## POLITECNICO DI TORINO

Master's Degree in Electronic Engineering

Master's Degree Thesis

# Actuation control for Flight Control Computer application by using parallel computation via AMD Xilinx Versal AI Engine



Supervisor:

Prof. Mario Roberto Casu

Co-supervisors:

Eng. Antonio Bruscino Eng. Riccardo Sticca Candidate
Paolo Molino

Academic Year 2024-2025

#### Abstract

The advent of high-performance computing on aircraft computer platforms introduces several design challenges. In particular, digital fly-by-wire Flight Control Computers rely on sophisticated actuation control loop algorithms, which demand significant hardware resources for efficient execution.

To address these challenges, the AMD Xilinx Versal ACAP (Adaptive Compute Acceleration Platform) provides a cutting-edge and promising architecture. This FPGA is the industry's first compute platform that integrates dedicated acceleration engines (AI Engines) for scalar and vector processing, tightly coupled with programmable logic and configurable on-chip connectivity. Such a heterogeneous architecture enables the design of customized, high-performance hardware solutions.

This thesis, carried out in collaboration with Leonardo Electronics, investigates the development of digital microarchitectures for airborne electronic modules on the Versal platform, and evaluates their hardware implementation both in standard DSP logic blocks and in AI Engines according to the required operation schedules. In particular, using Vitis Model Composer—a high-level development environment—an actuator control loop has been analyzed and mapped across AI Engines and Programmable Logic, interconnected through the AXI4 protocol. C++ code for AI Engines has been developed, while Programmable Logic VHDL code has been generated via graphical block design. Subsequently, the VCK190 Evaluation Board has been programmed accordingly, and performance has been assessed.

The results confirm the feasibility of the proposed implementation and highlight the potential to further generalize the algorithms in order to fully exploit the computational capabilities of the hardware.

# Acknowledgements

I would like to express my gratitude to Antonio and Riccardo for their help and for the opportunity and trust given to me. I would also like to thank my colleagues at Leonardo for making me feel part of the team.

## Contents

Li	st of	Tables	S	6
Li	${ m st}$ of	Figure	es	7
1	Intr	oducti	ion and thesis organization	11
	1.1	Organ	ization	11
<b>2</b>	Pre	limina	ry background	13
	2.1	Versal	ACAP	13
		2.1.1	The AI Engines	14
		2.1.2	The NoC	20
		2.1.3	AIE-ML	22
		2.1.4	Single Event Latch-up (SEL) and Single Event Upset (SEU)	23
	2.2	Fly by	wire control system	25
		2.2.1	The FCC	27
		2.2.2	The actuator control electronic (ACE)	27
		2.2.3	Parallelization for AI Engines	29
3	Des	ign an	d toolchain overview	31
	3.1	The de	esign flow	31
		3.1.1	Application Mapping	31
		3.1.2	Performance Modeling and Simulation	32
		3.1.3	HW only vs embedded systems	32
		3.1.4	Power and Thermal Planning	32

	3.2	Tools:	Vitis Model Composer	33
		3.2.1	HDL - AIE connections	33
		3.2.2	Features of the tool	34
	3.3	The A	XI4-Stream handshake	35
		3.3.1	AXI4 connections in complex DFG	35
		3.3.2	AXI4 connections in very complex DFG	37
	3.4	The b	oard	39
		3.4.1	The Versal Device	39
		3.4.2	The board	40
4	Imp	olemen	tation	43
	4.1	AI En	gines employment for the ACE	43
		4.1.1	AI Engine kernel - IIR order 2 filter, stream	43
		4.1.2	Enlarging the kernels	49
	4.2	The w	chole ACE on AIE + PL	52
		4.2.1	AIE and PL interconnection	52
		4.2.2	AIE subsystem	53
		4.2.3	PL subsystem	56
		4.2.4	Testbench & Simulink simulation	60
	4.3	Tests	on the board - simple example	66
		4.3.1	Procedure	66
		4.3.2	Debug of AXI handshake in PL	68
		4.3.3	Direct measure	72
		4.3.4	Other considerations	77
		4.3.5	Estimation of complete system behavior	80
	4.4	Code	generation and test of the ACE	82
		4.4.1	Results & performances	84
		4.4.2	Other reports	87
5	Con	clusio	ns	93

	5.1	Considerations on technology and tools	93
	5.2	Summary and Conclusions	94
$\mathbf{A}$	Mat	th demonstrations	97
	A.1	Two IIR order 2 filter compressed in a single IIR order $4 \ldots \ldots$	97
В	Mod	del composer generated code	99
	B.1	ACE host.cpp	99
	B.2	ACE AIE_subsystem.h	108
	В.3	ACE AIE_subsystem.cpp	113
	B.4	ACE s2mm.cpp	114
	B.5	ACE mm2s.cpp	116
$\mathbf{C}$	Ker	nel codes	119
	C.1	Single IIR filter order 2 - buffer	119
	C.2	Single IIR filter order 2 - stream (not optimal)	120
	C.3	Single IIR filter order 2 - stream	121
	C.4	Single IIR filter order 4 - stream (improved)	123
	C.5	Controller	124
	C6	Input block	190

# List of Tables

2.1	Summary of main differences between AIE and AIE-ML	23
3.1	Table summarizing main activities from design to bitstream generation. After the first generation, every source can be customized (but output files have to be generated via bash)	34
3.2	Most important features of the VC1902	39
4.1	Average latency of implemented custom functions. NN*: Not Needed; for the normal AIE, only the useful kernels for the Actuator Control Electronic have been tested	51

# List of Figures

2.1	ACAP scheme [1]	13
2.2	AI Engines array scheme [3]	15
2.3	Interface tiles [3]	16
2.4	AI Engines tile scheme [3]	17
2.5	AI Engines tile interconnections [3]	18
2.6	AI Engines scheme [3]	19
2.7	Scheme of the NoC [11]: NPI = Network Peripheral Interconnect, NMU = Network Master Unit, NSU = Network Slave Unit, NPS = Network Packet Switch	21
2.8	AIE-ML array scheme [5]	22
2.9	LVDT sensor: "the primary winding is illustrated in the center of the LVDT. Two secondary coils are wound symmetrically on each side of the primary coil as shown for "short stroke" LVDTs or on top of the primary coil for "long stroke" LVDTs. The two secondary windings are typically connected in "series opposing" (Differential)", [23]	26
2.10	Scheme of FCC's role in airplane flight. Actuator control electronic (red rectangle) is located in the FPGA Section of the interface board	27
2.11	DFG of the algorithm (Simulink model). Each small square (TFx) is an IIR filter of second order. In yellow are represented inputs, in red outputs, while blue labels correspond to tunable parameters	28
3.1	Interconnection between AIE and HDL blocks, [9]	33
3.2	Handshake mechanism example [9]	35
3.3	Simple PL subsystem with two AXI inputs and one output and its waveforms. Each color represents a group of data processed in the same temporal window or computational epoch	36
3.4	Example of complex AXI handshake protocol	37

3.0	a rising edge. The ALL_READY signal resets the flip-flop, allowing it to detect also situations in which the TVALID does not fall (the new data has already arrived)			
3.6	VCK190 Evaluation Board			
4.1	IIR filter order 2 - Simulink model			
4.2	Simulink model for the testbench of the single kernel			
4.3	Simulink testbench results			
4.4	Assembly code for AIE-ML architecture. Left: dense snippet, right: sparse snippet			
4.5	AIE-ML array. The device considered has a 17x2 AIE matrix (the lower one is the interface row)			
4.6	Trace results of the IIR order 2 custom function on an AI-ML tile			
4.7	Kernel computational structure of the "controller block"			
4.8	Kernel computational structure of the "input block"			
4.9	ACE functional map. Legend: purple = input; green = output; blue block = input of the PL; orange block = input of an AIE kernel			
4.10	Simulink block diagram of the ACE, top level. Yellow ovals are the inputs, while red ones are the outputs. All the interconnections between AIE and PL subsystems are managed through "AIE to HDL" and "HDL to AIE" blocks, which represent AXI4-Stream handshakes			
4.11	Simulink block diagram of the AIE subsystem			
4.12	AIE graph allocation on the array. Red rectangles are the used portions of the memory, while green arrows represent the data transfer			
4.13	Trace analyzer results for the AIE array			
4.14	Feeding of array tile 23			
4.15	Simulink block diagram of the PL subsystem. Blue rectangles contain the logic instantiation, yellow rectangles contain the input logic, and red rectangles contain the output logic			
4.16	Block diagram example for TVALID management in PL: in this case, the output TVALID is generated starting from an AND between all the inputs, delayed as the longest path			
4.17	TREADY generation for all the external world input blocks of the PL subsystem			

4.18	slock used to store the output of AIE inside the PL. It introduces I extra cycle of latency. For the 7 external inputs, these registers are not necessary.	60
4.19	Waveform for TVALID-TREADY internal management. It is supposed that external inputs are always valid, while inputs coming from AIE need the TREADY always high, so to keep the validity high, two internal signals are required (as shown in Figure 4.18). ALL_READY signal allows the reset of these registers and a new set of input data to enter the graph	60
4.20	Testbench of the ACE	61
4.21	Results of the testbench where both models receive the same input. x-axis is in $\mu s$ . In the legend of each graph, the first signal is the one coming from the PL subsystem, the second one comes from the reference model.	62
4.22	Debug probes position in the graph. The block CONTROLLER_1 is the one responsible for the integration	63
4.23	Slewer block diagram. Above: original model, below: PL one	64
4.24	Results of the testbench where each model receives its own feedback. x-axis is in $\mu s$ . In the legend of each graph, the first signal is the one coming from the PL subsystem, the second one comes from the reference model.	65
4.25	Block diagram of the simple system.	68
4.26	Initial block diagram of the PL region of the simple system	68
4.27	ABOVE: wrong initial ILA results. SLOTs 0 and 1 are the input, SLOT 2 is the output. BELOW: theoretical right ones. Each color represents a group of data processed in the same temporal window or computational epoch	69
4.28	Timing analysis results of the simple subsystem. Above the requirements are not met $(T_{ck} = 3.2 \text{ ns})$ , bottom the requirements are met $(T_{ck} = 6.4 \text{ ns})$ .	70
4.29	Correct block diagram of the PL region of the simple system	70
4.30	Correct behavior registered with ILA. Above: zoom on a single transition, below: subsequent transmissions	71
4.31	Settings of ILAs in order to measure with the same trigger. On the left, the first one, on the right, the second one, below the block diagram (axis_ila_1 is the first, while axis_ila_0 is the second)	73
4.32	Waveforms of the AIE array's input.	74
4.33	ILA slot positions in the system. SLOT 0-1: PL input; SLOT 2: PL output; SLOT 3-4: AIE input	75
4.34	Expanded waveforms of the system. Trigger condition: TDATA (slot 3) $!=0.\ldots\ldots\ldots\ldots$	75

4.33	dition: TDATA (slot 3) $!=0$
4.36	Block diagram to conFigure 1 ILA with AXI signal from different clock domains (zoom on Vitis region)
4.37	Chip plan for the simple example. The AIE region corresponds to the Y5 row of the chip, but on Vivado, even if the tiles are used, just the interface tiles are highlighted
4.38	NoC utilization for the simple example
4.39	Power consumption for the simple example
4.40	Resources utilization of the final implementation (absolute values on the left and $\%$ on the right)
4.41	Latency measured on the simple example
4.42	Latency estimation for the whole ACE. Blue number identifies the number of clock cycles to reach the destination (considering the input that generates it, as in Figure 4.16), while the green dotted line identifies the longest path
4.43	Simulink system to generate the .bin file
4.44	ACE project's block diagram.
4.45	Measurement of the latency with the ILA. The scale is in terms of clock cycles (1 cycle = 10 ns). The white checkpoint corresponds to the start of the elaboration, while the yellow line corresponds to the end of the elaboration
4.46	Periodicity of the system registered with the ILA
4.47	Synchronicity of the output validity registered with the ILA. Selected slots correspond to 3 of the 7 outputs (just because it is not possible to register more than 16 signals with ILA)
4.48	Chip plan of the final implementation. The AIE region corresponds to the Y5 row of the chip, but on Vivado, even if the tiles are used, just the interface ones are highlighted
4.49	NoC utilization of the final implementation. The 14 middle blue squares correspond to PL input/output of the system, while AIE-PL communication does not employ the NoC, and so is not present in the diagram
4.50	Power consumption estimation of the final implementation
4.51	Resources utilization of the final implementation (absolute values on the left and % on the right).

## Chapter 1

# Introduction and thesis organization

In modern aircraft, traditional direct connections between pilots' input and mechanical actuators are replaced by a *Fly-by-wire* system. This means that, before reaching the actuators in the real world, the input commands have to be processed by a *Digital Flight Control Computer* (DFCC). This process is capable of integrating autopilot and stabilization functions, and is usually performed by a CPU, which has naturally large overheads and is very inefficient for this kind of specialized computations.

An FCC realized with custom hardware is the best solution, especially given the growing complexity of control algorithms, and the performances can increase even by two orders of magnitude [21].

This electronic system is, however, *Safety-critical*: this means that a malfunction may cause important injuries to people, the environment, or important damage to property. For this reason, the custom hardware has to pass very severe certification processes, such as DO-254, which is the assurance guidance specific for airborne electronic hardware.

The safety-critical applications, such as FCCs, apply for the highest assurance level (Level A). For this reason, at the moment, the most popular choice for hardware implementation is the FPGA, which can provide both control and determinism.

This thesis, conducted in collaboration with Leonardo Electronics, focuses on the analysis of the choice of a recent Versal System on Chip (SoC) - Versal ACAP - in order to realize an actuator control loop, which can become a starting point to be taken into consideration for FCCs design in years to come.

## 1.1 Organization

The work begins with a description of the Versal ACAPs and their main characteristics, with a particular focus on the AI Engines. These engines represent one of AMD's most recent architectural innovations, capable of delivering high throughput through algorithm parallelization, especially in the context of DFCCs (Chapter 2).

In Chapter 3, following a brief overview of the design flow and a powerful tool that

significantly accelerates development—AMD Vitis Model Composer, which is based on Simulink and supports both simulation and code generation [9]—an analysis of the AXI4-Stream handshake protocol is provided. This analysis demonstrates how the protocol can facilitate the mapping of complex data flow graphs (DFGs) between various parts of the chip. The chapter concludes with a short introduction to the evaluation board used for testing the designs.

Chapter 4 details the steps taken to implement an example of Actuator Control Electronic (ACE) on the VCK190 evaluation board, using the aforementioned tool. This chapter also includes a simpler exercise implemented on the same board, whose results and performance are compared with those of the complete design.

Finally, Chapter 5 summarizes the results obtained and offers reflections on the potential future applications of the device in the avionics domain.

## Chapter 2

## Preliminary background

### 2.1 Versal ACAP

The computational demand in AI applications, real-time processing, Wireless 5G applications, modern Radars, and many other fields is continuously growing. One possibility to satisfy the need for performance are Versal adaptive System on Chips (SoCs). They are Heterogeneous devices in 7 nm FinFET technology that allow the user to exploit the advantages of both scalar and parallel computation, with a slice of the device reserved for programmable logic (PL) to provide greater flexibility.

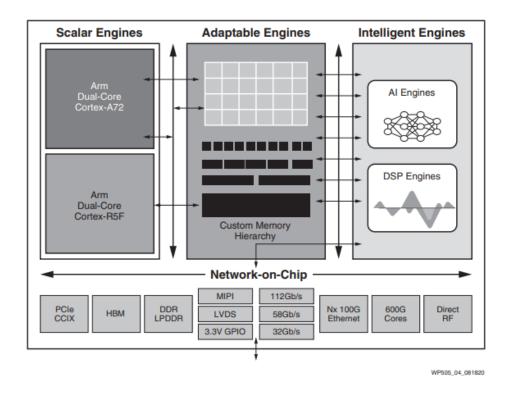


Figure 2.1. ACAP scheme [1].

Versal ACAPs (Adaptive Compute Acceleration Platform) contains (as shown in Figure 2.1) scalar processors (Arm Cortex-A72 and Arm Cortex-R5F), useful for complex algorithms; programmable logic ("Adaptable Engines" in the figure), which can be used also to obtain a custom hierarchical memory; and "Intelligent Engines": a portion of the silicon is occupied by DSP Engines, which usually are DSP58, optimized for  $27 \times 24$  bits accumulation [2]), and AI Engines, where AI stands for "Adaptive Intelligent", and are basically computational tiles with both a scalar and a vector processor units (their details are better explained in the Section 2.1.1). All three parts are interconnected by a high-bandwidth Versal Network on Chip (NoC, Section 2.1.2), which also provides connections with the external world by means of a series of I/O elements and interfaces.

It is clear that this kind of device, properly programmed, is capable of solving basically any kind of computational problem. Furthermore, the user can approach them at any level, from a low-level HDL description to a higher level description by means of C/C++ codes or with the help of tools like Vitis Model Composer, providing graphical solutions (this tool will be better described in Section 3.2).

#### 2.1.1 The AI Engines

The AI Engines are VLIW (Very Long Instruction Word) SIMD (Single Instruction Multiple Data) computational units. All the information related to AI Engines structure and features in this Section refers to [3].

#### The array

They are usually organized in a 2-dimensional array (Figure 2.2), in which the first row is composed of interface tiles, which allow communication from/to the NoC, communication from/to the Programmable Logic (PL), and a single configuration tile. An AI Engine array can contain 30-400 tiles.

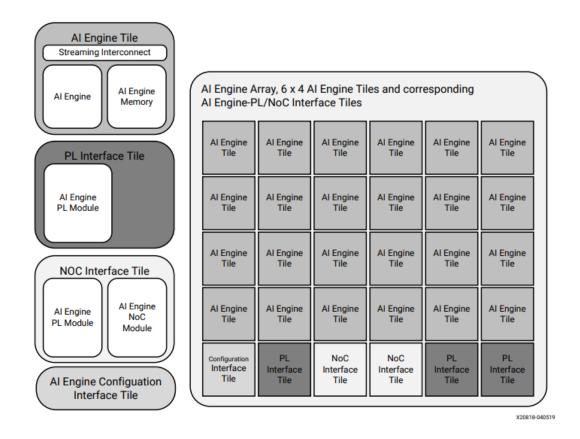


Figure 2.2. AI Engines array scheme [3].

#### The array interface

There are two kinds of array interface tiles:

- The NoC interface tiles are basically AXI4-Stream switches plus level shifters in AIE-NoC interface tiles (the NoC components are in different power domains).
- The PL interface tiles possess 8 AXI Stream ports in upstream and 6 in downstream, which offer a privileged path for PL-AIE communication (see Figure 2.3).

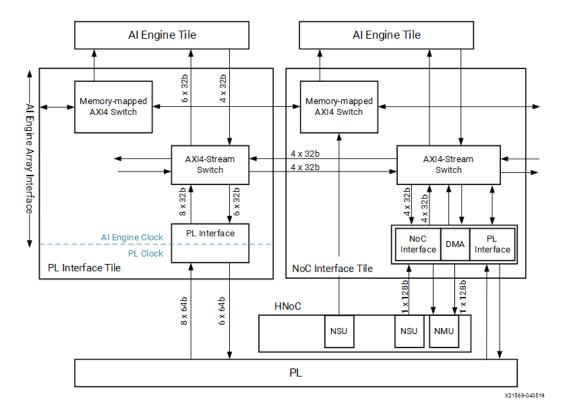


Figure 2.3. Interface tiles [3].

