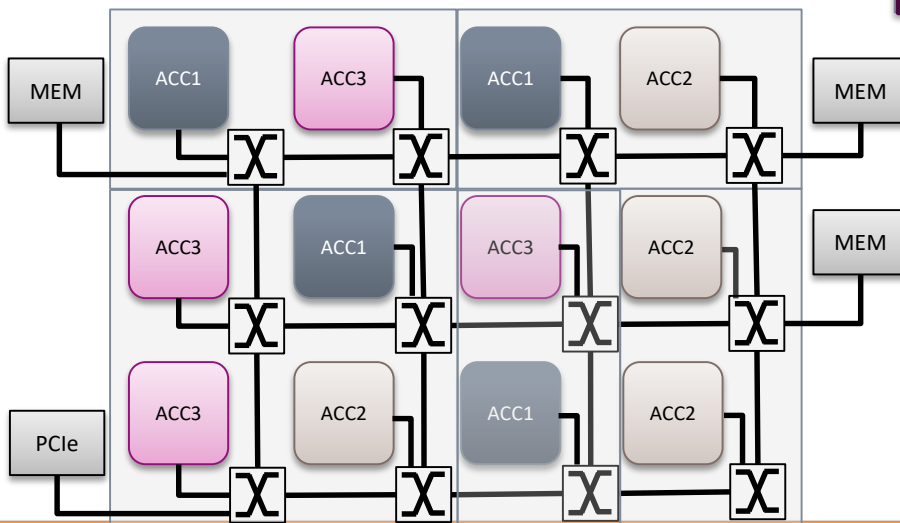# Heterogeneity

## Tile concept
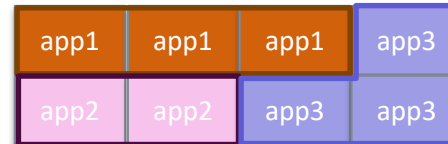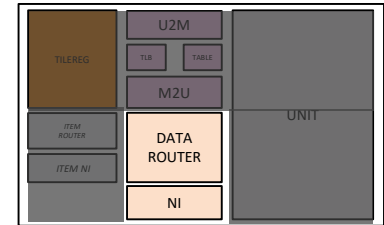


## Element connection

# Heterogeneity, but regular layout

# Data Network

o DATA Routers connected among them in a 2D mesh layout.

o NI decouples network from tile components.

o Support for different Virtual networks (VN) with different number of Virtual channels (VC).

o Support for dynamic assignment of bandwidth per VN.

o Support for capacity computing

# Control Network

- Used for configuration and monitoring
  - MANGO Infrastructure
  - Accelerators or units

- Flexible and generic to let the accelerators be configured based on their complexities
  - Ad-hoc protocols