As concerns the configuration tile, it contains the PLL for the AI Engine's clock, the POR (Power On Reset) unit, and registers for global features. It also contains the interrupt generation unit (interrupts can be generated from events, useful for debug or performance simulations).

#### The AI Engine tile

Each tile (Figure 2.4) contains:

- an AXI4 interconnect block
- a memory module
- the AI engine

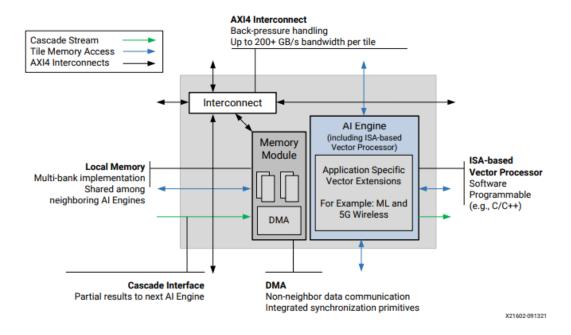


Figure 2.4. AI Engines tile scheme [3].

#### The interconnections

The mapping through AXI4 protocol provides the possibility to read/write the tile registers or their memory from an external master through the NoC. Also, an AXI4-Stream switch, which can be configured both for circuit-switched or packet-switched streams and is basically a programmable FIFO (4-deep FIFO with 2 cycles latency) with 6 programmable arbiters for the packet-switched stream mode.

The packet-switched streaming data transfer's latency is not deterministic, since resources can be contested, while circuit-switched streaming guarantees deterministic latency but not the best overall performance.

A cascade streaming interface allows a tile to write the results directly in the "next" tile's accumulator (a particular type of 384-bit register). Starting from the bottom-left tile, the "next" one is the one that follows a "snake" which moves in the right direction until the end of the row, then moves up of a tile and continues to the left until the end of the row, and so on...

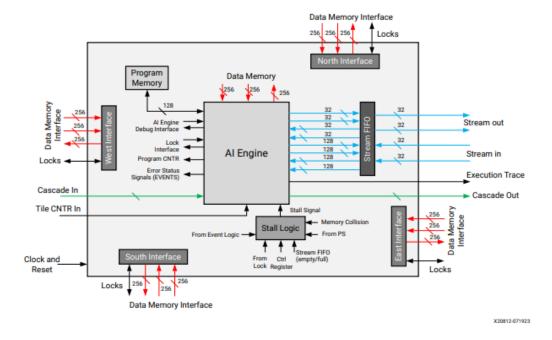


Figure 2.5. AI Engines tile interconnections [3].

#### The memory modules

The memory module can contain 32 kB of data. They are divided into 8 memory banks, in which the first two have 7-bit ECC for each 32-bit word, while the last two have a parity bit for each 32-bit word.

An additional 16 kB of program memory is present, with 7-bit ECC for each 128-bit instruction (it actually contains in fact  $1024 \times 128$ -bit instructions). The program memory can be accessed both from the memory-mapped AXI4 and from the AIE interface.

Despite the other ways, the main method for a couple of neighboring AIE tiles to communicate is the shared memory: each tile can access its 32 kB memory, but also the South tile's one, the North tile's ones and the Right or the Left tiles' one (depending on the position in the array), for a total of 128 kB of accessible memory.

To communicate with non-neighboring tiles, 2 x STMM (Stream to Memory Module) DMAs and 2 x MMTS (Memory Module to Streaming) DMAs are present. 16 buffer descriptors (buffers that contain the information regarding a DMA transfer) can be accessed, and a hardware synchronization (Locks) module behaves as an arbiter between simultaneous requests.

#### The engine

The AI Engine architecture and features are what make Versal ACAPs unique. The scheme of the computational unit is reported in Figure 2.6.

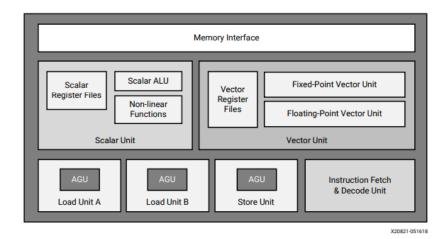


Figure 2.6. AI Engines scheme [3].

Up to 7 instructions can be issued simultaneously: 2 loads, 1 store, 1 vector operation, 1 scalar operation, 2 move operations. Who writes the kernel, which has to be in c/c++, must take this into account, and can take advantage of the AI Engine API [7], which is a set of functions and pragmas to program the engines.

In order to use these functions in an efficient way to reach the best performances possible, a good understanding of the architecture is required.

The AI Engine contains both a scalar and a vector unit.

#### Scalar unit

It can perform standard scalar operations such as sum, subtraction, multiplication, comparisons, bit-wise operations, but also more complex non-linear functions such as sin/cos, absolute value, count leading zeros, square root, inverse and inverse square root, or conversion float2fix and fix2float.

All these operations are performed with a throughput of 1 cycle, but the latency depends on the function; for example, multiplication is performed in 3 cycles, while non-linear functions require up to 4 cycles.

The scalar unit DOES NOT contain a floating point (FP) unit: FP operations are performed through emulation. Furthermore, the floating point arithmetic used in the AIEs is not fully compliant with the IEEE one, for 2 reasons:

- in case of saturation, 0 is returned instead of max/min values;
- if float2fix returns  $-2^{-31}$ , which is still a value within the valid range, the OVFL exception is raised.

#### Vector unit

The vector unit contains three different data paths:

- Multiply Accumulator (MAC) Path
- Upshift Path
- Shift-round Saturate (SRS) Path

The last one is needed since the accumulator registers are 48-bit or 80-bit, so this operation is needed in order to store the results in normal registers, which contain locations with power of two number of bits.

The engine contains scalar registers, special registers, and vector registers. These last in particular are  $16 \times 128$ -bit wide registers, which can be combined to create 256-bit, 512-bit, or 1024-bit vector registers.

The number of MACs depends on the data type and is strictly linked to the vector dimension: for what concerns fixed-point data, 128 8-bit multipliers can be post-added and combined into 16 or 8 accumulator lanes of 48 bits: the 384-bit accumulator register can be seen as 8 lanes of 48 bits. So it can perform up to 128 x 8-bit MAC, but also, for example, 8 x 32-real MACs or 2 x 32-complex MACs with a single instruction.

As concerns FP arithmetic, 8 lanes of SPFP MACs are provided, with a 3-cycle latency and 1-cycle Initialization Interval.

The vector operations allowed include comparisons, min/max (element-wise), and fix2float conversions (for float2fix, the scalar unit is used instead).

The AI Engines do not support Half Precision Floating Point (HPFP) numbers or other custom formats of floating point numbers.

Furthermore, 8 exception bits can be converted into events and raise interrupts.

#### Boot sequence

Here are reported the main steps of the boot sequence, which involves the PMC (Platform Management Controller), which is independent from the PS (Processing System), and is responsible for the boot and configuration of the device.

- 1. Power-on and power-on-reset
- 2. The PMC provides the AIE array configuration using the NPI interface (data comes from a flash device)
- 3. PLL is enabled
- 4. Reset is released
- 5. Programming of the AIE array over the memory-mapped AXI4 (from the NoC)

#### 2.1.2 The NoC

Even if the distances across various parts of the chip are reducing, the delays are increasing due to scaling. This leads to the need to design a proper way to manage data transfer, capable of sustaining large bandwidths and programmable. The information in this Section can be found in [22].

Versal Network on Chip (NoC) globally interconnects all the blocks through a single memory mapping: the master connects to an entry point, called Network Master Unit (NMU), and the slave connects to an exit point called Network Slave Unit (NSU). Routes are programmable and re-programmable, and the data width is (32-512 bits).

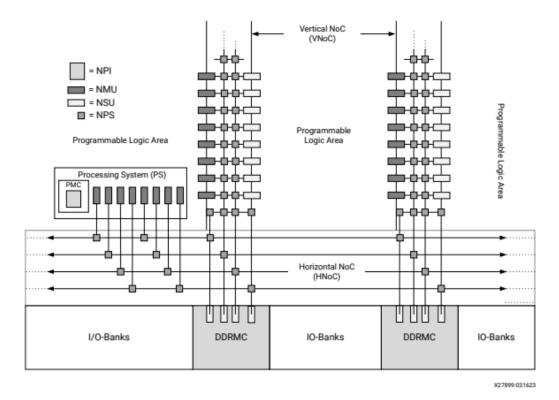


Figure 2.7. Scheme of the NoC [11]: NPI = Network Peripheral Interconnect, NMU = Network Master Unit, NSU = Network Slave Unit, NPS = Network Packet Switch.

The basic elements (NoC Switch) have 4 inputs and 4 outputs, with a minimum latency of 2 cycles, and are arranged into horizontal (HNoC) and vertical (VNoC) sub-blocks.

The network employs *Wormhole Routing* [19]: a packed switched technique in which the packet is divided into *flits*. The first one is the header, and the last one is the tail. The header contains information about the path that is to be followed, all the next ones follow it like a "worm", and the channel is set free when the tail is passed.

There is no need to wait for the entire packet to arrive in a node to start the next movement; on the contrary, this form of routing is like a simple FIFO.

Thanks to the routing scheme, the maximum number of input ports is 4096 and the latency can be deterministic.

In order to satisfy the bandwidth request and the determinism, no link has to be oversubscribed, and the routing graph has to be acyclic, such that deadlocks never occur. In Versal NoC, this is obtained by means of virtual channels and constrained routing (for example, allocating the horizontal portion of the flow always before the vertical).

Furthermore, the traffic from the starting node to the destination can easily be divided into different flows to satisfy the request more easily.

A compiler receives all the user information and constraints about data movements across the chip, and provides a report with the performance of the solution found, making this as user-friendly as efficient.

#### 2.1.3 AIE-ML

The AIE-ML is a different and optimized version of the initial AIE. It almost doubles the computational capability, improves connectivity between different tiles, and in the array organization itself. Furthermore, the memory contained in each tile is doubled, some types of data are not supported anymore, while some others are introduced, and enhanced DMAs are introduced. For any details about them, [5] can be referred.

Since the AIE features have already been described, the remaining part of this Section will be focused on the main differences between the two versions of the engines.

#### New features with respect to AIE

IN AIE-ML, there is no native INT32 support: INT32 multiplications are obtained through decomposition into 16 x 16-bit simpler multiplications. Also, FP32 is supported through emulation of bfloat16 operations.

Bfloat16 and INT4 types of data are supported by this version of the engine.

The spatial disposition of the matrix is different, as shown in Figure 2.8.

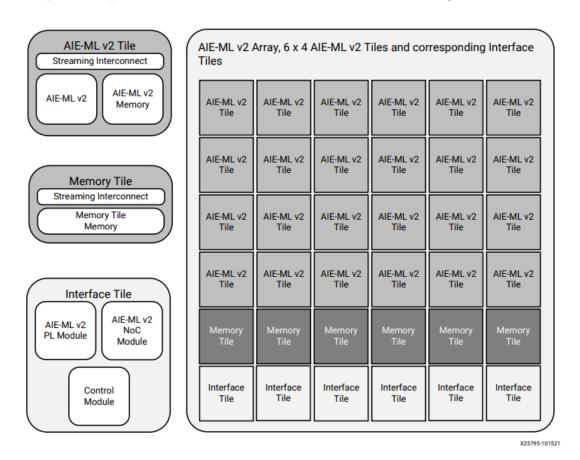


Figure 2.8. AIE-ML array scheme [5].

Each column is equipped with an interface tile, which communicates with both PL and Noc. Furthermore, additional memory tiles are added, which are capable of storing 512 kB of data with ECC and contain 12 DMAs, and support up to 30 GB/s read and write

in parallel.

Also, the data memory inside a normal tile is doubled, passing from 32 kB to 64 kB. This new array structure supports an improved cascade mechanism: 512-bit cascade from left to right and 384-bit cascade from top to bottom (there is no more the "snake" configuration).

The operations' parallelism is doubled: it passes from 128 8b x 8b operations to 256, and introducing the INT4 data type, 512 4b x 8b MAC are supported in parallel. For floating point, the MAC parallelism is 128 bfloat16 x bfloat16 with FP32 accumulation. FP32 ops/tile grows from 16 to 42  $^{1}$ , but are obtained through emulation of FP16 operations: this means that the user can choose whether to have FP32 standard accuracy or to be content with less precision in favor of faster calculation.

Sparsity support is added, while some other features are removed, such as scalar nonlinear functions and not aligned memory accesses.

Table 2.1 summarizes	the main	differences	between	AIE a	and AIE-ML.

	AIE	AIE-ML
Array structure	Checkerboard	All lines identical
Cascade interface	384-bits wide Horizontal direction	512-bits wide Horizontal and vertical directions
Tile stream interface	$2 \times 32$ -bit in and $2 \times 32$ -bit out	$1 \times 32$ -bit in and $1 \times 32$ -bit out
Memory load/store per cycle	512/256 bits	512/256 bits
Advanced DSP functionality	Yes	No
INT4 operations/tile	256	1024
INT8 operations/tile	256	512
INT16 operations/tile	64	128
INT32 operations/tile	16	32
Bfloat16 operations/tile	-	256
FP32 operations/tile	16	42
Data memory/tile	32 kB	64 kB
Program memory/tile	16 kB	16 kB
Memory tiles	-	512 kB
Programmable logic (PL) to AIE array bandwidth	1X	1X
Tile local memory DMA	-	3D addressing mode, S2MM finish on TLAST, out of order packets, compression/decompression
Local memory locks	Boolean	Semaphore

Table 2.1. Summary of main differences between AIE and AIE-ML.

#### 2.1.4 Single Event Latch-up (SEL) and Single Event Upset (SEU)

In radiation-prone environments, such as aerospace or high-altitude avionics, electronic components are subject to Single Event Effects (SEE), which occur when a high-energy particle strikes a sensitive region of the device. Among these effects, Single Event Upset

 $<sup>^{1}</sup>$ Recall that 1 MAC = 2 operations

(SEU) and Single Event Latch-up (SEL) are two of the most critical phenomena to consider during system design and verification.

An SEU is a non-destructive event in which a charged particle alters the logical state of a memory cell, flip-flop, or register, typically causing a bit-flip (i.e., from '0' to '1' or vice versa). The underlying hardware remains physically intact, but the corrupted data may propagate and affect system behavior or software execution.

An SEL, by contrast, is a potentially destructive condition resulting from the triggering of a parasitic structure within the silicon substrate. This creates a low-impedance path between the power supply and ground, leading to an abnormally high current draw. If not promptly interrupted by power cycling, the overcurrent condition can cause permanent damage to the device.

Trying to test SEL and SEU countermeasure proposed by Versal, a study [17] reports: "No SEL events were observed in Xilinx 7 nm XCVC1902 at VCCmax and Tj up to  $120^{\circ}C$  in for both neutron & 64 MeV proton testing for an equivalent of > 10 million years of neutron irradiation at NYC sea level for a total fluence of 2x1012 protons/ $cm^2$  and 1x1012 neutrons/ $cm^2$ ". And "7 nm SEU innovations combined with 7 nm process lead to approximately 10X reduction in CRAM SEU cross Section relative to Xilinx 16nm Kintex UltraScale+TM CRAM SEU cross section. In addition, zero uncorrectable events were observed in 7nm UltraScale CRAM and URAM at sea level with built-in interleaving scheme. And 0.2% of events are uncorrectable in BRAM".

Xilinx also provides a pre-verified Single Event Mitigation firmware (XilSEM). XilSEM on Versal<sup>TM</sup> ACAP performs CRAM error detection and correction (via integrated support for ECC and CRC). Single-bit errors are correctable thanks to ECCs, while multi-bit ones are only detectable. During another experiment, using NYC sea level neutron flux of 12.9  $neutron/cm^2/hr$ , "no uncorrectable CRAM multiple-bit upset (MBU) events (i.e., within the same word) were observed and reported by the XilSEM tool. This result confirms that the built-in interleaving solution is more than adequate for this generation" [14].

## 2.2 Fly by wire control system

In fly-by-wire control systems, the input signal is transmitted to the actuator electronically rather than mechanically. These kinds of systems overcame traditional mechanical ones for several improvements concerning their weight, reliability, and maintenance (since there are fewer moving parts), precision of the response, safety, and, above all, reduction of pilot workload.

As explained in [20], the pilot has basically 3 functions:

- sensing
- regulation
- · decision making.

The first two functions can be defined to be "high bandwidth", because there is little to think about, while decision making is "low bandwidth" since it requires more concentration.

FCC can allow moving the high bandwidth functions to the autopilot, since they are just like reflexes, in such a way that humans' jobs become easier. Of course, computers are not intelligent and do not possess emotions, but, reducing their domain to just some algorithmic functionality, they can achieve operative equivalence.

Basically, in this kind of system, the electronic signal is transmitted to an actuator, for example, a Direct Drive Valve (DDV), controlling the fluid flux and allowing to move the RAM (the aerodynamic surface).

In order to monitor the position of the aerodynamic surfaces across the aircraft, several sensors are placed, for example, LVDT (Linear Variable Differential Transformer) sensors [23], depicted in Figure 2.9.

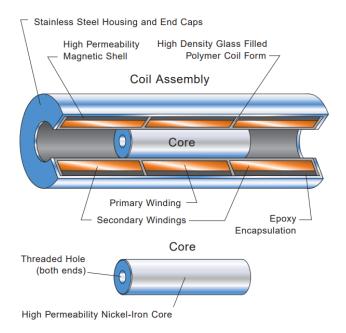


Figure 2.9. LVDT sensor: "the primary winding is illustrated in the center of the LVDT. Two secondary coils are wound symmetrically on each side of the primary coil as shown for "short stroke" LVDTs or on top of the primary coil for "long stroke" LVDTs. The two secondary windings are typically connected in "series opposing" (Differential)", [23]

The primary winding is excited by a sinusoidal signal. The secondary coils provide a differential output, whose value depends on the position of the core (which strongly affect the coupling), allowing a very high precision of position detection.

In particular, for the FCC considered in the following sections, these sensors provides a ratiometric feedback<sup>2</sup>, necessary to stabilize the control loop. For each actuator, LVDT sensors are used to monitor both DDV and RAM position, and the DDV valve is controlled through an H-bridge, that allows the current to flow in both directions through the actuator coil.

The flight control computer is divided in a series of electronic boards with different purposes. In the next paragraphs the various components of the FCC are quickly described, and a schematic drawing is also reported in Figure 2.10 in order to provide a better understanding.

 $<sup>^2</sup>$ Ratiometric feedback means  $V_{fb} = \frac{V_1 - V_2}{V_1 + V_2} \eqno(2.1)$ 

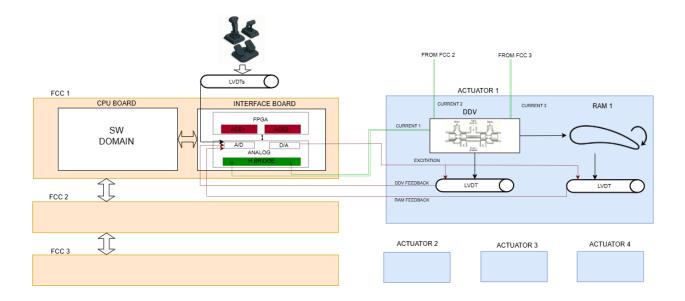


Figure 2.10. Scheme of FCC's role in airplane flight. Actuator control electronic (red rectangle) is located in the FPGA Section of the interface board.