# Accelerator interface

o Decouple the UNIT from the rest of the MANGO platform

o Allow the implementation of an unique interface for every UNIT

o Unify memory access

    – Byte, Half (16 bits), word (32 bits) & block (512 bits) memory access types

o Allow to map synchronization registers in the virtual memory address space

# MANGO Design Flow



Architecture definition template

Architecture Definition file

Verilog Code

Bitstream file

proFPGA Configuration file

tools

# Architecture Definition File

**ARCHITECTURE ID**
CONCEPT

o Sections
   o General
   o FPGA/multi-FPGA
   o Memory devices
   o I/O devices
   o Network

```
// * MULTI_FPGA__SECTION,
 define N_FPGAS                          3              // How many FPGAs are used for this architecture de.
 define TILE_IDS_PER_FPGA_VECTOR         /*FPGA_2*/10'd9,10'd8,/*FPGA_1*/10'd7,10'd6,10'd3,10'd2,/*FPGA_0*/1b
 define TILE_TYPE_VECTOR                 `TILE_NONE_TYPE,`TILE_NONE_TYPE,`TILE_NONE_TYPE,`TILE_NONE_TYPE,`TILE_N

 define TILE_IDS_PER_FPGA_OFFSET_VECTOR  10'd80,10'd40,10'd0
 define MULTI_FPGA
 define USE_DAISY_CHAIN
 define DAISY_CHAIN_VECTOR               `MB0_A3__ID,`MB0_A1__ID,`MB0_C1__ID
 define DAISY_CHAIN_DIR_FROM_FIRST_VECTOR `NORTH_DIR,`WEST_DIR,`SOUTH_DIR
 define DAISY_CHAIN_DIR_FROM_LAST_VECTOR  `NORTH_DIR,`WEST_DIR,`EAST_DIR
 define MB0_A1__MB0_A3_PINES             10'd124        // number of pines available to route VN wires between two fpgas :
 define MB0_A3__MB0_A1_PINES             10'd124        // number of pines available to route VN wires between two fpgas :
 define MB0_A1__MB0_C1_PINES             10'd90         // number of pines available to route VN wires between two fpgas :
 define MB0_C1__MB0_A1_PINES             10'd90         // number of pines available to route VN wires between two fpgas :
 define CONNECTIVITY_PORTS               10'd4
 define FPGA_CONNECTIVITY_OUT_VECTOR     /*FPGA_2*/1'b1,1'b0,1'b0,1'b0,/*FPGA_1*/1'b0,1'b0,1'b1,1'b0,/*FGPA_0*/1'b0,1'b1,
 define FPGA_CONNECTIVITY_IN_VECTOR      /*FPGA_2*/1'b1,1'b0,1'b0,1'b0,/*FPGA_1*/1'b0,1'b0,1'b1,1'b0,/*FGPA_0*/1'b0,1'b1,
 define FPGA_NEIGHBOURS_VECTOR           /*FPGA_2*/`MB0_A1__ID,10'd0,10'd0,10'd0,/*FPGA_1*/10'd0,10'd0,`MB0_A1__ID,10'd0,
 define FPGA_PINOUT_VECTOR               /*FPGA_2*/`MB0_A3__MB0_A1_PINES,10'd0,10'd0,10'd0,/*FPGA_1*/10'd0,10'd0,`MB0_C1_
 define FPGA_PININ_VECTOR                /*FPGA_2*/`MB0_A1__MB0_A3_PINES,10'd0,10'd0,10'd0,/*FPGA_1*/10'd0,10'd0,`MB0_A1_

// * CLOCK__SECTION
 define MMCM_CLKIN_FREQ                  100.0          // Frequency of the input oscilator for the Primary MMCM
 define ACCELERATOR_FREQ                 5.0            // Accelerator frequency (Mhz)

// * TOPOLOGY__SECTION
 define MB0_A1__TOPOLOGY                 `TOPOLOGY_2x2_1_4   // Mother board 0, FPGA A1 implements a 2x2_1_4 topology co
 define MB0_A3__TOPOLOGY                 `TOPOLOGY_2x1_1_2   // Mother board 0, FPGA A3 implements a 2x1_1_2 topology co
 define MB0_C1__TOPOLOGY                 `TOPOLOGY_2x2_1_4   // Mother board 0, FPGA C1 implements a 3x2_1_6 topology co
 `fine GLOBAL_TOPOLOGY                   `TOPOLOGY_4x3_1_12
   `e TOPOLOGY_PER_FPGA_VECTOR           `GLOBAL_TOPOLOGY,`MB0_A3__TOPOLOGY,`MB0_C1__TOPOLOGY,`MB0_A1__TOPOLOGY
    `OPOLOGY_PER_FPGA_VECTOR_w           (`MESH_2D_w*4)
      `OGY_PER_FPGA_OFFSET_VECTOR        10'd120,10'd80,10'd40,10'd0
       `POLOGY_WIDTH_VECTOR              `MESH_2D_w,`MESH_2D_w,`MESH_2D_w,`MESH_2D_w
```
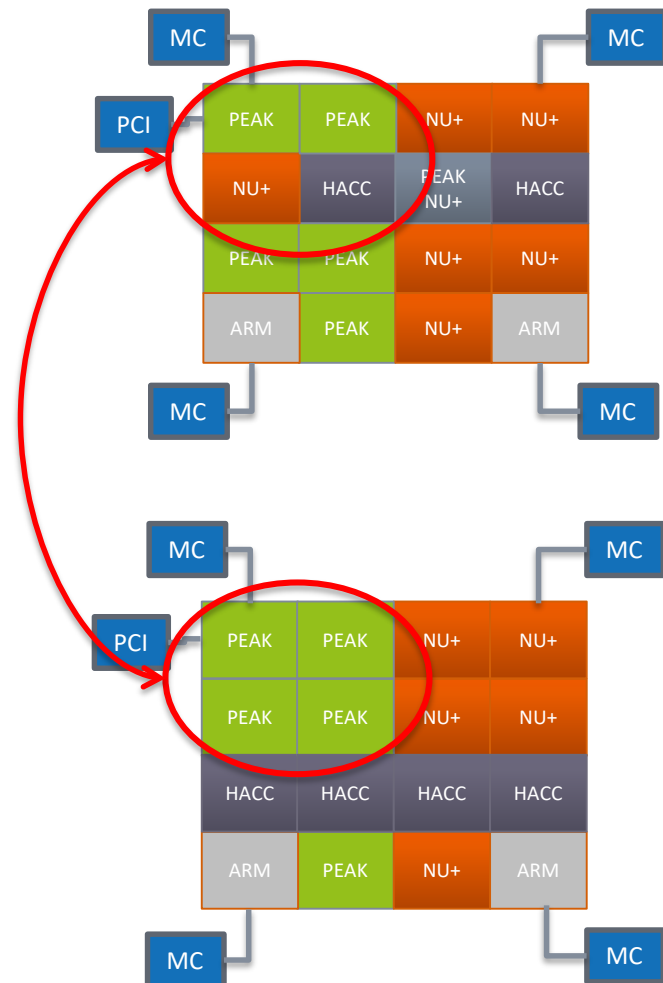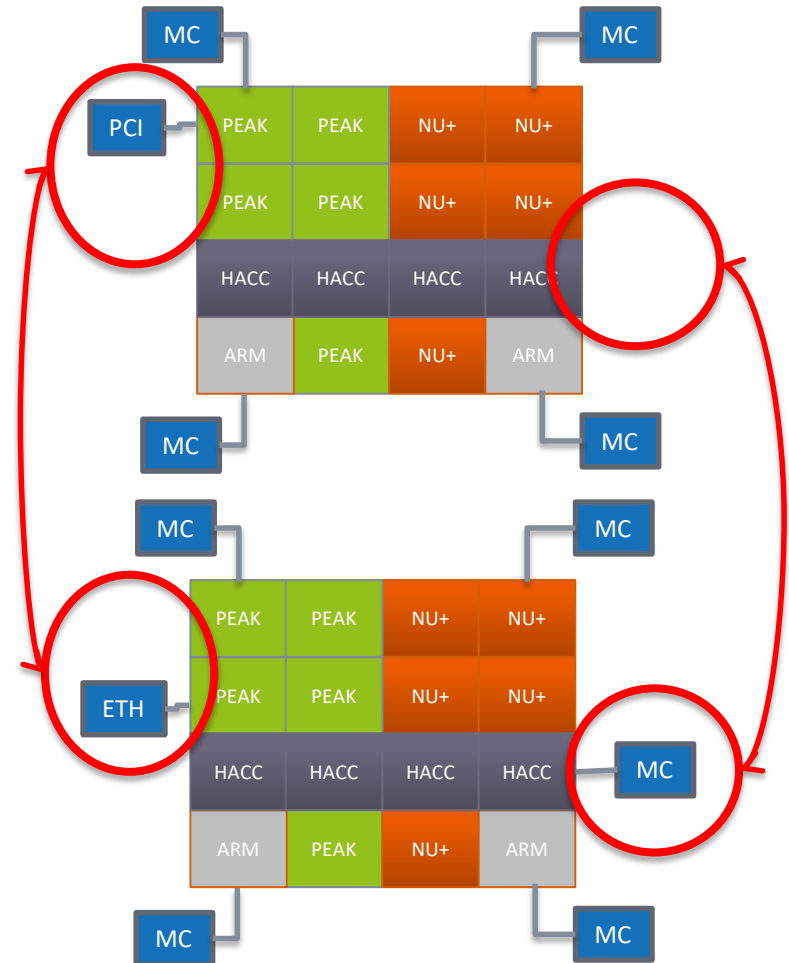
# Architecture Definition Template