#### 2.2.1 The FCC

The FCC receives an external command every 10 ms, but reads the sensors and updates the PWM value to control an hydraulic valve every 100  $\mu s$ .

This kind of safety-critical apparatus are usually redundant: that means that the same FCC's hardware has 3, 4 or more copies that operates receiving input from the same number of copies of sensors. If some fault is detected from the software, the failed one can be turned off, giving more weight to the other ones. In the actuator, more than one coupled coil is present, one for each copy of the FCC, and their contributes are summed.

### 2.2.2 The actuator control electronic (ACE)

The case study algorithm taken as example is reported in Figure 2.11.

The ACE can be divided into 3 main steps: the first one that receives the SETPOINT command and, given the difference between the desired RAM position and the actual one (feedback), performs a filtering through the first Proportional Integral Derivative (PID) module (that will be called "controller block" later on).

The filtered output is the translation in DDV domain of the needed RAM movement, and so this time is the difference between that signal and the DDV feedback to be filtered and to proceed into the third section of the graph.

The last step is in the current domain, and after another difference with the ACTUATOR CURRENT signal, the execution proceeds through the last PID module, in order to be ready for the PWM generator.

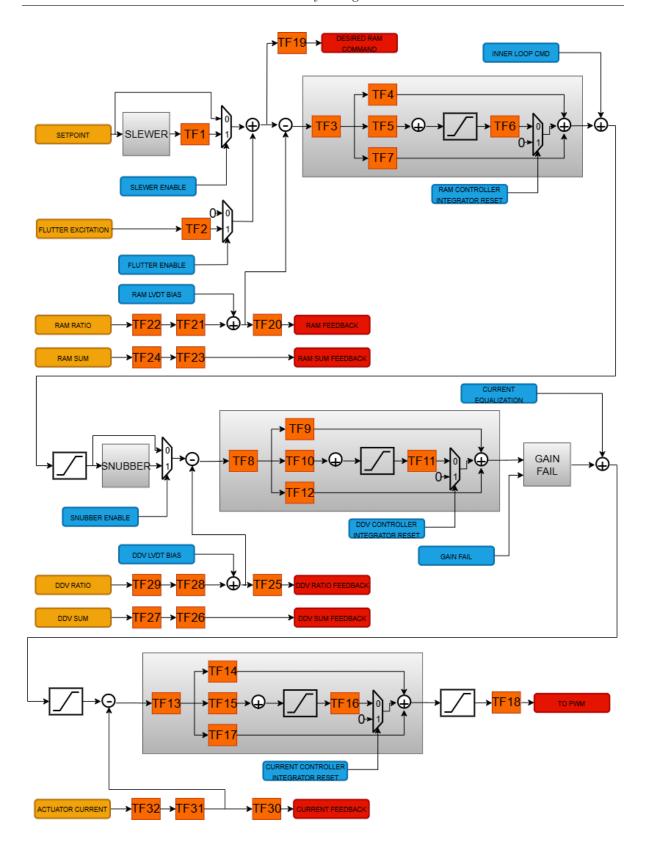


Figure 2.11. DFG of the algorithm (Simulink model). Each small square (TFx) is an IIR filter of second order. In yellow are represented inputs, in red outputs, while blue labels correspond to tunable parameters.

This system has 7 main inputs:

- SETPOINT, which is the position toward which the RAM has to move. This command is generated by the human input and is elaborated by the software domain;
- FLUTTER EXCITATION;
- RAM RATIO, which contains information about the RAM position;
- RAM SUM, that corresponds to E1+E2, and should be constant. If it is not constant, probably something is malfunctioning;
- DDV RATIO;
- DDV SUM;
- ACTUATOR CURRENT, which is the measured current value inside the DDV actuator.

Furthermore, many other control parameters can be used in order to activate/deactivate some functionality or to add some biases. Given this explanation, here is a list of the tunable parameters (they can be tuned by the software domain):

- SLEWER ENABLE: the slewer is a block that operates directly on the input, since its sample frequency is 100 times smaller and its variation may be seen from the system as a step, the slewer can help to mitigate this effect;
- SNUBBER ENABLE;
- FLUTTER ENABLE:
- GAIN FAIL: the sum of the gain fails in the various copies of the FCCs is usually 1, and is equally distributed between the working ones;
- RAM LVDT BIAS, that can be helpful to neutralize bias differences between the copies of the same sensor;
- DDV LVDT BIAS;
- RAM CONTROLLER INTEGRATOR RESET, to reset the RAM integrator loop (see Figure 2.11);
- DDV CONTROLLER INTEGRATOR RESET;
- CURRENT CONTROLLER INTEGRATOR RESET;
- INNER LOOP CMD;
- CURRENT EQUALIZATION.

As concerns the outputs, they are:

- DESIRED RAM COMMAND, which is basically the setpoint passed through the slewer block and filtered:
- RAM POSITION FEEDBACK SUM, which should be 0 and goes to the SW domain;
- RAM POSITION FEEDBACK, which gives the SW info about RAM position;
- DDV POSITION FEEDBACK SUM, which should be 0 and goes to the SW domain;
- DDV POSITION FEEDBACK, which gives the SW info about DDV position;
- CURRENT FEEDBACK AMP, which gives the SW info about H-bridge current in Ampére;
- TO PWM, which is the digital signal that controls the PWM duty cycle.

#### 2.2.3 Parallelization for AI Engines

As previously described, AI Engines are basically vector processors, and the maximum arithmetic intensity (OP/Byte) is obtained through massive parallelization. The most efficient examples suitable for these architectures are vector (or matrix) multiplications, and this is the reason why they are so widely used in machine learning (e.g., CNNs

(Convolutional Neural Networks)).

In the case of an algorithm such as the one described in Section 2.2.2, there is no possibility of parallelization on the input data (for example, providing multiple input samples and elaborating all of them together), since the whole DFG has to be traversed before reading the next input set.

Furthermore, there is a long series of IIR-order 2 filters (13 filters), over which applying some sort of pipeline is not possible since many feedback inputs are present, and it is not possible to apply a pipeline to feedback loops.

Considering this factor, this specific may not be the best kind of algorithm to run on these devices. Anyway, the focus of the thesis is to evaluate its feasibility and, given the astonishing performance of the hardware, suggest some modifications that can lead to an increase in performance or flexibility.

## Chapter 3

## Design and toolchain overview

## 3.1 The design flow

In UG1504 [10], AMD advises following a method in order to help streamline the design on Versal ACAPs.

#### 3.1.1 Application Mapping

The prerequisite is having a good understanding of the requirements of the system, and it can be achieved with a functional behavioral model of the algorithm. The most relevant pieces of information needed about the algorithm are:

- Sampling rates
- Data types
- Storage requirements
- Compute requirements
- Graph form (branches or no, pipeline, feedback,...)

All of them have to be decided on the basis of throughput and latency requirements, maximum bandwidth, power requirements, and precision of the output (this concerns mainly the data type, but more complex data types may lead to more computational time, especially with AI Engines).

After that, the next thing to do is the so-called *Application Mapping*: it is necessary to associate each part of the algorithm to some part of the SoC.

Usually, the more compute-intensive parts are mapped to the accelerators (AI Engines, DSP Engines,...), having in mind what kind of operations they allow, and the data flows and manipulations are left to the programmable logic (PL).

For example, since AI Engines are SIMD vector processors, mapping a kernel on them is efficient only if the algorithm is parallelizable in some way. If it is not the case, the overhead in terms of clock cycles due to the data movements may neatly overcome the resource saving.

#### Data movement

Also, the data movement across the chip has to be properly planned: the AI Engines, for example, have a direct connection with the external DDR memory, but for large bandwidth it is more efficient to pass through the PL.

Furthermore, the PL contains Ultra RAM and Block RAM, which can be employed to build a custom memory hierarchy (that is especially effective in large matrix multiplication, for example, in Convolutional Neural Networks applying tiling techniques[16]).

Data re-usage is the key to bandwidth savings, especially considering the shared memory mechanism of the AI Engines array. Also, the cascade streaming can be exploited to save multiple reading/writing operations, if the data flow consents.

#### 3.1.2 Performance Modeling and Simulation

In order to monitor the traffic flow performances, it is very useful to perform simulations of the NoC dataflow (for example, with AMD generator IP). In this way, the designer can understand if the throughput and latency requirements of the systems are met. Some tutorials are accessible on GitHub from the User Guide.

#### 3.1.3 HW only vs embedded systems

This phase is different depending on the basis of the system type.

For hardware-only systems, embedded or external processors are not used, but only the PMC (*Platform Management Controller*) is used to program and control the design.

As concerns embedded systems, the starting focus should be the type of accelerators that are to be used (PL or AIEs). Generally speaking, both PS, PL, and AIEs can work together, according to the application mapping, which has to be performed manually.

Both AIEs and PL will be controlled by the software stack. The PMC boots the device, and different boot modes are supported. The Arm Cortex-R5F works in a low power domain (LPD), is also ASIL-C safety compliant, and is a good choice if looking for determinism and real-time applications. For more complex applications, the Linux OS can be the best solution, but it requires the full-power domain (FPD).

#### 3.1.4 Power and Thermal Planning

Since restarting from zero is very time and money consuming, inserting power and thermal considerations in the board planning phase is generally a good idea.

AMD provides the Versal adaptive SoC Power Design Manager (PDM) tool, in order to run simulations and obtain power reports, which may help to understand if some constraints are not met.

Thermal simulations allow to get  $\theta_{JA}$ , and in addition to the maximum ambient temperature, that information allows to evaluate the worst-case power consumption.

## 3.2 Tools: Vitis Model Composer

"AMD Vitis™ Model Composer offers a high-level graphical entry environment based on Simulink® from MathWorks for simulation and code generation of designs that include AI Engine, HLS, and RTL components. For more information, see the Vitis Model Composer User Guide (UG1483 [9])" (UG1076 [6]).

This tool takes on another level of abstraction the development of a complex system. It is still needed to write by hand the kernel C++ code for the AI Engine, but it is capable to generate automatically the graph description, mapping the kernels on the array and organizing the data movement between them.

As concerns the Programmable Logic (PL), a whole library of basic blocks is already included in the tool; however, it is possible to insert custom HDL components (even if they should be a single entity) or generate HDL code starting from MATLAB functions.

It also supports HLS, so if the user wants to further raise the abstraction level, it is possible to include custom HLS blocks.

#### 3.2.1 HDL - AIE connections

To exploit both AIEs and PL in the same project, some considerations have to be made:

- all the PL blocks have to be inserted in a subsystem, which has to contain only blocks of the specific HDL library;
- all the AIE kernels have to be grouped in another subsystem, containing only blocks of the AIE library;
- the link between the two subsystems has to be managed passing through "AIE to HDL" or "HDL to AIE" blocks, which are capable of managing AXI4 interconnections by means of a handshake protocol.

These are summarized in Figure 3.1



Figure 3.1. Interconnection between AIE and HDL blocks, [9]

All the simulations are bit-accurate, which means that the results obtained are not approximated but are exactly the ones generated by the hardware. The only problem is that AIE kernel simulations are not cycle accurate: the latency of the kernels and of data transfer between AIE and PL is not considered by the simulator. This does not allow the user to estimate easily latency performances, and some measurements have to be performed on real HW.

#### 3.2.2 Features of the tool

Model Composer does not stop at a simple behavioral simulation, but is capable of also performing a series of steps that normally have to be performed by hand, guiding the user all the way through the flow up to the bitstream generation.

The tool contains a constraints editor, in which it is possible to control which tiles have to be employed to run kernels, FIFO depths in blocks at AIE or PL input, and kernel Runtime Ratios (fraction of available execution cycles that are used by the kernel during operation).

#### **DMAs**

To test the system on hardware with real samples, an input and a golden output files are added to the bitstream, with HLS DMAs responsible for the data movement from DDR memory to input ports and for checking the real output against the Simulink golden reference (it should not happen, but it is possible that on real HW the actual behavior does not correspond to simulations). These two DMAs should be found in the Vivado project's block diagram (as in Figure 4.36).

HLS codes for the DMAs are located in the datamover/src/ directory, and the system behavior is controlled by PS.

By default, in automatically generated PS code (host.cpp), DMAs are programmed and only then the AIE graph is initialized and started (an example is reported in appendix B.1), but to measure latencies (maybe modifying the block diagram on Vivado), the user can customize the C++ code as please. However, in this case, it is necessary to repeat the .elf creation with makefile commands (not relying on automatic code generation offered by the tool, since the host.cpp will be overwritten).

Table 3.1 summarizes how the tool helps.

	User (by hand)	Model Composer
Kernels (.cpp and .h)	X	
Block diagram (PL & AIE)	X	
HDL description & IP creation (PL)		X
AI subsystem graph (.cpp and .h)		X
libdaf.a		X
connectivity file (.cfg)		X
input.txt / reference_output.txt (from Simulink)		X
datamovers (s2mm and mm2s) for in/out		X
platform (.cpp and .h)		X
host.cpp		X
lscript.ld		X
sd_card.img / BOOT.bin		X

Table 3.1. Table summarizing main activities from design to bitstream generation. After the first generation, every source can be customized (but output files have to be generated via bash).

### 3.3 The AXI4-Stream handshake

The communication between the AIE and the PL, in both directions, uses the AXI4-Stream protocol. For each 32-bit data transfer, the TVALID and TREADY signals are used to manage the handshake between the transmitter and the receiver.

UG1483 [9] reports the following:

"The TREADY/TVALID handshake is a fundamental concept in AXI to control how data is exchanged between the master and slave allowing for bidirectional flow control. TDATA, and all the other AXI4-Stream signals (TSTRB, TUSER, TLAST, TID, and TDEST) are all qualified by the TREADY/TVALID handshake. The master indicates a valid beat of data by the assertion of TVALID and must hold the data beat until TREADY is asserted. TVALID once asserted cannot be de-asserted until TREADY is asserted in response (this behavior is referred to as a "sticky" TVALID). AXI also adds the rule that TREADY can depend on TVALID, but the assertion of TVALID cannot depend on TREADY. This rule prevents circular timing loops."

The timing diagram in Figure 3.2 provides an example of the TREADY/TVALID handshake.

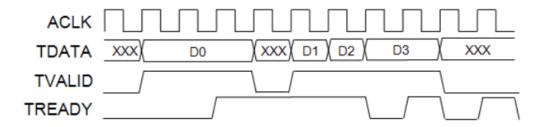


Figure 3.2. Handshake mechanism example [9].

### 3.3.1 AXI4 connections in complex DFG

Usually, data coming from a PL kernel is processed by a single AI Engine (AIE) kernel. In such cases, the simple sequential nature of the dataflow can be implemented using straightforward AXI4-Stream interfaces.

However, when the graph becomes more complex—for example, when the output of multiple AIE kernels must be combined in the PL—it is necessary to ensure that all input data remain valid until all required data are available. To achieve this, the TREADY signal generated in the PL should be the logical AND of all TREADY signals from the PL output interfaces, as well as all TVALID signals from the input interfaces of that Section of the algorithm.

Some hints on how to manage this kind of MISO handshake protocols can be found in [13], where the same concepts have been applied in the context of latency-insensitive protocols (even if in the context of this thesis, the complexity of the suggested adaptive protocols is completely exaggerated given the requirements of the application).

In Figure 3.3 is reported a simple example is reported, in which two input vectors are sent to PL with an AXI handshake protocol just to be summed together.

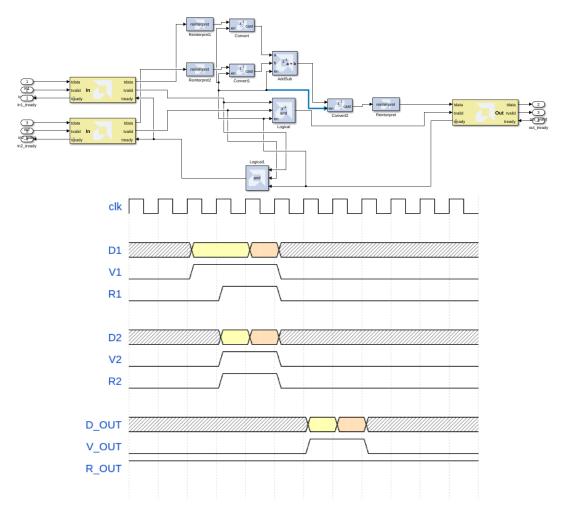


Figure 3.3. Simple PL subsystem with two AXI inputs and one output and its waveforms. Each color represents a group of data processed in the same temporal window or computational epoch.

A first important thing that needs to be underlined concerning complex DFGs is the TVALID pipeline: TVALID latency should last exactly as the main data pipeline (in Figure 3.3 it's 3 cycles).

Then, here is a list of logic functions explained to produce output TVALID, input TREADY, and the ENABLE of the logic blocks:

- Since the pipeline has to be propagated only if all the necessary inputs are valid,  $TVALID_{out} = \prod_{i=1}^n TVALID_i$
- As stated in [9], it is just the TVALID that cannot be generated from the TREADY, but not vice versa, so as TREADY is used

$$TREADY_i = (\prod_{i=1}^n TVALID_i) \cdot TREADY_{out}$$

In this way, data are sampled when the output is ready to receive them.

• The problem arises when the output is no longer ready to receive data: since in the pipeline there are already three valid samples that need to be stopped to not lose

information, all the registers need an ENABLE, which should be  $ENABLE = TREADY_{out}$ 

The waveforms registered on the VCK190 evaluation board with ILA are reported and explained in Section 4.3.2.

Unfortunately, on Model Composer, external IPs inclusion in the design is not so easy, anyway on Vivado there are some Xilinx IPs that may help to fork or combine more AXI4 Streams, for example the AXI4-Stream Broadcaster and the AXI4-Stream Combiner (this one reproduces the same mechanism explained in the previous paragraphs)[4], while the broadcaster is extremely useful when transmitting the same AXI Stream to more destinations.

### 3.3.2 AXI4 connections in very complex DFG

What happens if the mapping of the algorithm leads to an even more complex form of AXI-based PL subsystem?

An example representative of the algorithm analyzed in this thesis work, even if much simpler, is reported in Figure 3.4. Here, data entering the PL are both sent to an AIE kernel and kept inside the PL, and the TVALID should be kept high at least until the output of the AIE is ready, to maintain data coherency.

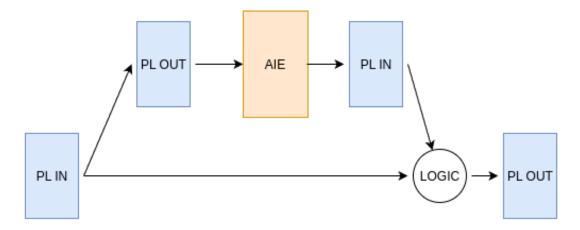


Figure 3.4. Example of complex AXI handshake protocol.

Sampling the input when both TVALID and TREADY are high is no longer sufficient, as the execution dependency on the upper branch must also be considered.

The canonical and most correct approach to this problem involves the use of Vivado IPs: in such a case, a broadcaster distributes the input to both branches, while a combiner merges the output of the upper and lower branches.

However, this solution comes with increased complexity in both the Vivado design and the overall system architecture, and it cannot be simulated within Model Composer.

A simpler and original alternative, designed to stay entirely within Model Composer, can be applied to the specific case in which the complete graph—possibly much more complex—must process one input sample at a time.

Since the case study algorithm is structured to start processing the next sample only after the previous one has fully propagated through the entire graph, the chosen strategy is to generate a single TREADY signal. This signal is asserted only when the full computation has completed.

As a result, an *edge detector* block (Figure 3.5) is required, because the TVALID signal remains high throughout the entire graph execution, while the data must be sampled only once.

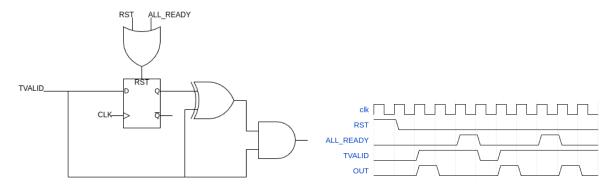


Figure 3.5. Edge detector block: it becomes high only when the input (TVALID) has a rising edge. The ALL\_READY signal resets the flip-flop, allowing it to detect also situations in which the TVALID does not fall (the new data has already arrived).

### 3.4 The board

The chosen board in order to experiment with the adaptive engine's effectiveness on the FCC algorithm is the *VCK190 Evaluation Board* [8] developed by *AMD*, which mounts a VC1902 device.

### 3.4.1 The Versal Device

In particular, the complete part number of the chip is XCVC1902-2MSEVSVA2197, where:

- XC is the commercial Xilinx prefix;
- VC identifies the family (Versal Core);
- 1902 identifies the devices (dimension, resources,...);
- 2 is the speed grade (medium);
- M identifies the voltage, which is medium;
- S identifies the static screen (standard);
- E identifies the temperature grade (0 to 110°);
- VSVA2197 identifies the package: Ball Grid Array (BGA) with 2197 balls;

The device contains 400 x AIE engines. That is precious information to be taken into account in order to properly think about how to efficiently exploit them. They are arranged as a 50x8 tile matrix, with 39 PL interface tiles [3]. With a simple calculation, the total bandwidth for Pl-AIE communication is

$$BW = f_{ck} \cdot 39 \cdot n \cdot 8 \, Bytes/s \tag{3.1}$$

where  $f_{ck}$  is 0.5 GHz (maximum PL frequency), and n is the number of streaming ports (8 for input, 6 for output).

So the actual AIE-PL bandwidth for this specific device is 1.560 TB/s for input streaming and 1.170 TB/s for output streaming.

The device contains 34 Mb of Block RAM, 130 Mb of Ultra RAM, 4xDDR Memory Controller (DDRMC) with a DDR Bus Width of 256 bits. The overall DDR4 bandwidth is 102.4 GB/s, and the overall SRAM Bandwidth is 188 TB/s.

In table 3.2 are summarized other important features of the Versal device.

AI Engines	400
DSP Engines	1968
System Logic Cells (k)	1968
LUTs	899840
Application Processing Unit	Dual-Core Arm® Cortex®-A72
Real-Time Processing Unit	Dual-Core Arm® Cortex®-R5F
Maximum I/O Pins	770
Programmable NoC Ports	28
Integrated Memory Controllers	4

Table 3.2. Most important features of the VC1902.

### 3.4.2 The board

In Figure 3.6, the board is reported.

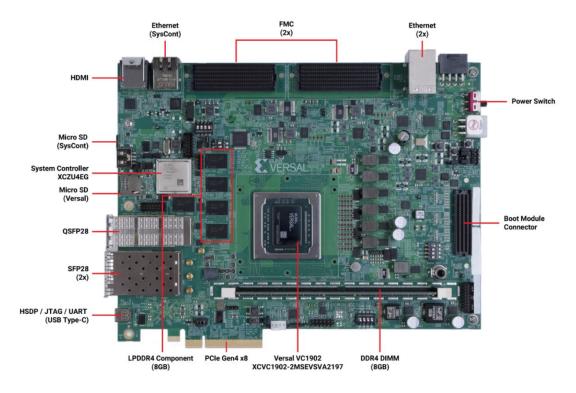


Figure 3.6. VCK190 Evaluation Board.

As concerns the configuration of the Versal device, both JTAG, QSPI32, and SD1\_3.0 modes are supported, and DIP switch SW1 is responsible for the choice of the boot configuration (as shown in [8]). In particular, to boot from an SD card, it is sufficient to:

• copy a valid SD image in the Versal SD (see Figure 3.6). This can be done with the following command (on a terminal inside the folder containing the sd\_card.img:

where "sdc" is the name of the SD. It is strongly suggested to check the right name of the SD card with the command

lsblk

- Set SW1 to SD1 $\_3.0$  configuration (Mode Pins [3:0] = 1110).
- Power the board (SW2) or press the POR pushbutton.

Then, if the serial output has to be displayed on a terminal, it is sufficient to install *Screen* (Linux environment) and type the command

sudo screen /dev/ttyUSB1 115200

In order to connect the board to the PC through JTAG, it can be used the USB-C input port. It is just necessary to tune the SW1 in JTAG configuration (Mode Pins [3:0] = 0000). This is very useful, especially because, with a JTAG connection, it is possible to debug the device with the ILA (Integrated Logic Analyzer).

# Chapter 4

# Implementation

# 4.1 AI Engines employment for the ACE

The first step in implementing the algorithm on this SoC is mapping: it is important to identify which parts of the graph are suitable for execution on the AI Engines.

Based on the considerations discussed in sub Section 2.2.3, and intending to explore the technology as well as evaluate latency for this class of algorithms, the approach is to focus on repetitive patterns within the algorithm. These patterns are typically the computationally intensive parts of the graph, preferably those that can be parallelized.

The first obvious choice, looking at Figure 2.11, is the IIR filter order 2, which appears to be repeated 32 times in the computational graph. An alternative could be a group, for example, 8 IIR filters together in a single kernel: in this way, the parallelism would have been much better exploited. The only problem with this idea is that the longest chain of filters in sequence is 13 filters, so it is not possible to group 8 filters at a time because data dependencies would not be respected.

A first attempt was therefore made with a single second-order IIR filter.

### 4.1.1 AI Engine kernel - IIR order 2 filter, stream

In Figure 4.1, the Simulink model of the filter is represented, which has been described in C++ using the proper APIs.

As shown in the Figure, the already implemented model adopts a signed fixed-point data type format. This format ensures no precision loss during intermediate computations, while the output is finally floored to a 32-bit word with the binary point positioned after 16 bits.

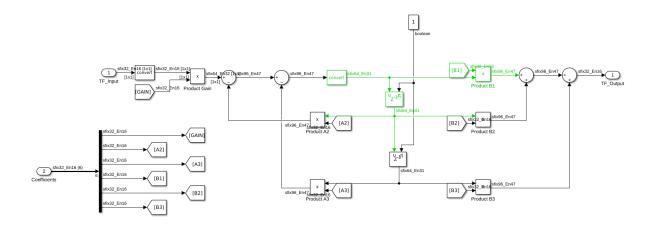


Figure 4.1. IIR filter order 2 - Simulink model.

### Floating point data format for AIE

With AIE, a variable-length fixed-point representation is not possible. For this reason, Single Precision Floating Point (SPFP) or FP32 has been adopted. SPFP, according to IEEE 754, has

- 1 x sign bit
- 8 x exponent bits (biased form)
- 23 x mantissa bits (implicit 1)

Let's evaluate the pros and cons of this choice:

- For what concerns the quantization error, with floating point numbers it is not constant in absolute value, but considering that the first bit of the mantissa is always 1, the relative error is always  $2^{-23}$ , which is worse than the fixed point one (32-bit) only if the whole dynamic is exploited efficiently in the whole algorithm.
- Numbers dynamics increase a lot in terms of orders of magnitude.
- To work in PL, it is more efficient to use fixed-point numerical formats, so a cast operation is required.
- IEEE standards impose an unbiased rounding, so no bias is integrated along the algorithm graph (as it would be, for example, if the rounding technique adopted was biased).

#### The kernel

On the web, several IIR filter parallelizations are available to be used on AIEs [25]. However, none of them can be used in this case, since the examples suppose that parallelism can be exploited by operating on an already given input sequence. As concerns the ACE, only one sample can be elaborated at a time.

The code has been developed from zero, and is reported in the appendix section. Some different versions of it have been studied in order to analyze performances.

Also, the assembly code generated has been compared to understand what could be the best implementation of the single filter.

The experiments on this first simple example enlighten a couple of key points to be taken into account concerning kernel development and feeding of AIE:

- simple scalar mathematical operations on floating point numbers are translated into a call of the FP32 specific function (FPADD or FPMUL), causing large overheads (50-70 cycles). If the same operation is transformed into a vectorial one, even if just the first position of the vector contains meaningful data, the jump to the FP function is avoided, and the performances increase (let's recall that no scalar floating point unit is present in AI Engines).
- if buffer data type is chosen, the AXI4 handshake protocol has to be handled manually. For this reason, it has been chosen to save coefficients as constants inside the internal AIE's memory, and partial products as a static variable (initialized to 0). In this way, just one input and one output are present, and it is possible to adopt streaming data transport (with handshake protocol automatically managed). Both implementations have been analyzed, but only the stream has a future.

To exploit internal parallelism, the IIR filter has been re-conducted as a simple vectorial operation:

$$y_{n} = \begin{pmatrix} x_{n-1} & x_{n-2} & y_{n-1} & y_{n-2} & x_{n} & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} b_{2} \\ b_{3} \\ -a_{2} \\ -a_{3} \\ b1 \\ 0 \\ 0 \\ 0 \end{pmatrix} = x_{n-1} \cdot b_{2} + x_{n-2} \cdot b_{3} - y_{n-1} \cdot a_{2} - y_{n-2} \cdot a_{3} + x_{n} \cdot b_{1}$$

$$(4.1)$$

With subsequent update of the row vector <sup>1</sup>.

### Assembly code generation and simulation

In order to proceed with code generation, the hardware selection inside the Model Composer Hub block is fundamental. If the target is just a simple AIE compilation, it is important to know at least if the chosen hardware contains AIE or AIE-ML, especially if working with floating-point data types.

The reason is that the architectures differ so much that the assembly code can be completely different. For example, AIE-ML does not support FP32 data operations naturally, and those have to be decomposed into FP16 operations "emulating" FP32. There are different kinds of emulation accuracy:

- safe: this is the highest level of accuracy, and the result of the emulated operation perfectly matches the normal one. To emulate an FP32 operations in this way 9 cycles of MACs of bfloat16 are needed;
- fast (default one): slightly worse accuracy, but 6 cycles of MACs required;
- low: worst accuracy (however, it's better than bfloat16), but only 3 cycles of MACs.

<sup>&</sup>lt;sup>1</sup>the value  $x_n$  is the input value multiplied by g. In this way, when the row value is updated,  $x_{n-1}$  is already the one multiplied by g.

The floating-point accuracy mode can be set through the following pre-processor directives:

```
-DAIE2_FP32_EMULATION_ACCURACY_FAST
```

and has to be set on a single kernel basis.

The AIE code verification is composed of three phases:

- 1. Compiling the AI Engine graph design.
- 2. Running simulation using the AI Engine simulator.
- 3. Verifying the simulation results by comparing the output with the golden reference output (generated by Simulink).

If the directive is not properly set, the third phase of the code verification compares AIE results obtained with the default *fast* accuracy, with normal FP32 Simulink's result, and this leads to a result mismatch (even if probably in proximity of mantissa's LSBs) <sup>2</sup>.

The difference between a fast accuracy mode IIR filter order 2 and the same filter generated for a standard AIE tile can be 155 instructions vs only 80. In order to work with SPFP numbers, given the obtained result, it would probably be more efficient to use normal AIE.

### Test and comparison

The code has been simulated with Simulink, and SPFP results compared with the ones obtained by the original model. The testbench prepared is reported in Figure 4.2.

<sup>-</sup>DAIE2\_FP32\_EMULATION\_ACCURACY\_SAFE

<sup>-</sup>DAIE2\_FP32\_EMULATION\_ACCURACY\_LOW

<sup>&</sup>lt;sup>2</sup>This behavior was not well documented. After opening a ticket with Xilinx support, their response was: "We are also discussing with the development team to document this behavior for cases where simulation mismatches occur on AIE-ML devices, so that users can be aware and apply the pre-processor flag themselves."

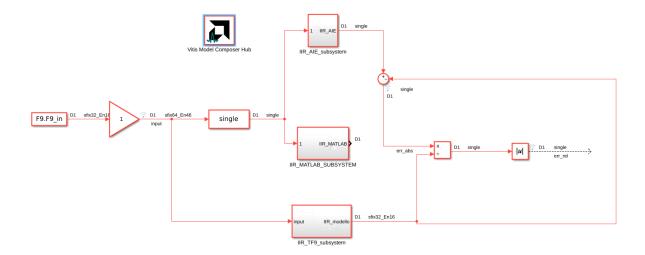


Figure 4.2. Simulink model for the testbench of the single kernel.

As a device, initially the ve2302 was selected, containing a 17x2 matrix of AIE-ML<sup>3</sup>.

The filter TF09 of the Actuator Control Loop has been selected as a subject for the test. The same input has been fed to both the AIE kernel, the MATLAB standard filter block, and the original fixed-point filter. Then, absolute and relative error between AIE and original solutions has been evaluated.

The MATLAB block and the AIE kernel did not present any differences, while, given the dynamics of the input, no noticeable absolute error was detected between SPFP and fixed point models, as shown in Figure 4.3.

<sup>&</sup>lt;sup>3</sup>Initial tests were performed supposing to use the Trenz TE0950 evaluation board [24]; however, due to bureaucratic issues, the hardware that eventually arrived for testing was the VCK190 evaluation board. As a result, the necessary kernels were also generated for the standard AIE architecture (see Table 4.1).

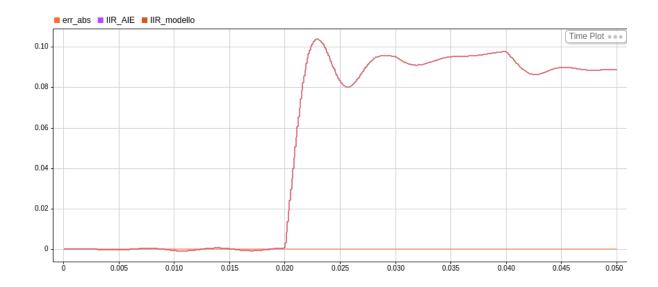


Figure 4.3. Simulink testbench results.

### **Analysis**

Thanks to Model Composer, it is possible to have a look at the generated assembly code.

The obtained latency for the final kernel (code in C.3) is 155 cycles (which means 155 ns). The main problem of this kernel, looking at the assembly code, is that, even if some parts are quite dense (let's recall that the AIE are VLIW processors, able to issue up to 7 instructions per clock cycle), when it comes to sum together the 5 components of the vector data dependencies has to be respected. For this reason, several NOPs are allocated since the processor architecture is pipelined, and at least 4 cycles are needed until the results can be re-utilized for another calculation. The portion of the assembly code dedicated to that operation (Figure 4.4) is extremely inefficient.



Figure 4.4. Assembly code for AIE-ML architecture. Left: dense snippet, right: sparse snippet.

It is also possible to look at the kernel location inside the array (which can be modified with constraints), as shown in Figure 4.5. Moving the cursor on the green lines, the latency from the interface tile to the destination appears: in the picture is 12 ns.

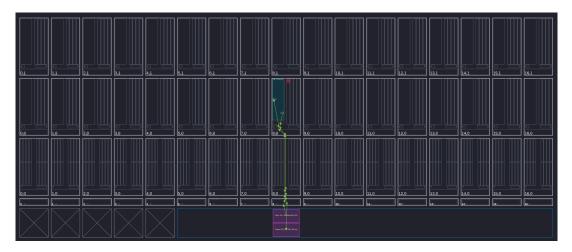


Figure 4.5. AIE-ML array. The device considered has a 17x2 AIE matrix (the lower one is the interface row).

With the trace analyzer (Figure 4.6), it is possible to have a look inside the AIE array, and to know in each clock cycle what is happening in every tile.



Figure 4.6. Trace results of the IIR order 2 custom function on an AI-ML tile.

After a first stall due to missing input, from the second call on, the latency is periodic. The blue segment in the first row is the main function, called every time before the IIR2 function, that has a latency of 155 cycles.

### 4.1.2 Enlarging the kernels

Since moving continuously from the AIE array to PL and vice versa may not be optimal and causes a non-negligible overhead, it is useful to think about how to enlarge the kernel in some way. Furthermore, if the kernel is bigger but the added operations can be executed in parallel, the overall latency may not increase linearly but in a slower way.

#### Enlarging the kernel merging filters

According to appendix A.1, two consecutive IIR filters can be considered as a single IIR filter, whose degree is the sum of the initial two. This idea can be used both to increase the flexibility of the system, for example, enlarging the degree of each filter, or maybe to reduce their number, since every time two (or more) of them are connected in series, they can be merged into a larger one.

This has been investigated: starting from 2 IIR filters of order 2, a larger IIR filter of order 4 has been created and, adopting a similar strategy as before, 169 cycles latency has been obtained. This value is very similar to the order 2 one, making this new approach very promising.

### Enlarging the kernel in parallel

Since in the DFG a group of three filters to be executed in parallel appears many times, the next trial has been to create a kernel enlarged in that direction. In particular, since a sub-block of the graph is repeated three times, it seemed a good idea to try to fit inside an AIE as a kernel the whole block, from now on named "controller block" (Figure 4.7).

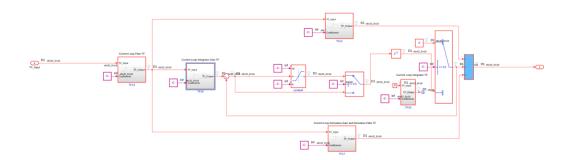


Figure 4.7. Kernel computational structure of the "controller block".

Here, after a first filter executed normally, 3 filters are computed in parallel, with a similar scheme as the one already described. Before summing all the results together, in the second lane, an integrator with a limiter is added in the code. Since the limiter implies an if...else statement (and so the possibility of taking jumps in the code), determinism has been momentarily put aside to investigate possible performances of this block, which resulted in 567 cycles of Initialization Interval.

Since if statements are not ideal when aiming for deterministic execution time, an optimal version of this kernel should be modified to always last the same, but anyway, the first version is reported in the appendix (code in C.5).

### Feedback input

Another repeated block can be found in the group of 4 filters and a sum concerning the elaboration of the input, the "input block" (Figure 4.8).

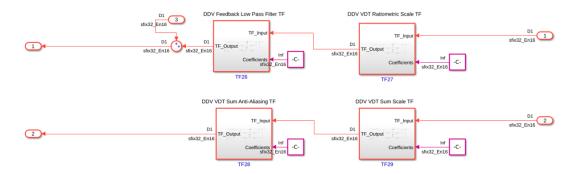


Figure 4.8. Kernel computational structure of the "input block".