## Key aspects

- Ease to configure the Unit type for every tile

# Architecture Definition Template

## Key aspects

- Ease to configure the Unit type for every tile

- Ease to attach Memory Controllers (MC) and I/O devices (PCIe, Ethernet) to any tile
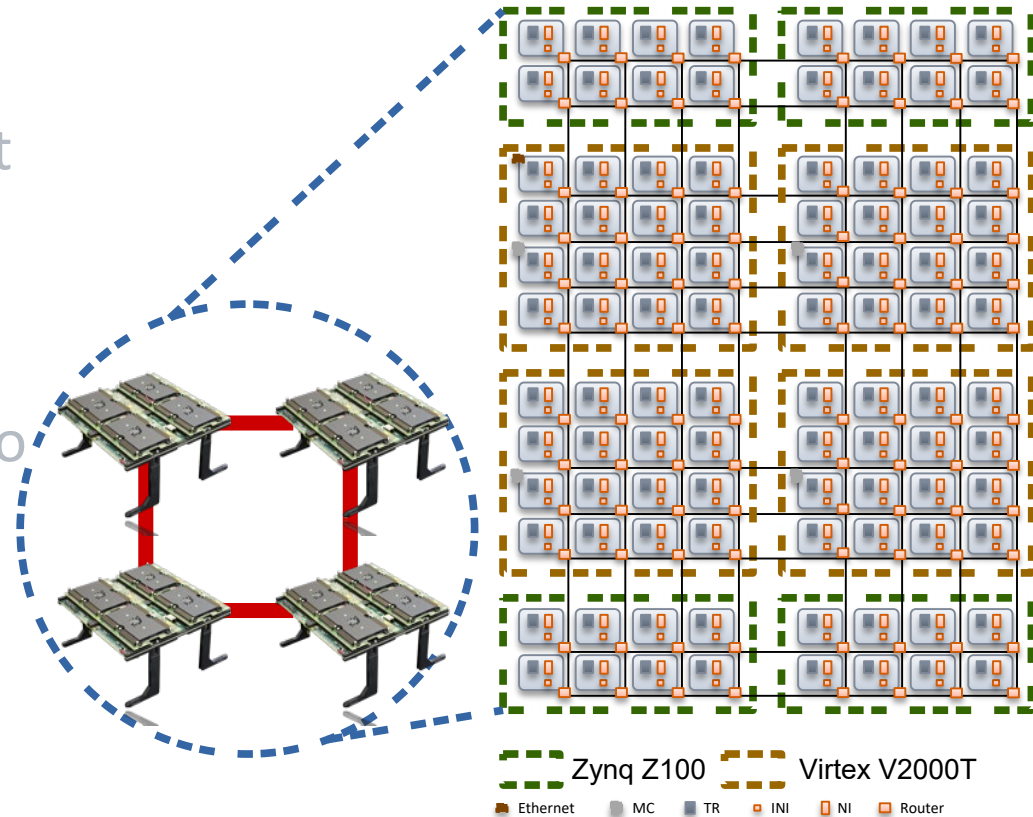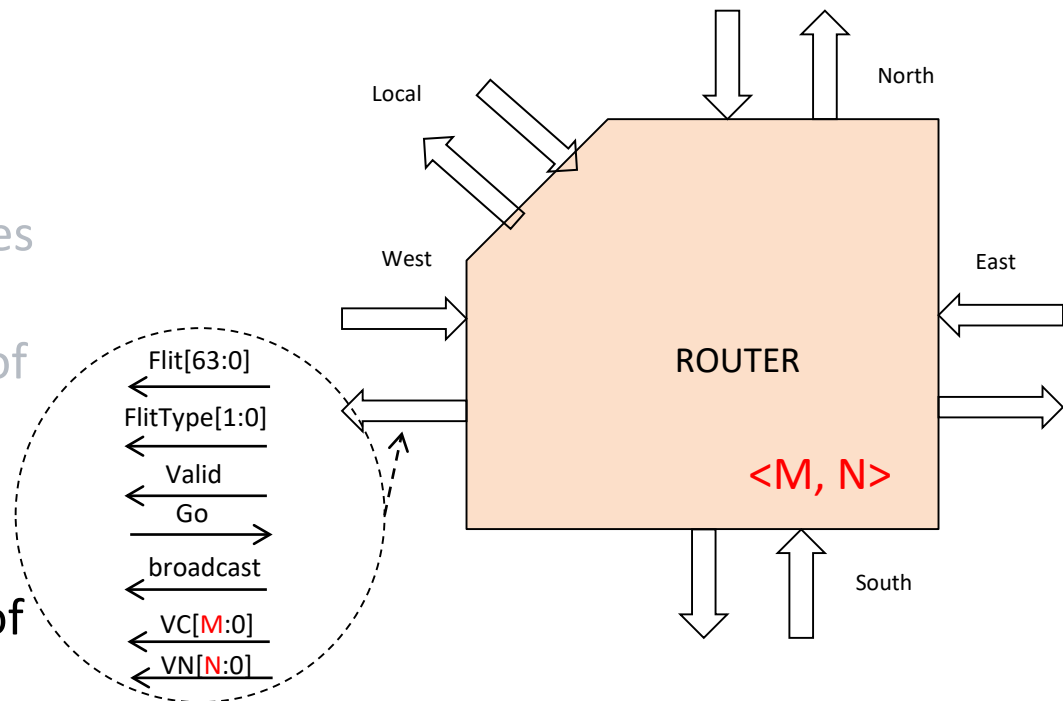
# Architecture Definition Template