This block has two inputs and two outputs, which is not a problem on standard AIE tiles: they can receive two streams as input.

As concerns AIE-ML, only one streaming input is present.

Since the first tests have been performed on AIE-ML tiles, the two 32-bit numbers have been joined to a single 64-bit word input.

This has been implemented, again exploiting a larger parallelism, obtaining an Initialization Interval of 350 cycles.

	Average latency (# assembly instructions), AIE-ML	Average latency (# assembly instructions) AIE
IIR filter order 2 - buffer input	210	NN*
IIR filter order 2 - streaming input (no code optimizations)	220	NN*
IIR filter order 2 - streaming input	155	74
IIR filter order 4	169	NN*
Controller block	567	254
Input block	350	197

Table 4.1. Average latency of implemented custom functions. NN\*: Not Needed; for the normal AIE, only the useful kernels for the Actuator Control Electronic have been tested.

## 4.2 The whole ACE on AIE + PL

Now that the key kernels have been identified, here is a first approach to recreate the whole system. It contains 13 computational kernels:

- 3 x controller blocks,
- 3 x input blocks,
- 7 x second-order IIR filter blocks.

Data enters PL, then moves towards AIEs and comes back to PL. All the passages between PL and AIE have been arranged employing AXI4 handshakes, with the help of the custom block "edge detector" described in sub Section 3.3.1.

Figure 4.9 represents a map of the algorithm that shows how the graph is distributed across the AIE array and PL.

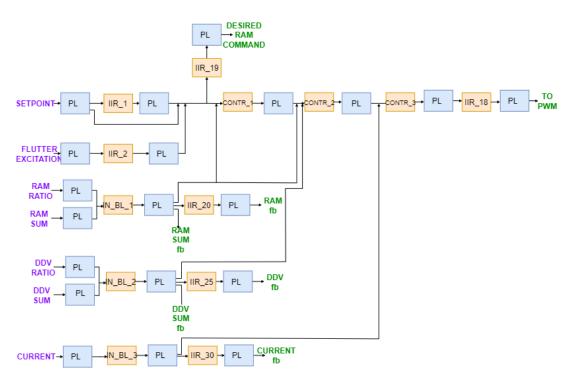


Figure 4.9. ACE functional map. Legend: purple = input; green = output; blue block = input of the PL; orange block = input of an AIE kernel.

This system is translated into two subsystems: the AIE subsystem, where the 13 kernels are instantiated, and the PL subsystem, which contains everything not implemented within the AIE.

### 4.2.1 AIE and PL interconnection

The overall system's block diagram is depicted in Figure 4.10.

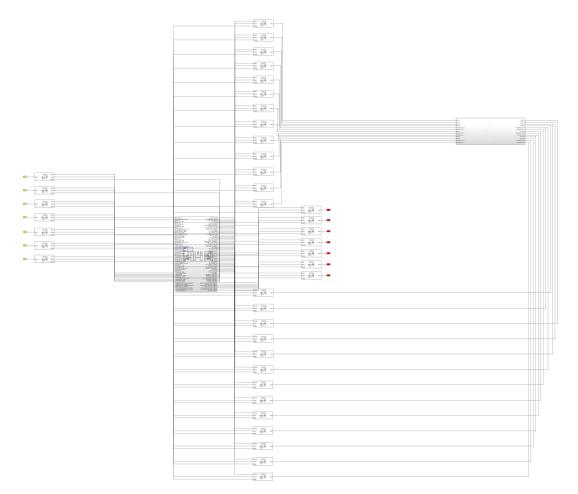


Figure 4.10. Simulink block diagram of the ACE, top level. Yellow ovals are the inputs, while red ones are the outputs. All the interconnections between AIE and PL subsystems are managed through "AIE to HDL" and "HDL to AIE" blocks, which represent AXI4-Stream handshakes.

Since on the AIE array 13 kernels are instantiated, 26 AXI-Stream interconnections are used to connect the AIE array and the PL. The problem is that in this system, only one sample will be sent to AIE for elaboration at a time, but as shown in Figure 2.3, AIE-PL streaming communication interfaces are 64-bit wide. If just one 32-bit sample is sent, the interface waits for the second one to arrive, and so the whole graph is paused. To overcome this problem, the AIE kernels have been modified in order to work on 64-bit input words, of which just the lowest 32 are the real input. The first thing performed in the kernel is an unmasking of the input, while the last is the composition of a 64-bit output word, of which just the lowest 32 are valid.

This may seem a large number of interconnections, but considering that in the xcvc1902's AIE array there are 39 PL interface tiles, with a total of 234 32-bit AXI4-Stream input and 312 32-bit output, it is less than 6% of the overall resources.

### 4.2.2 AIE subsystem

The AIE subsystem is pretty simple: the graph is composed of just 13 kernels in parallel, each of which contains the C++ code reported in the appendix (properly modified in

order to host the correct filter(s)).

Each kernel input and output pass through a PLIO block, in which it is possible to set the interface frequency.

A picture of the AIE subsystem is reported in Figure 4.11, which in the AIE array is translated into the representation in Figure 4.12.

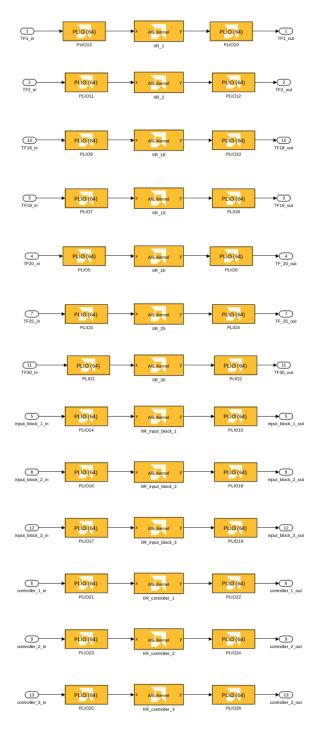


Figure 4.11. Simulink block diagram of the AIE subsystem.

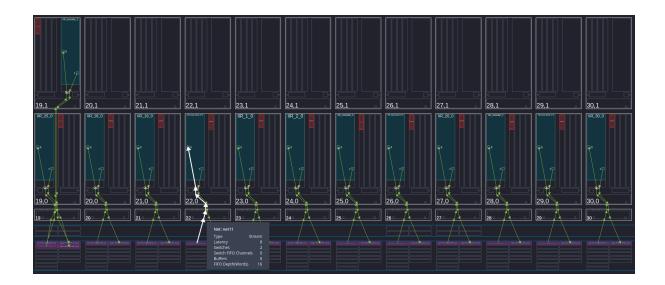


Figure 4.12. AIE graph allocation on the array. Red rectangles are the used portions of the memory, while green arrows represent the data transfer.

As shown in Figure 4.12, the graph is mainly allocated in the central columns of the first row of the array, and data transfer from/to the interface tiles lasts 8 ns for what concerns the first row kernels and 12 ns for the ones located in the second row.

An AIE simulation has been performed, and the trace analyzer results are reported in Figure 4.13. However, since this simulation does not take into account the PL latencies but feeds the array as if all the inputs are furnished in a stream way (Figure 4.14), this is not what happens during the real execution, where only one sample is elaborated at a time.

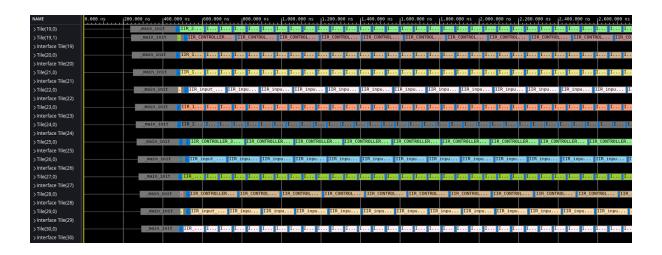


Figure 4.13. Trace analyzer results for the AIE array.

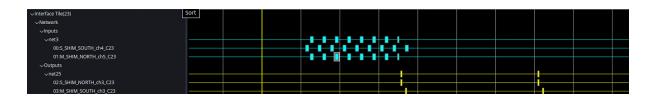


Figure 4.14. Feeding of array tile 23.

As concerns kernel latencies, they are the ones reported in table 4.1 (the table ones do not consider the initial setup). The small "main" calls (blue) between kernel executions are given by data transfers.

### 4.2.3 PL subsystem

The most challenging part of the design, once all the kernels were functioning as expected, has been the PL subsystem — particularly the complex mechanism chosen for managing handshakes.

The subsystem, which is reported in Figure 4.15, contains 20 AXI-Stream inputs (7 from the external world and 13 from AIEs).

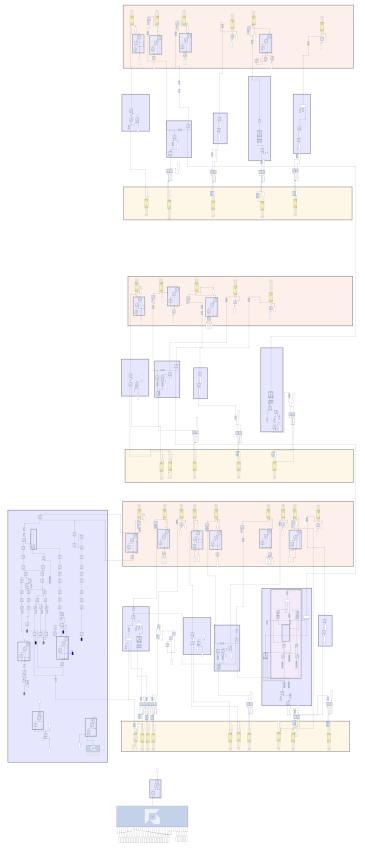


Figure 4.15. Simulink block diagram of the PL subsystem. Blue rectangles contain the logic instantiation, yellow rectangles contain the input logic, and red rectangles contain the output logic.

Data arrives from the external world and AIEs as floating point, but since it is more efficient to work with fixed point data types inside the PL, it has been chosen to cast all the input into signed fixed point of 32 bits (16 bits of decimal part). Then, data width is managed in order not to lose any information until the final re-conversion to the floating world.

### **TVALID-TREADY** management

In order to speed up the clock frequency, a pipeline approach has been introduced. This approach complicated a lot the TVALID management in handshakes: since at each PL output port TVALID should be raised for just 1 cycle, the scheme for its generation should be the one reported in Figure 4.16: a logical AND between all the input from which the output depends, delayed by the same number of cycles as the longest combinatorial path, to be fed to an edge detector block.

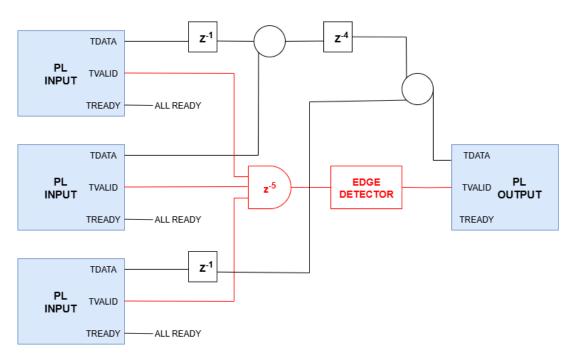


Figure 4.16. Block diagram example for TVALID management in PL: in this case, the output TVALID is generated starting from an AND between all the inputs, delayed as the longest path.

Figure 4.17 shows the TREADY generation: since it is a TREADY whose assertion allows the new set of input to enter the graph, it is generated as a logical AND of all the TREADY (all outputs must be accepted) and all the TVALID (all the graph must have finished execution).

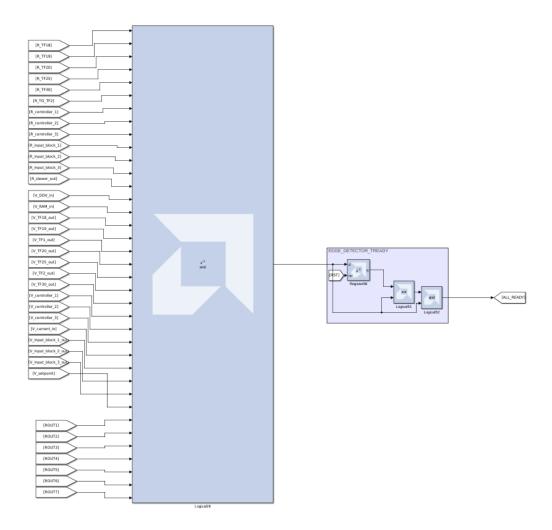


Figure 4.17. TREADY generation for all the external world input blocks of the PL subsystem.

This TREADY signal is used just for the 7 input interfaces coming from the external world, while for the 13 outputs of AIE, it is fixed to a logical '1', since if just one sample is elaborated at a time, the PL is always ready to accept a new value. Since the input TVALID should stay high at least until the data is no longer necessary, the ALL\_READY signal is used to reset a register for TVALID coming from each AIE output, as shown in Figure 4.18.

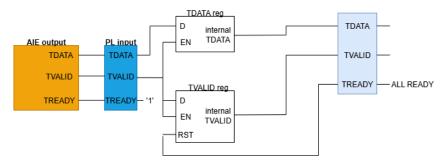


Figure 4.18. Block used to store the output of AIE inside the PL. It introduces 1 extra cycle of latency. For the 7 external inputs, these registers are not necessary.

In Figure 4.19 is reported a simple example waveform of TVALID/TREADY management.

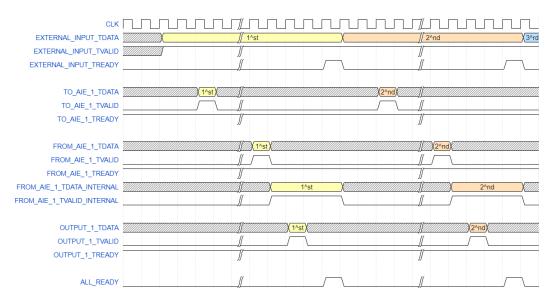


Figure 4.19. Waveform for TVALID-TREADY internal management. It is supposed that external inputs are always valid, while inputs coming from AIE need the TREADY always high, so to keep the validity high, two internal signals are required (as shown in Figure 4.18). ALL\_READY signal allows the reset of these registers and a new set of input data to enter the graph.

It is supposed that in the real application, a new set of input arrives every 100 us, but for simulation purposes, the new elaboration will start just as the previous one has finished, and so the sampling time is given by a logical AND of the validity of the 7 inputs.

### 4.2.4 Testbench & Simulink simulation

As a first approach to debug the system, a Simulink testbench has been created, which contains both the original model and the Model Composer one. Both the models receive the same setpoint, but each one has its own actuator model to generate the feedback signals (Figure 4.20).

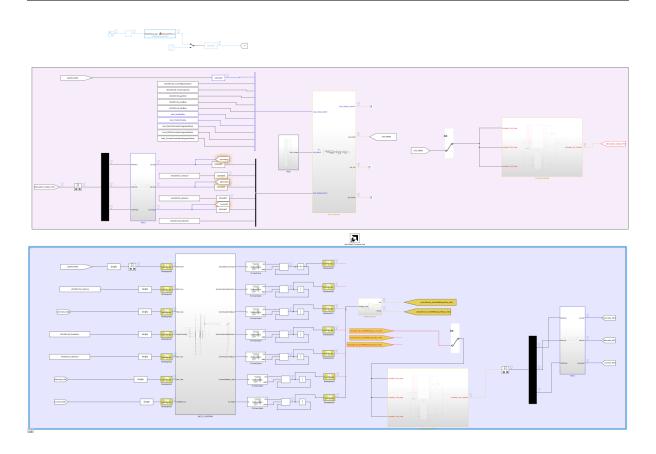


Figure 4.20. Testbench of the ACE.

The first simulation run connected the feedback generated by the original model to the input of both models, such that all 7 inputs were the same, and the result of this simulation is shown in Figure 4.21.

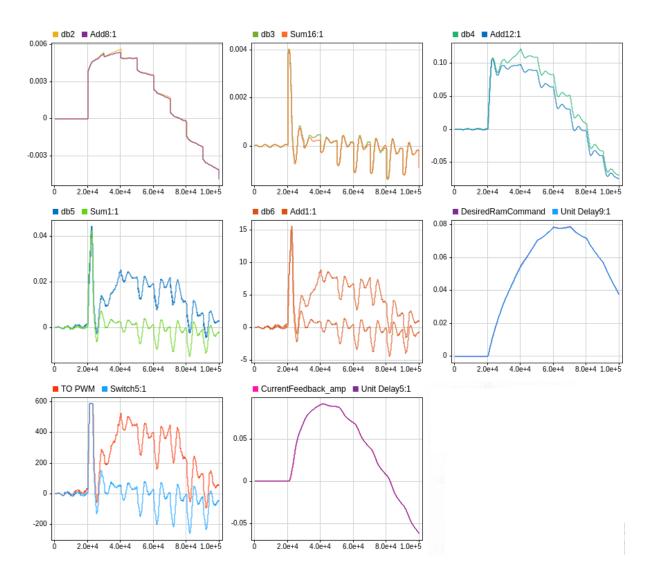


Figure 4.21. Results of the testbench where both models receive the same input. x-axis is in  $\mu s$ . In the legend of each graph, the first signal is the one coming from the PL subsystem, the second one comes from the reference model.

It is quite evident that the results are quite divergent, so some debug probes have been inserted in the PL design in order to understand this behavior (Figure 4.22).

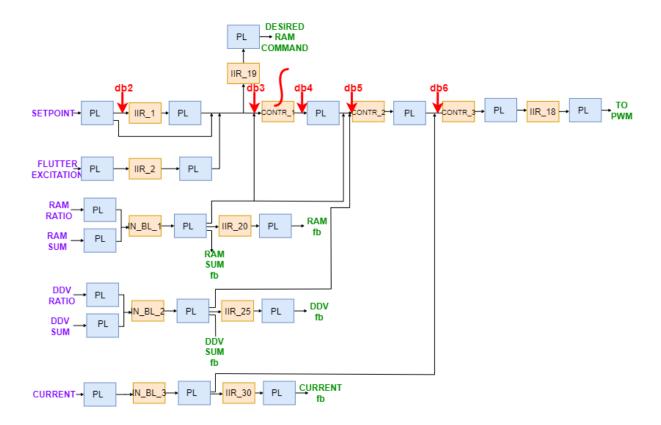


Figure 4.22. Debug probes position in the graph. The block CONTROLLER\_1 is the one responsible for the integration.

A little error, introduced by the slewer logic, is integrated in db4; then, after a subtraction, the relative amplitude of the error increases in db5, and is translated to TO\_PWM output, which can be considered as the ultimate output of the ACE since, to generate it, the entire graph must have been traversed.

The original cause of the divergence is the slewer, but why? The slewer is a block needed in order to mitigate step-like behavior given by the fact that the setpoint changes once every 100 executions of the graph. In Figure 4.23 are reported the block diagrams of the slewer model and implementation on the PL.

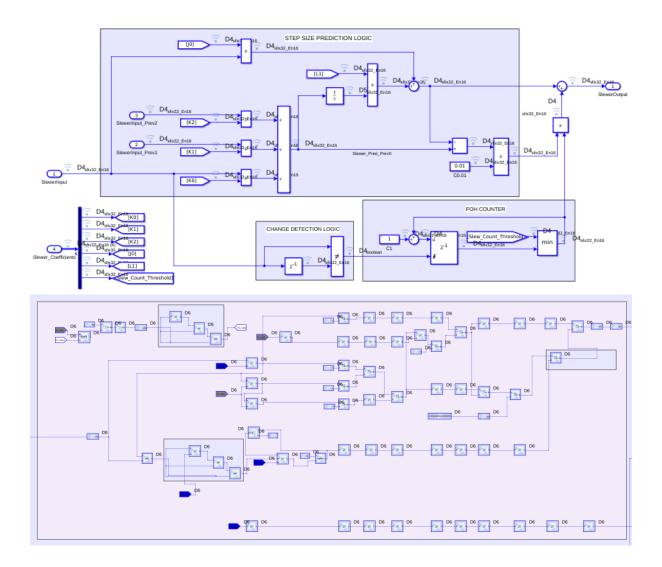


Figure 4.23. Slewer block diagram. Above: original model, below: PL one.

The difference between the two systems is the data type: in the original model, all the operations are performed on sfixdt(1,32,16), and even after multiplication, the output is floored to the nearest representable value (only 16 bits of decimal part). The represented behavior is not correct, since values are typically small numbers, which means that always keeping 16 bits of decimal part may cause significant relative errors. Since the result is rounded towards zero, when the value is positive, the rounding error is always negative, while when the output is negative, the rounding error is positive. This also explains why the gap between the two curves gets narrower when the input becomes negative.

This behavior is not acceptable because such different values for the PWM will cause the system to behave unpredictably.

Things change when the feedback of the Model Composer system is connected to the one generated by its actuator. Such a simulation produces the output shown in Figure 4.24. As shown there, the bias error has disappeared because of the negative feedback

■ db3 ■ Sum16:1 ■ db4 ■ Add12:1 0.006 0.004 0.10 0.003 0.05 0.002 0 -0.003 -0.05 4.0e+4 6.0e+4 8.0e+4 1.0e+5 ■ DesiredRamCommand ■ Unit Delay9:1 ■ db5 ■ Sum1:1 ■ db6 ■ Add1:1 0.08 15 0.04 0.06 0.02 0.04 0.02 8.0e+4 1.0e+5 ■ TO PWM ■ Switch5:1 ■ CurrentFeedback\_amp ■ Unit Delay5:1 400 0.05 200

effect, and now the numerical behavior is the same.

-200

Figure 4.24. Results of the testbench where each model receives its own feedback. x-axis is in  $\mu s$ . In the legend of each graph, the first signal is the one coming from the PL subsystem, the second one comes from the reference model.

4.0e+4

6.0e+4

8.0e+4 1.0e+5

2.0e+4

-0.05

8.0e+4 1.0e+5

At this point, since the behavior of the system is correct, the only remaining thing is to generate a BOOT.BIN to be loaded into the device and check whether the behavior is the same also on the real hardware, and then check its performance against the normal FPGA design.

# 4.3 Tests on the board - simple example

Since Model Composer's simulations are only bit-accurate concerning the AI Engines (AIEs), some tests were performed on the VCK190 board to estimate actual latencies and verify result correctness.

Manually managing correct AXI-Stream handshakes is not trivial, but it allows for greater control over the synthesized circuit. As the initial hardware results did not match those from the Simulink simulations, a simpler system was chosen to ease the debugging process.

This simplified system integrates both the Programmable Logic (PL) and the AI Engine (AIE), effectively serving as a scaled-down representation of the complete project architecture. It enables early validation of core functionalities and design choices.

#### 4.3.1 Procedure

Model Composer automatically generates software for the Processing System and creates two HLS DMAs in order to move input from DDR toward the IP under test and to compare them with the golden output.

The problem is that with this procedure, it can be established only the correctness of the output, but not the latency of the computation. Furthermore, in case the system does not work correctly, debugging is extremely difficult.

For this reason, it is useful to open the Vivado project generated by Model Composer and integrate an ILA (Integrated Logic Analyzer) in order to take a look at the real waveforms. If simulation results are different from hardware simulation, probably the error lies in the AIE-to-PL interfaces, since Simulink does not take into account the AXIS interface's latency or AIE's latency.

Here is reported just an example of serial test results generated according to Model Composer's PS software:

```
Mismatch found for output signal gateway_out_axis in sample 11.Hardware

→ output: 3f2b6000,Model Composer output:3eaf3800

Mismatch found for output signal gateway_out_axis in sample 12.Hardware

→ output: 3f5f9a00,Model Composer output:3efc1000

Mismatch found for output signal gateway_out_axis in sample 13.Hardware

→ output: 3f5f9a00,Model Composer output:3efc1000

Mismatch found for output signal gateway_out_axis in sample 14.Hardware

→ output: 3f8d5b00,Model Composer output:3f2b6000

...
```

As shown there, it seems like some hardware result is repeated, but in order to understand what is happening, ILA is necessary.

Here are the steps to be followed to introduce the ILA block inside the design:

- Open the Vivado project generated by model composer (the project is located in code/platform/work/\_x/link/vivado/vpl/prj/prj.xpr
   ).
- 2. Add the ILA IP from the catalog inside the block diagram.
- 3. Double-click the ILA and set it as necessary. To monitor PL interfaces, it should be set as "interface monitor" for axis RTL.
- 4. Connect the clock, the reset, and the input signals.
- 5. Generate a new device image (this may take a while).
- 6. Copy the just-created vitis\_design\_wrapper.pdi file in the /code/hw folder.
- 7. Create a copy of boot\_img.bif and rename it (for example boot\_img\_custom.bif).
- 8. Inside the boot\_img\_custom.bif, change the .bin file name, substituting ". with '\_', also removing the path. Then, locate it and copy it inside the /code/hw folder.
- 9. Repeat the procedure for the .pdi file, but in this case it has to be substituted with the just created vitis\_design\_wrapper.pdi.
- 10. Run the Petalinux settings.sh script, to use it from the shell.
- 11. Launch

```
bootgen -arch versal -image <path to /code/hw>/boot_image_custom.bif \ \hookrightarrow \  -o <path to /code/hw>/BOOT.BIN -w on
```

Following these steps, thanks to the ILA, the system can be easily debugged to find problems.

#### Trigger at startup

If the whole execution of the system is very short, and the system is just prepared to run once, the ILA has to be triggered at startup.

To do so, as shown in UG908 [12], the following passages have to be added, to be performed on the TCL console in Vivado and the VCK190 programmed via JTAG with the previously produced BOOT.bin:

- 1. open the ILA dashboard and set the trigger condition
- 2. export the register map of the ILA core with

```
run hw ila -file ila_trig_1.tas [get_hw_ilas hw_ila_1] -force
```

- 3. open the implemented design and apply the trigger design settings with apply\_hw\_ila\_trigger ila\_trig\_1.tas
- 4. save the modified constraints file
- 5. write the device image from the TCL console with write\_device\_image vitis\_design\_wrapper.pdi -force
- 6. re-perform the steps 6-11 of the previous section.

### 4.3.2 Debug of AXI handshake in PL

This procedure has been applied to a simple not-working example of Model Composer project, reported in Figure 4.25.



Figure 4.25. Block diagram of the simple system.

Two identical FP32 inputs are provided to two identical AIE kernels. Those kernels are second-order IIR filters, whose outputs are sent to PL. Here, they are summed together after a conversion to fixed-point numbers. The result is again converted into SPFP and then sent to an AXI4 interface. The adder is implemented with 1 cycle of latency, to break the critical path, and as clock frequency 100 MHz was chosen. The PL subsystem is reported in Figure 4.26.

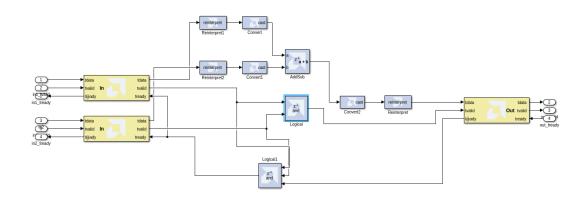


Figure 4.26. Initial block diagram of the PL region of the simple system.

ILA tests in this case helped to discover a problem in TREADY/TVALID management: every two inputs, three outputs were produced instead of two (as shown in Figure 4.27).

The solution to this problem is to sample the output TVALID as a logic AND of the input block's TVALID and TREADY. In this example, 1 cycle latency has been introduced, which is, of course, wrong, but it is also a good example for showing the testbench reaction to errors and ILA debugging of wrong waveforms.

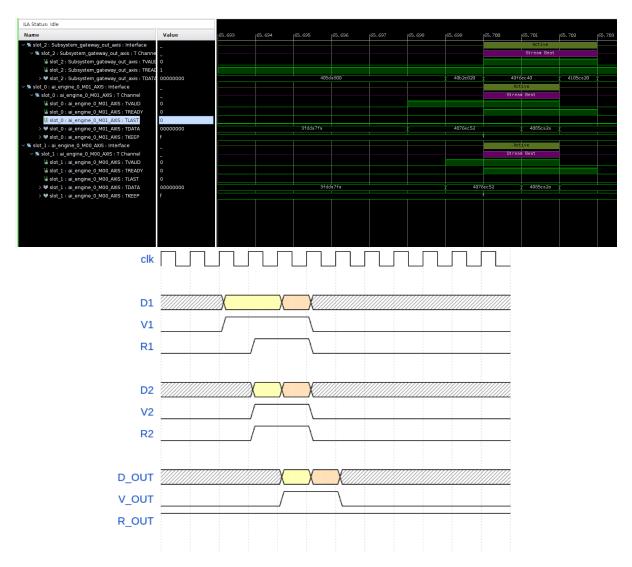


Figure 4.27. ABOVE: wrong initial ILA results. SLOTs 0 and 1 are the input, SLOT 2 is the output. BELOW: theoretical right ones. Each color represents a group of data processed in the same temporal window or computational epoch.

To eliminate timing violations problems (and to find the highest frequency at which the PL can work) from possible errors, a timing analysis has been performed.

Performing a timing analysis with a 3.2 ns PL period results in a timing violation, even introducing other registers in the cast blocks to further break the critical path (now there is a 3-cycle latency), but with 6.4 ns, all timing requirements are met, as shown in Figure 4.28.

The last modification was the introduction of an enable block for the internal PL logic. Without it, if TREADY was de-asserted, the three samples in the pipeline would have continued propagating to the output, which would not have been ready to accept them.

	Violation type: (settup v										
	Slack (ns)	Delay (ns)	Logic Delay (ns)	Routing Delay (ns)	Levels of Logic	Source	Destination	Source Cloc	k Destination Cl	ck Path Constraints	
-1	-2.2720	5.1450	2.13	50 3.0100		23 base_kernel_finale/AIE_PL/PL/AddSub	base_kernel_finale/AIE_PL/PL/Convert2	clk	clk	create_clock -name clk -period 3.2 [get_ports clk]	
2	1.3440	1.5190	0.73	90 0.7800		8 base_kernel_finale/AIE_PL/PL/Conver	1 base_kernel_finale/AIE_PL/PL/AddSub	clk	clk	create_clock -name clk -period 3.2 [get_ports clk]	
3	2.2150	0.6170	0.38	70 0.2300		0 base_kernel_finale/AIE_PL/PL/Logical	base_kernel_finale/AIE_PL/PL/Logical	clk	clk	create_clock -name clk -period 3.2 [get_ports clk]	
4	2.3700	0.5200	0.14	10 0.3790		1 base_kernel_finale/AIE_PL/PL/Logical	L base_kernel_finale/AIE_PL/PL/Logical	clk	clk	create_clock -name clk -period 3.2 [get_ports clk]	
5	2.5110	0.3620	0.09	00 0.2720		0 base_kernel_finale/AIE_PL/PL/AddSub	base_kernel_finale/AIE_PL/PL/Convert2	clk	clk	create_clock -name clk -period 3.2 [get_ports clk]	
Viol											
	Slack (ns)	Setup ▼  Delay (ns)	Logic Delay (ns)	Routing Delay (ns)	Levels of Logic	Source	Destination	Source Clock	Destination Clo	Path Constraints	
1		Delay (ns)		Routing Delay (ns)	_	Source base_kernel_finale/AIE_PL/PL/AddSub	Destination base_kernel_finale/AIE_PL/PL/Convert2		Destination Clo	Path Constraints create_clock -name clk -period 6.4 [get_ports clk]	
1 2	Slack (ns)	Delay (ns) 5.1450	2.1350		23		base_kernel_finale/AIE_PL/PL/Convert2				
1 2 3	Slack (ns) 0.928	Delay (ns) 30 5.1450 40 1.5190	2.1350 0.7390	3.0100	23	base_kernel_finale/AIE_PL/PL/AddSub	base_kernel_finale/AIE_PL/PL/Convert2	clk	clk	create_clock -name clk -period 6.4 [get_ports clk]	
1 2 3 4	Slack (ns) 0.928 4.544	Delay (ns) 30 5.1450 40 1.5190 50 0.6170	2.1350 0.7390 0.3870	3.0100 0.7800	23 8 0	base_kernel_finale/AIE_PL/PL/AddSub base_kernel_finale/AIE_PL/PL/Convert1	base_kernel_finale/AIE_PL/PL/Convert2 base_kernel_finale/AIE_PL/PL/AddSub base_kernel_finale/AIE_PL/PL/Logical	clk	CIK CIK	create_clock -name clk -period 6.4 [get_ports clk] create_clock -name clk -period 6.4 [get_ports clk]	

Figure 4.28. Timing analysis results of the simple subsystem. Above the requirements are not met  $(T_{ck} = 3.2 \text{ ns})$ , bottom the requirements are met  $(T_{ck} = 6.4 \text{ ns})$ .

With these couple of modifications, with the PL subsystem reported in Figure 4.29, the serial test result of the VCK190 reports:

\*\*\*\*\*\* Model Composer and Hardware output match for all 47 samples for output

ignal gateway\_out\_axis

\*\*\*\*\* TEST PASSED

\*\*\*\*\* VMC\_TEST\_DONE

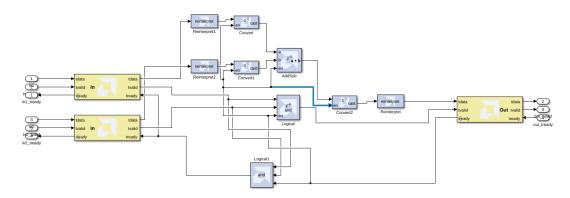


Figure 4.29. Correct block diagram of the PL region of the simple system.

As concerns the ILA in this working experiment, in Figure 4.30 is reported the correct AXIS handshake waveform.

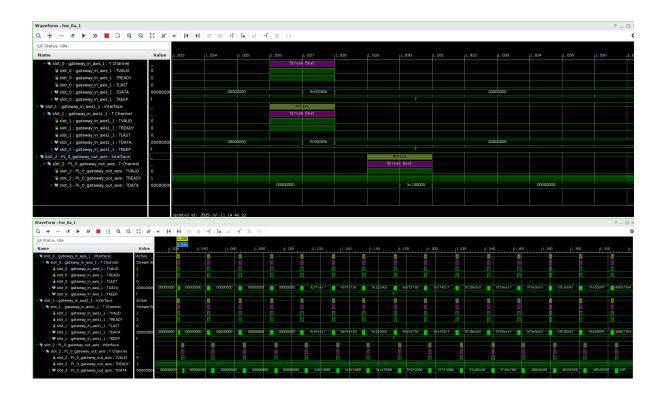


Figure 4.30. Correct behavior registered with ILA. Above: zoom on a single transition, below: subsequent transmissions.

A couple of observations deserve to be pointed out about these waveforms:

- the latency now is 3 cycles, since two more stages have been introduced;
- stream beats are composed of 2 samples. This is because the AIE-PL streaming interface is 64-bit wide, so 2 is the minimum number of 32-bit words that can be packed together;
- after the first four beats, in idle condition, the bus is not empty anymore, but hosts an old value: it is as if the compiler used 8 different paths to stream the data cyclically, and so when it is time for the ninth data, before the actual transmissions, the bus still contains the first data;
- trying to measure throughput, the 2-sample beat is received every 21 or 22 clock cycles. This may be due to the different frequencies: the AIE interface is working at 312.5 MHz, while PL is at 156.25 MHz. The sampling may happen sometimes before and sometimes after the arrival time from the different clock domains.

#### Throughput and latency

In order to estimate the throughput of the system, it is known that  $f_{ck} = 156,25\,MHz$ , so  $T_{ck} = 6.4\,ns$ . Two outputs are ready in around 21.5 cycles, this means that the Iteration Interval  $II = \frac{6.4ns \cdot 21.5}{2} = 68.8\,ns$ . This should also be the latency of AIE computational kernels, since for a new execution to start, the previous must end.

Looking at the assembly code of an IIR filter order 2, the number of assembly instructions is 80, the clock frequency is  $f_{ck} = 1.250 \, GHz$ , and so the esteemed execution time inside the AIE is 64 ns, and since that should be the main contribution of the II, things

make sense.

Now, it is interesting to take a look at the overall latency from input to output introduced by the whole system. The components of the latency are:

- Memory to AIE data transfer (NoC).
- Intra-AIE NoC's latency.
- AIE kernel computational latency.
- Intra-AIE NoC's latency.
- AIE to PL data transfer latency.
- PL latency.

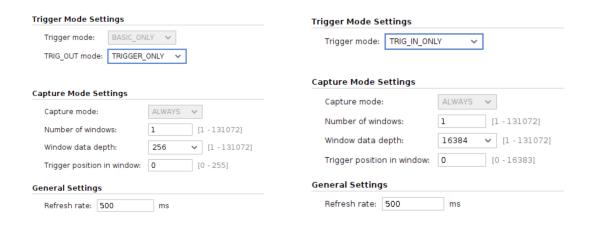
As concerns intra-AIE NoC latency, its contribution can be found thanks to Vitis Analyzer, and is  $8 \times 0.8$  ns = 6.4 ns, while PL latency is 3 cycles (19.2 ns). The NoC and the AIE to PL data latency cannot be measured, but their value is probably negligible with respect to the rest.

#### 4.3.3 Direct measure

A direct measure of latency given by the AIE array (2 x Intra-AIE streaming + kernel execution) can be obtained by inserting two ILAs: one at the input of AIE and one at the output. The first one can be used to trigger the second one.

The problem is that the trigger should be the same for both the ILAs to measure how many cycles of intercourse occur between the input and output of the AIE array.

To do so, the ILAs should be configured as shown in Figure 4.31.



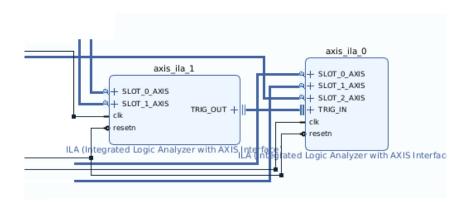


Figure 4.31. Settings of ILAs in order to measure with the same trigger. On the left, the first one, on the right, the second one, below the block diagram (axis\_ila\_1 is the first, while axis\_ila\_0 is the second).

The trigger condition can be set as TDATA! = 0. Then, the same procedure applied for the trigger at startup has to be repeated, doubling each command to configure the registers of both ILAs, and the constraints file is automatically updated, configuring the ILAs to arm the trigger at startup.