**Key aspects**

- Ease to configure the Unit type for every tile

- Ease to attach Memory Controllers (MC) and I/O devices (PCIe, Ethernet) to any tile

- Ease to configure the number of tiles

# Architecture Definition Template

## Key aspects

- Ease to configure the Unit type for every tile

- Ease to attach Memory Controllers (MC) and I/O devices (PCIe, Ethernet) to any tile

- Ease to configure the number of tiles

- Ease to implement multi-FPGA designs



Zynq Z100    Virtex V2000T

Ethernet    MC    TR    INI    NI    Router

# Architecture Definition Template

## Key aspects

- Ease to configure the Unit type for every tile

- Ease to attach Memory Controllers (MC) and I/O devices (PCIe, Ethernet) to any tile

- Ease to configure the number of tiles

- Ease to implement multi-FPGA designs

- Easy to configure the number of virtual networks and virtual channels to achieve QoS guarantee

# Architecture Definition Template

## Key aspects

- Ease to configure the Unit type for every tile

- Ease to attach Memory Controllers (MC) and I/O devices (PCIe, Ethernet) to any tile

- Ease to configure the number of tiles

- Ease to implement multi-FPGA designs

- Easy to configure the number of virtual networks and virtual channels to achieve QoS guarantee

# Architecture Definition File: Example



```
// * MEMORY__SECTION
`define N_MEMORIES                           1
`define FPGA_WITH_MEMORY_VECTOR              `MB0_A1__ID
`define TYPE_MEMORY_VECTOR                   `EB_PDS_DDR3_R2_TYPE
`define TILES_CONNECTED_TO_MC_VECTOR         10'd0
`define NUM_MEMORIES_PER_FPGA_VECTOR         10'd1
```

ARCHITECTURE ID    0

ARCHITECTURE ID    1

# FPGA Resource utilization

- Xilinx Virtex V2000T speedgrade -1

- Tile with Accelerator:
  - MIPS-based cache coherent 2-core accelerator
    - 32K L1I
    - 512K L1D
    - 1MB L2D

## % Utilization

| | Tile (without accelerator) | Tile (with accelerator) |
|---|---|---|
| FF | 1.34 | 5.23 |
| LUTRAM | 2.47 | 4.94 |
| LUT | 3.29 | 10.72 |

LUT ■ LUTRAM ■ FF

# Conclusions

o The MANGO approach for supporting the implementation of multiple accelerators on a multi-FPGA platform

- High customization of the cluster
    - Flexibility for architecture exploration
- Flexibility in configuring an architecture
    - Rapid architecture exploration
- Effectivity in the system communications
    - QoS guarantee
- Effectivity in monitoring the system

o Percentage of resources needed per tile quite slow

# Thank you for your attention

o Contact us at

– [www.mango-project.eu](www.mango-project.eu)

o Other directly related EU project

– [www.recipe-project.eu](www.recipe-project.eu)

# Architecture Definition File



oSections
- o**General**
- oFPGA/multi-FPGA
- oMemory devices
- oI/O devices
- oNetwork

# Architecture Definition

oSections
  oGeneral
  o**FPGA/multi-FPGA**
  oMemory devices
  oI/O devices
  oNetwork



```
MULTI_FPGA__SECTION,
  – N_FPGAS
  – TILE_IDS_PER_FPGA_VECTOR
  – TILE_TYPE_VECTOR

  – TILE_IDS_PER_FPGA_OFFSET_VECTOR
  – MULTI_FPGA

  – USE_DAISY_CHAIN
  – DAISY_CHAIN_DIR_FROM_FIRST_VECTOR

  – DAISY_CHAIN_DIR_FROM_LAST_VECTOR
  – DAISY_CHAIN_VECTOR

  – CONNECTIVITY_PORTS
  – FPGA_CONNECTIVITY_OUT_VECTOR

  – FPGA_CONNECTIVITY_IN_VECTOR
  – FPGA_NEIGHBOURS_VECTOR

  – FPGA_PINOUT_VECTOR

  – FPGA_PININ_VECTOR

....

CLOCK__SECTION,

  – MMCM_CLKIN_FREQ
  – ACCELERATOR_FREQ
  ...

TOPOLOGY__SECTION,
  – MB__FPGA_SLOT__TOPOLOGY
  – GLOBAL_TOPOLOGY
  – TOPOLOGY_PER_FPGA_VECTOR
  – TOPOLOGY_PER_FPGA_VECTOR_w
  – TOPOLOGY_PER_FPGA_OFFSET_VECTOR
  – FPGA_TOPOLOGY_WIDTH_VECTOR
```

# Architecture Definition File

o Sections
- o General
- o FPGA/multi-FPGA
- o **Memory devices**
- o I/O devices
- o Network



```
MEMORY__SECTION
    - N_MEMORIES
    - FPGA_WITH_MEMORY_VECTOR

    - TYPE_MEMORY_VECTOR

    - TILES_CONNECTED_TO_MC_VECTOR
    - NUM_MEMORIES_PER_FPGA_VECTOR
```

# Architecture Definition File

o Sections
   o General
   o FPGA/multi-FPGA
   o Memory devices
   o **I/O devices**
   o Network

```
IO__SECTION
  - N_IO_DEVICES
  - FPGA_WITH_IO_DEVICES_VECTOR

  - TYPE_IO_DEVICES_VECTOR

  - TILES_CONNECTED_TO_IO_DEVICE_VECTOR
  - NUM_IO_DEVICES_PER_FPGA_VECTOR
  - USED_FOR_VECTOR
```

# Architecture Definition File

oSections

    oGeneral

    oFPGA/multi-FPGA

    oMemory devices

    oI/O devices

    o**Network**

```
DATANET__SECTION

   – VN_WITH_PRIORITIES
   – DATA_NET_NUM_VC_PER_VN
   – MANGO_DATA_NET_NUM_VN
   – DATA_NET_FLIT_w
   – DATA_NET_PHIT_w
   – MANGO_DATA_NET_VNID_VECTOR
   – MANGO_DATA_NET_VN_PRIORITY_VECTOR_w
   – MANGO_VN_WEIGHT_PRIORITIES
```

# DMA Transfers

o DMA controller attached to every memory controller.

o Programming parameters: base address, source tile, size of the transfer, destination of the transfer, target address.
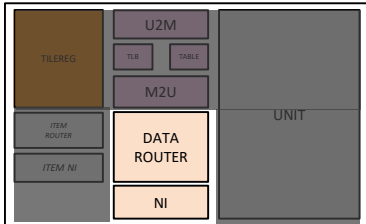
o Possibilities:

# Data Network



- DATA Routers connected among them in a 2D mesh layout.
- NI decouples network from tile components.
- Supported different Virtual networks (VN) with different number of Virtual channels (VC).
- Supported dynamic assignment of bandwidth per VN.

| app1 | app1 | app1 | app3 |
|------|------|------|------|
| app2 | app2 | app3 | app3 |

| T0 | T1 | T2 | T3 |
|----|----|----|----|
| T4 | T5 | T6 | T7 |

VN 0 (10%)

VN 1 (80%)

VN 2 (10%)