Measuring the input of the array with the ILA, the waveform in Figure 4.32 is measured.

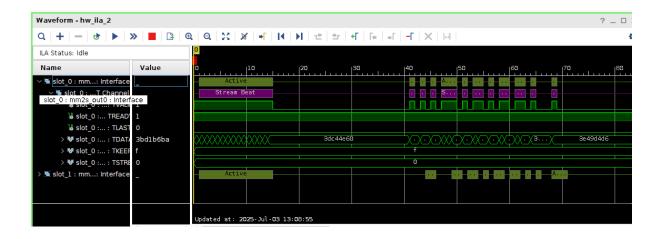


Figure 4.32. Waveforms of the AIE array's input.

The data transfer registered here is between the DDR and the AIE array, performed with the Model Composer's automatically generated HLS DMA mm2s. Probably the compiler of the NoC retained the most efficient start with a 16-sample beat (that probably is filling an input buffer of the AIE array), to wait a little and start again with smaller beats later.

The second ILA was set up to trigger at the same time as the first one, but looking at the waveforms, it is clear that the AIE elaboration's start is not registered (no rising edge of TVALID detected at the output of the AIE, at least in the 0.85 ms after the trigger condition).

Looking at the code host.cpp running on the PS (generated by Model Composer), the reason for that is clear: Model Composer automatically programs the DMAs before initializing and starting the graph, as shown here:

```
Xil_Out32(mm2s_base + M_INO_OFFSET, (uint32_t) mem_inOAddr);
Xil_Out32(mm2s_base + M_INO_OFFSET + 4, 0);
Xil_Out32(mm2s_base + M_INO_SIZE_OFFSET, input0_size);

Xil_Out32(mm2s_base + M_IN1_OFFSET, (uint32_t) mem_in1Addr);
Xil_Out32(mm2s_base + M_IN1_OFFSET + 4, 0);
Xil_Out32(mm2s_base + M_IN1_SIZE_OFFSET, input1_size);

Xil_Out32(s2mm_base + S_OUT0_OFFSET, (uint32_t) mem_outOAddr);
Xil_Out32(s2mm_base + S_OUT0_OFFSET + 4, 0);
Xil_Out32(s2mm_base + S_OUT0_SIZE_OFFSET, output0_size);

Xil_Out32(mm2s_base + M_PULSE_OFFSET, 2);
Xil_Out32(mm2s_base + M_CTRL_OFFSET, 1);
Xil_Out32(s2mm_base + S_CTRL_OFFSET, 1);

printf("AI Engine graph init\n");
mygraph.init();
```

```
printf("AI Engine graph run\n");
mygraph.run();
```

and the delay given by graph initialization is greater than the maximum ILA window at  $156.25~\mathrm{MHz}.$ 

Before measuring the result, another modification was made: instead of inserting two ILAs with a cascade trigger, just one ILA was inserted with a clock conversion and AXI broadcasters, to visualize all the waves on the same window.

In Figure 4.36, the block diagram set up for this measurement is reported, while in Figure 4.33, Figure 4.34, and Figure 4.35 they are displayed the waveforms of the final modified system.

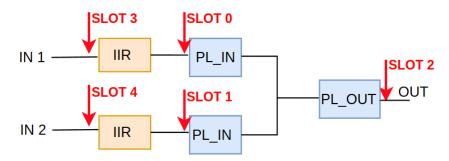


Figure 4.33. ILA slot positions in the system. SLOT 0-1: PL input; SLOT 2: PL output; SLOT 3-4: AIE input.

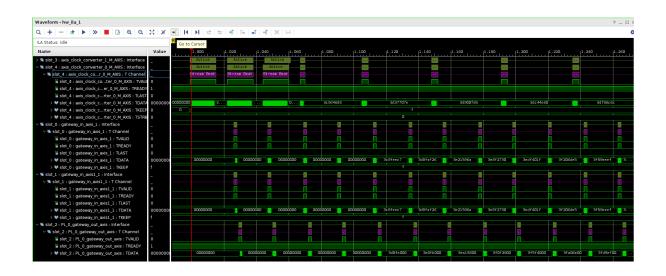


Figure 4.34. Expanded waveforms of the system. Trigger condition: TDATA (slot 3) != 0.



Figure 4.35. Measurement of the latency on the waveforms of the system. Trigger condition: TDATA (slot 3) !=0.

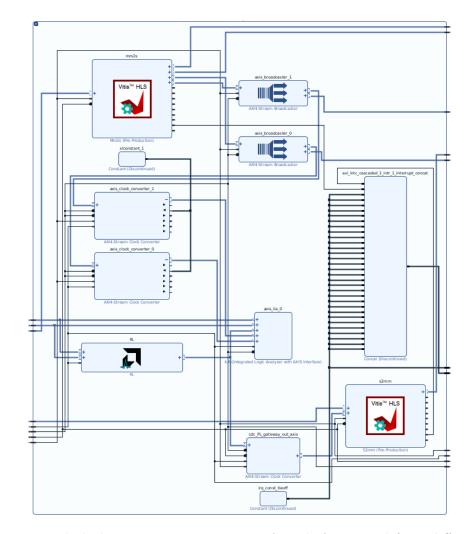


Figure 4.36. Block diagram to conFigure 1 ILA with AXI signal from different clock domains (zoom on Vitis region).

The latency between the input of the PL and the input of AIE is 27 cycles at 156.25 MHz, which corresponds to 172.8 ns, and in this period, the first two outputs are provided. These 172.8 ns are probably given by 59.2 ns for the first kernels' execution, 59.2 ns for the second, and the remaining 54.4 ns for the data transport.

Given this measure, the latency of AIE-PL transport can be estimated to be in the  $(35 \div 50)$  ns range, since two data are transferred, but the 54.4 ns include both AIE-PL data

transfer of the two samples and intra-AIE data transfer, and also considering that the sampling frequency is not the same frequency of the AIE array interface.

It is important to remember that ILA data are registered in PL, so also the delay due to their movement to PL has to be considered; however, these values are precious estimations for what concerns the graph mapping on the SoC.

#### 4.3.4 Other considerations

On Vivado it is possible to take a look at many interesting report, such as the chip planner (Figure 4.37), that shows which portion of the chip are used, a NoC utilization scheme (Figure 4.38) and power reports (Figure 4.39). These power consumptions will be compared with the ACE implementation's one in subsection 4.4.2.

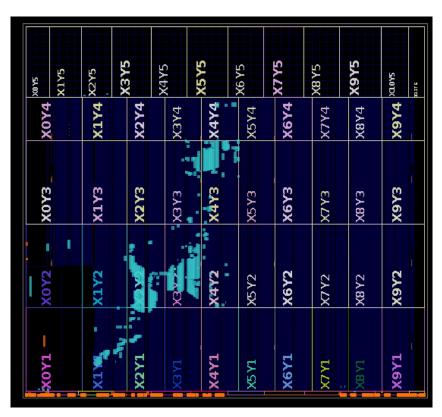


Figure 4.37. Chip plan for the simple example. The AIE region corresponds to the Y5 row of the chip, but on Vivado, even if the tiles are used, just the interface tiles are highlighted.

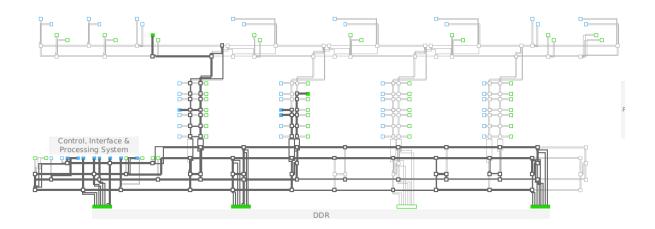


Figure 4.38. NoC utilization for the simple example.

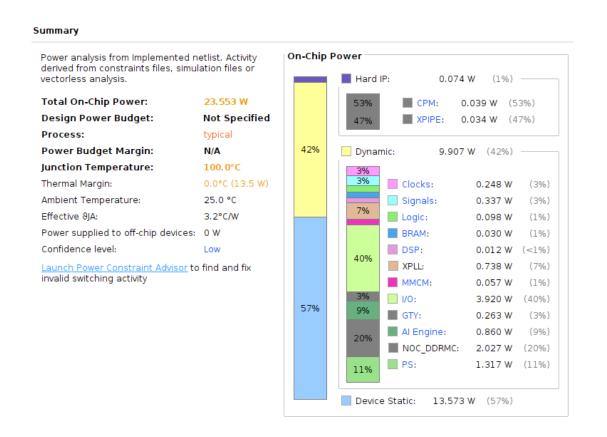


Figure 4.39. Power consumption for the simple example.

Figure 4.40 shows the resource utilization of this simple project.

BLI Registe rs (88264	144	0	144	144	0	0	0	0	0	0	0	0	0	0	0	BLI Registe rs (88264	0.16%	0.00%	0.16%	0.16%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
XPIPE Reg	1	0	-	0	0	0	0	0	0	0	0	0	0	0	0	KPIPE Reg QUA rs D (4) (88;	25.00% 0.1	0.00% 0.0	25.00% 0.1	0.00% 0.1	0.00% 0.0	0.00% 0.0	0.00% 0.0	0.00% 0.0	0.00% 0.0	0.00% 0.0	0.00% 0.0	0.00% 0.0	0.00% 0.0	0.00% 0.0	
	-	0		0	0	0	0	0	0	0	0	0	0	0	0																%00.0 %
PS9 (1)																PS9 (1)	100.00%	0.00%	100.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
GTYES QUAD (11)	-	0	1	0	0	0	0	0	0	0	0	0	0	0	0	GTYES QUAD (11)	9.09%	0.00%	9.09%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
DDR MC_RI U (4)	m	0	m	0	0	0	0	0	0	0	0	0	0	0	0	DDR MC_RI U (4)	75.009	0.00%	75.009	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
DDRM C (4)	m	0	m	0	0	0	0	0	0	0	0	0	0	0	0	DDRM C (4)	75.00%	%00.0	75.00%	0.00%	%00.0	0.00%	%00.0	%00.0	%00.0	0.00%	0.00%	%00.0	0.00%	%00.0	0.00%
CPM MAIN (1)	-	0	1	0	0	0	0	0	0	0	0	0	0	0	0	CPM MAIN (1)	100.005	0.00%	100.00	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
CPM EXT_	-	0	-	0	0	0	0	0	0	0	0	0	0	0	0	CPM_ EXT_ (1)	100.00	%00.0	100.00	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	%00.0	0.00%
NOC Slave (16)	-	0	п	-	0	0	0	0	0	0	0	0	0	0	0	NOC Slave (16)	6.25%	0.00%	6.25%	6.25%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
PL Slav e (312	2	0	2	2	0	0	0	0	0	0	0	0	0	0	0	PL Slav e (312	0.64%	0.00%	0.64%	0.64%	0.00%	0.00%	0.00% 0.00%	0.00% 0.00%	0.00%	0.00%	0.00%	0.00% 0.00%	0.00% 0.00%	0.00% 0.00% 0.00%	0.00% 0.00% 0.00% 0.00%
PL Mast er (234	2	0	2	2	0	0	0	0	0	0	0	0	0	0	0	PL Mast er (234	0.85%	0.00%	0.85%	0.85%	0.00%	0.00%			0.00%	0.00%	0.00%				0.00%
NOC Slave 128 bit (22)	-	0	-	0	0	0	0	0	0	0	0	0	0	0	0	NOC Slave 128 bit (22)	4.55%	0.00%	4.55%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
NOC Maste r 128 bit (26)	œ	0	00	0	0	0	0	0	0	0	0	0	0	0	0	NOC Maste r 128 bit (26)	30.77%	0.00%	30.77%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
NOC Slave 512 bit (28)	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	NOC Slave 512 bit (28)	3.57%	%00.0	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
NOC Maste r 512 bit (28)	е	0	m	0	0	0	0	0	0	0	0	0	0	0	0	NOC Maste r 512 bit (28)	10.71%	%00.0	10.71%	0.00%	%00.0	0.00%	%00.0	0.00%	%00.0	0.00%	0.00%	%00.0	%00.0	%00.0	0.00%
MMCM (12)	н	0	-	0	0	0	0	0	0	0	0	0	0	0	0	MMCM (12)	8.33%	%00.0	8.33%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	%00.0	0.00%
XPLL (24)	O	0	o	0	0	0	0	0	0	0	0	0	0	0	0	XPLL (24)	37.509	0.00%	37.509	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
BUFG CE/MB UFGC E E (296)	т	0	m	0	0	0	0	0	0	0	0	0	0	0	0	BUFG CE/MB UFGC E E (296)	1.01%	0.00%	1.01%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	%00.0
BUFG PS/MB UFG P S (12)	п	0	г	0	0	0	0	0	0	0	0	0	0	0	0	BUFG PS/MB UFG P S (12)	8.33%	%00.0	8.33%	0.00%	0.00%	0.00%	0.00%	0.00%	%00.0	0.00%	%00.0	0.00%	%00.0	%00.0	%00.0
Bond ed 108 (692)	382	0	0	0	0	0	0	0	0	0	0	0	0	0	0	Bond ed 108 (692)	55.20%	%00.0	%00.0	%00.0	%00.0	%00.0	0.00%	%00.0	%00.0	%00.0	%00.0	%00.0	%00.0	%00.0	0.00%
DSP Slice s (196 8)	9	0	9	0	0	0	0	9	0	0	0	0	9	0	0	DSP Slice s (196 8)	0.30%	%00.0	0.30%	0.00%	0.00%	0.00%	%00.0	0.30%	0.00%	0.00%	%00.0	0.00%	0.30%	%00.0	0.00%
Block RAM Tile (967)	63	0	6	0	0	0	0	6	0	0	92	0	1	0	0	Block RAM Tile (967)	9.62%	0.00%	9.62%	0.00%	0.00%	0.00%	0.00%	9.62%	0.00%	0.00%	9.51%	0.00%	0.10%	0.00%	0.00%
CLB Register s (179968 0)	16242	268	15310	176	419	420	2691	11519	83	83	3746	75	5739	132	1657	CLB Register s (179968 0)	0.90%	0.05%	0.85%	<0.01%	0.02%	0.02%	0.15%	0.64%	<0.01%	<0.01%	0.21%	<0.01%	0.32%	<0.01%	0.09%
SUCE (11248 0)	2870	167	2751	33	104	102	569	1950	15	14	674	12	928	104	274	SLICE (11248 0)	2.55%	0.15%	2.45%	0.03%	%60.0	0.09%	0.51%	1.73%	0.01%	0.01%	%09.0	0.01%	0.83%	%60.0	0.24%
LOOKAH EAD8 (11248 0)	498	7	491	0	0	0	S	486	0	0	28	0	402	28	28	LOOKAH EAD8 (11248 0)	0.44%	<0.01%	0.44%	0.00%	0.00%	0.00%	<0.01%	0.43%	0.00%	0.00%	0.02%	0.00%	0.36%	0.02%	0.02%
LUT as Memory (449920)	1153	0	1152	40	0	0	205	903	0	0	155	0	394	-	353	LUT as Memory (449920)	0.26%	0.00%	0.26%	<0.01%	0.00%	0.00%	0.05%	0.20%	0.00%	0.00%	0.03%	0.00%	0.09%	<0.01%	0.08%
LUT as Logic (899840	11027	488	10528	74	378	378	1986	6292	41	42	1987	41	4021	524	1019	LUT as Logic (899840	1.23%	0.05%	1.17%	<0.01%	0.04%	0.04%	0.22%	0.85%	<0.01%	<0.01%	0.22%	<0.01%	0.45%	0.06%	0.11%
CLB LUTs (899840)	12180	488	11680	114	378	378	2191	8582	41	42	2142	41	4415	525	1372	CLB LUTs (899840)	1.35%	0.05%	1.30%	0.01%	0.04%	0.04%	0.24%	0.95%	<0.01%	<0.01%	0.24%	<0.01%	0.49%	0.06%	0.15%
Register s (17996 80)	16386	897	15454	320	419	420	2691	11519	83	83	3746	75	5739	132	1657	Register s (17996 80)	0.91%	0.05%	0.86%	0.02%	0.02%	0.02%	0.15%	0.64%	<0.01%	<0.01%	0.21%	<0.01%	0.32%	<0.01%	0.09%
Name 1	N vitis_design_wrapper	> Te axi_dbg_hub (axi_dbg_hub_axi_dbg_hub_0)	i vitis_design_i (vitis_design)	> I al_engine_0 (vitis_design_ai_engine_0_0)	> I axi_intc_cascaded_1 (vitis_design_axi_intc	> I axi_intc_parent (vitis_design_axi_intc_pare	> Taxi_smc_vlp_hler (axi_smc_vlp_hler_imp_G0	▼ I VitisRegion (VitisRegion_imp_9ASBXH)	> I axis_clock_converter_0 (vitis_design_axi	> [] axis_clock_converter_1 (vitis_design_axi	> TE axis_ila_0 (vitis_design_axis_ila_0_0)	> [ cdc_PL_gateway_out_axis (vitis_design_	> I mm2s (vitis_design_mm2s_0)	> I PL_1 (vitis_design_PL_1_0)	> I s2mm (vitis_design_s2mm_0)	Name	∨ N vitis_design_wrapper	> Te axi_dbg_hub (axi_dbg_hub_axi_dbg_hub_0)	✓ I vitis_design_i (vitis_design)	> I ai_engine_0 (vitis_design_ai_engine_0_0)	> I axi_intc_cascaded_1 (vitis_design_axi_intc	> I axi intc_parent (vitis_design_axi_intc_pare	> I axi_smc_vip_hier (axi_smc_vip_hier_imp_G0	▼ I VitisRegion (VitisRegion_Imp_9ASBXH)	> 🔳 axis_clock_converter_0 (vitis_design_axi	> I axis clock converter 1 (vitis design axi	> 1 axis_ila_0 (vitis_design_axis_ila_0_0)	> I cdc PL gateway out axis (vitis design	> I mm2s (vitis_design_mm2s_0)	> I PL 1 (vitis_design_PL_1_0)	> I s2mm (vitis_design_s2mm_0)
	>	^	>														>	^	>												

Figure 4.40. Resources utilization of the final implementation (absolute values on the left and % on the right).

#### 4.3.5 Estimation of complete system behavior

Based on measurements of this simple example, a possible latency estimation of the whole ACE can be made.

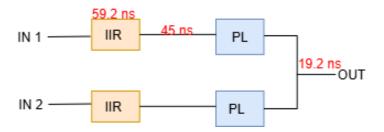


Figure 4.41. Latency measured on the simple example.

In Figure 4.41 are reported the latency measurements performed on the simple system. The kernel latency can be evaluated as  $\#assembly\ instructions \cdot 0.8ns$ , and their values can be found in table 4.1.

As concerns the PL latency, it is known to be exactly  $T_{ck} \cdot \#pipeline\ stages$ .

For intra-AIE + AIE-PL data transfer, it is reasonable to assume 45 ns as an estimation, since probably around 10 ns are spent for intra-AIE transport and  $(35 \div 50)$  for data streaming from AIE to PL or vice-versa (for simplicity, also the latency from the AIE array interface tile to the kernel tile is included in the 45 ns).

The graphical calculation of the overall latency is reported in Figure 4.42.

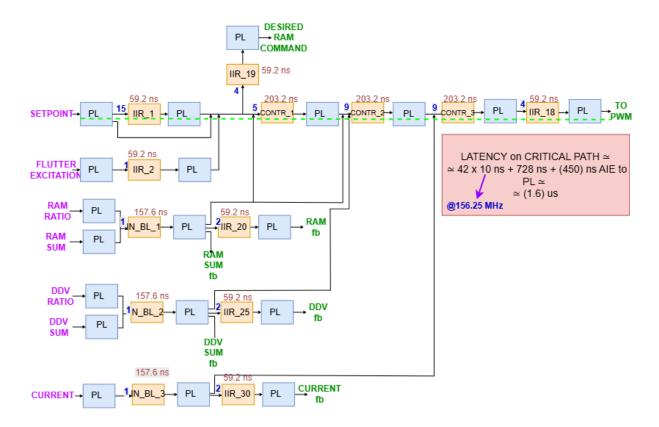


Figure 4.42. Latency estimation for the whole ACE. Blue number identifies the number of clock cycles to reach the destination (considering the input that generates it, as in Figure 4.16), while the green dotted line identifies the longest path.

Even if this algorithm is not much parallelizable, its implementation on Versal Adaptive SoC makes easily possible to obtain a minimum latency of the whole ACE of less than  $2\mu s$ , which is less than 2% of the requirement.

## 4.4 Code generation and test of the ACE

In order to proceed with a test of the ACE model described in Section 4.2, some modification to the Simulink testbench has to be performed: on Model Composer 2025.1, even if not directly synthesized, testbench block such as buses or discrete time integrator breaks the timing analysis and so the whole code generation itself.

Since neither the PL nor the AIE subsystem is a problem itself, once the numerical behavior of the system has been approved, they have been transferred into a new blank project, and the input has been generated randomly, just to produce the golden reference data to populate the testbench for the simulation.

The new Simulink system is reported in Figure 4.43.

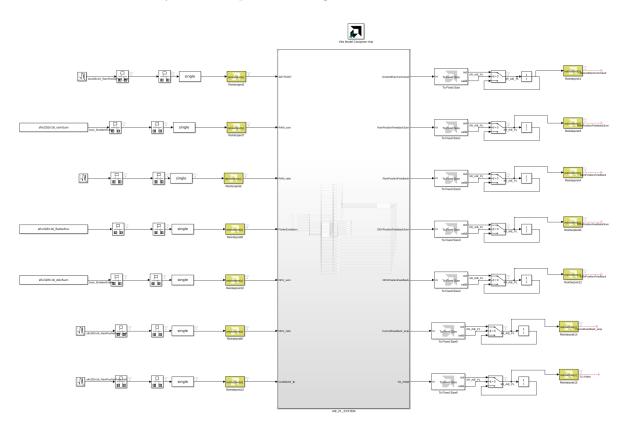


Figure 4.43. Simulink system to generate the .bin file.

The complete serial output generated by the host code (appendix B.1) is:

```
[3.667] *********Boot PDI Load: Started*******
[3.725]Loading PDI from SBI
[3.751] Monolithic/Master Device
[4.009]0.299 ms: PDI initialization time
[4.045]+++Loading Image#: 0x1, Name: lpd, Id: 0x04210002
[4.091]---Loading Partition#: 0x1, Id: 0xC
[58.415] 54.280 ms for Partition#: 0x1, Size: 10544 Bytes
[63.423] --- Loading Partition#: 0x2, Id: 0xB
[104.556] 37.257 ms for Partition#: 0x2, Size: 63008 Bytes
[107.065]+++Loading Image#: 0x2, Name: pl_cfi, Id: 0x18700000
[112.295] --- Loading Partition#: 0x3, Id: 0x3
[2114.011] 1997.753 ms for Partition#: 0x3, Size: 3189232 Bytes
[2116.767] --- Loading Partition#: 0x4, Id: 0x5
[3040.845] 920.029 ms for Partition#: 0x4, Size: 1438112 Bytes
[3043.527]+++Loading Image#: 0x3, Name: cpm, Id: 0x04218007
[3048.760] --- Loading Partition#: 0x5, Id: 0x6
[3052.994] 0.190 ms for Partition#: 0x5, Size: 2224 Bytes
[3057.886]+++Loading Image#: 0x4, Name: aie_subsys, Id: 0x0421C005
[3063.718] --- Loading Partition#: 0x6, Id: 0x7
[3070.677] 2.912 ms for Partition#: 0x6, Size: 1936 Bytes
[3072.930]+++Loading Image#: 0x5, Name: fpd, Id: 0x0420C003
[3078.160]---Loading Partition#: 0x7, Id: 0x8
[3084.568] 2.362 ms for Partition#: 0x7, Size: 4544 Bytes
[3087.286]+++Loading Image#: 0x6, Name: aie_dev_part, Id: 0x18800000
[3093.291] --- Loading Partition#: 0x8, Id: 0x0
[3099.165] 1.826 ms for Partition#: 0x8, Size: 7712 Bytes
[3102.421]+++Loading Image#: 0x7, Name: aie_image, Id: 0x18800000
[3108.167] --- Loading Partition#: 0x9, Id: 0x0
[3160.557] 48.344 ms for Partition#: 0x9, Size: 82400 Bytes
[3163.264]+++Loading Image#: 0x8, Name: default subsys, Id: 0x1C000000
[3169.161]---Loading Partition#: 0xA, Id: 0x0
[4433.118] 1259.907 ms for Partition#: 0xA, Size: 2011344 Bytes
[4435.931] **********Boot PDI LInitializing AIE driver...
Initializing ADF API...
XAIEFAL: INFO: RePLMrce group Avail is created.
XAIEFAL: INFO: Resource group Static is created.
XAIEFAL: INFO: Resource group Generic is created.
Beginning test
in0=188748
in1=188648
in2=1ec790
in3=1ecb90
in4=188e48
in5=1eca90
in6=188d48
out0=244720
out1=244830
out2=244940
out3=244a50
out4=244b60
out5=244c70
out6=244d80
Starting test w/ cu
```

```
AI Engine graph init
Initializing graph mygraph...
AI Engine graph run
Calling adf::graph::run() without specifying number of iterations ... It will
\rightarrow either run AIE cores indefinitely, or run AIE cores for t.
Enabling core(s) of graph mygraph
50 out of 64 samples have been processed for new_desired_ram_cmd...
64 out of 64 samples have been processed for new_ram_sum_out...
64 out of 64 samples have been processed for new_ram_pos_fb...
64 out of 64 samples have been processed for new_ddv_sum_out...
64 out of 64 samples have been processed for new_ddv_position_feedback...
64 out of 64 samples have been processed for new_currentfeedback_amp...
64 out of 64 samples have been processed for new_to_pwm...
***************************** Test Results ***********************
***** Model Composer and Hardware outputs match for all 64 samples for output
\hookrightarrow signal new_desired_ram_cmd
***** Model Composer and Hardware outputs match for all 64 samples for output
\hookrightarrow signal new_ram_sum_out
***** Model Composer and Hardware outputs match for all 64 samples for output
\rightarrow signal new_ram_pos_fb
***** Model Composer and Hardware outputs match for all 64 samples for output
\hookrightarrow signal new_ddv_sum_out
***** Model Composer and Hardware outputs match for all 64 samples for output

→ signal new_ddv_position_feedback

***** Model Composer and Hardware outputs match for all 64 samples for output

→ signal new_currentfeedback_amp

***** Model Composer and Hardware outputs match for all 64 samples for output
\hookrightarrow signal new_to_pwm
***** TEST PASSED
***** VMC_TEST_DONE
```

#### 4.4.1 Results & performances

In order to measure the real latencies on the board, the generated Vivado project has been modified, adding an ILA triggered at startup as shown in Section 4.3.

The new Vivado block diagram is reported in Figure 4.44.

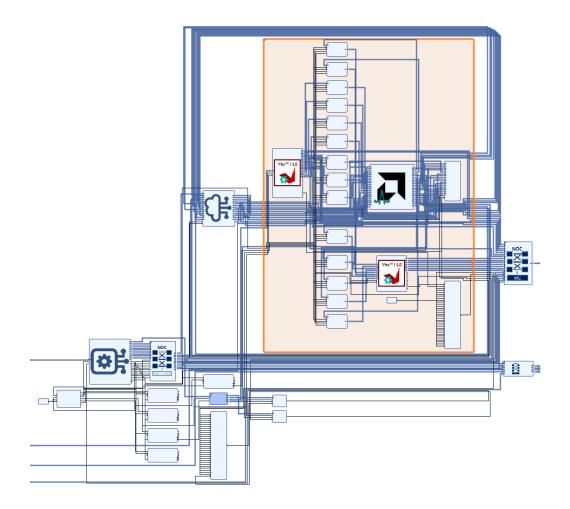


Figure 4.44. ACE project's block diagram.

The base blocks are the same as the simple example one, but the PL subsystem (central block in the orange square) is much more complex.

#### Latency measurement

The ILA is connected to PL-relevant input/output, and therefore, differently from Section 4.3, there is no extra latency due to measurements made outside the PL.

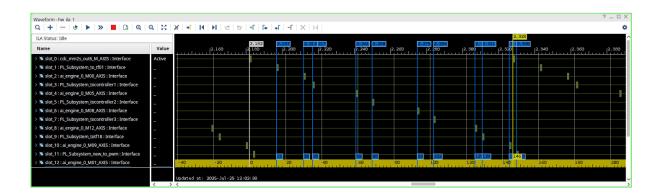


Figure 4.45. Measurement of the latency with the ILA. The scale is in terms of clock cycles (1 cycle = 10 ns). The white checkpoint corresponds to the start of the elaboration, while the yellow line corresponds to the end of the elaboration.

In Figure 4.45, the latency measurement is reported. This is a list of the ILA slots (the ones on the longest path of Figure 4.42) and their latencies from the setpoint ingress inside the PL in terms of PL clock cycles (the frequency is 100 MHz):

- SLOT\_0: SETPOINT input -> 0 cycles;
- SLOT\_1: IIR\_1 input  $\rightarrow$  15 cycles;
- SLOT\_2: IIR\_1 output -> 30 cycles;
- SLOT 3: CONTR 1 input -> 35 cycles;
- SLOT\_4: CONTR\_1 output -> 59 cycles;
- SLOT\_5: CONTR\_2 input -> 68 cycles;
- SLOT\_6: CONTR\_2 output -> 93 cycles;
- SLOT\_7: CONTR\_3 input -> 102 cycles;
- SLOT\_8: CONTR\_3 output -> 125 cycles;
- SLOT\_9: IIR\_18 input -> 129 cycles;
- SLOT\_10: IIR\_18 output -> 144 cycles;
- SLOT\_11: TO\_PWM output -> 148 cycles.

All the latencies are cyclically repeated until all the testbench input data are elaborated, as shown in Figure 4.46.



Figure 4.46. Periodicity of the system registered with the ILA.

The ALL\_READY signal rises before the validity of the output (this event corresponds to the ingress of a new series of input, and in Figure 4.45 corresponds to the yellow line with 146 clock cycles of latency), because all is still needed is for the last branch to flush, and the system is ready to read other inputs.

Since the Initialization Interval is not exactly  $100 \,\mu s$ , both for simplicity and to reduce occupation of the ILA waveforms, the set of 7 outputs is sent outside the PL once every time a new input is received (a logical AND of the 7 inputs TVALID is used as a temporization signal, in absence of an external one).

In Figure 4.47, the synchronicity of the output validity is shown.

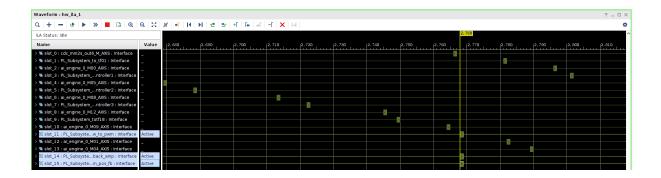


Figure 4.47. Synchronicity of the output validity registered with the ILA. Selected slots correspond to 3 of the 7 outputs (just because it is not possible to register more than 16 signals with ILA).

The overall latency to traverse all the graph is therefore 1.48  $\mu s$ , which is very similar to the estimation made in Section 4.3.5.

#### 4.4.2 Other reports

In Figure 4.48 the chip plan is reported, in Figure 4.49 the NoC utilization, and in Figure 4.50 the power report.

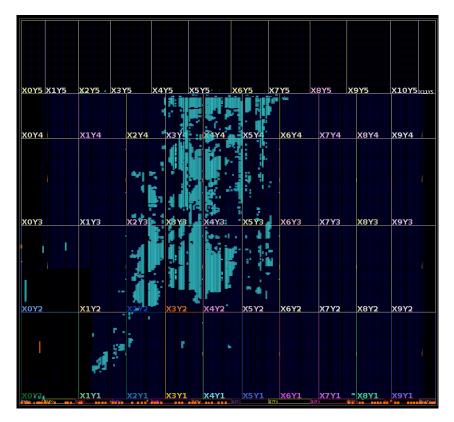


Figure 4.48. Chip plan of the final implementation. The AIE region corresponds to the Y5 row of the chip, but on Vivado, even if the tiles are used, just the interface ones are highlighted.

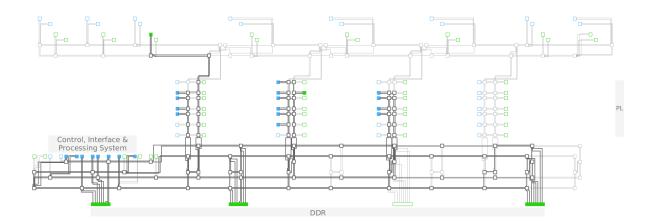


Figure 4.49. NoC utilization of the final implementation. The 14 middle blue squares correspond to PL input/output of the system, while AIE-PL communication does not employ the NoC, and so is not present in the diagram.

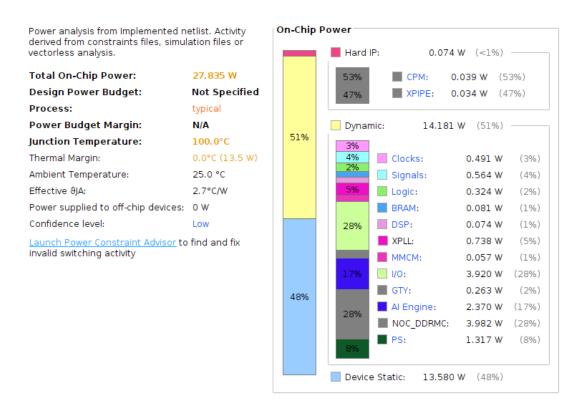


Figure 4.50. Power consumption estimation of the final implementation.

It is interesting to compare the result of the complete ACE implementation with the simple example one (subsection 4.3.4): the chip plan is much more colored, and in the NoC, much more data is traveling (from external DDR to PL input).

As concerns power utilization, the most relevant increases are DSP power (6.2x), logic power (3.3x), AIE (2.8x), BRAM (2.8x), clock (2x), NoC DDRMC (2x), and signals power (17x2). All other contributions remain more or less the same, and for this reason, the overall power consumption switching from the simple design to the complete ACE is just 1.2x.

The power consumption is much greater than the one obtained by implementing the ACE on a traditional FPGA (around 4 W), but by increasing the complexity of the algorithm, Versal ACAPs can achieve far better power consumptions, especially if a higher level of parallelism allows a better exploitation of the AI Engines. The same conclusions have been reached by [15].

However, due to the highly concentrated computation on a single chip, its power consumption may also reach 100 W. Such a delivery, with currents that may exceed 100 A, can be challenging to minimize voltage drop, losses, and temperature rises, as explained in [18].

The overall resource utilization is reported in Figure 4.51.

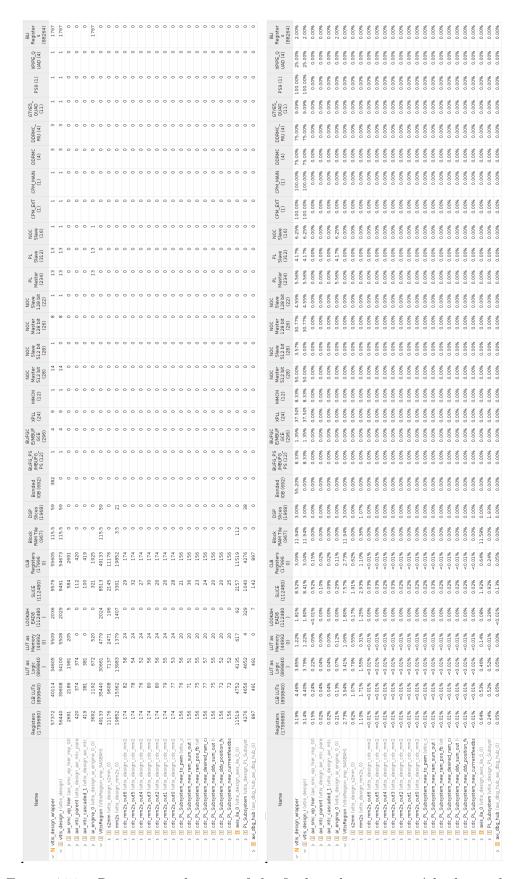


Figure 4.51. Resources utilization of the final implementation (absolute values on the left and % on the right).

It is interesting to look at the 13 PL Master and PL Slave, which are the interfaces between AIE and PL, and at the BRAM tiles, in which the ILA data is stored.

From the low % values, it is clear that the chip is largely unused, and even if it may seem that the NoC master ports are almost fully employed, many of them have a low bandwidth occupation.

## Chapter 5

## Conclusions

## 5.1 Considerations on technology and tools

It is important to highlight that the toolchain used to program this technology is still relatively immature and not yet fully refined.

During the development of this thesis, several issues and errors were encountered—many of which have already been discussed in the previous chapters. The following are some of the most significant problems that led to major delays in the development process:

- "HDDMProto::readMessage failed: bad length" This error occurred when pressing the "Analyze" button to launch an AIE simulation in Model Composer. It was resolved by uninstalling and reinstalling the tool.
- Limited documentation on AIE-ML floating-point accuracy This caused several days of delay and required the opening of a support ticket with AMD. The issue was clarified once a more detailed version of the manual was provided.
- Timing analysis in Model Composer This was not possible due to the presence of blocks external to the AMD subsystem. According to available information, the issue may be resolved in version 2025.2 of the tool.
- Simulation freeze Even on a correctly configured machine (correct tool versions, valid licenses, and all required libraries), simulations sometimes froze indefinitely without displaying error messages. Although they appeared to be running, no actual CPU usage was observed. This issue remains unresolved.
- Hardware emulation in Vivado This failed with the error: "Neither AIE\_WORK\_DIR
   (AIE work directory) is set nor AIESIM\_CONFIG. AIE simulation won't run
   without setting any one of these. Exiting simulation". No useful documentation was found, and a support ticket was opened.

While these issues highlight the current limitations of the toolchain, it is reasonable to expect that programming SoCs with these tools will become increasingly streamlined soon, as the ecosystem matures and documentation improves. However, at present, the development process can still be quite challenging and, at times, frustrating.

## 5.2 Summary and Conclusions

This research aimed to explore the capabilities of a new and versatile technology: the ACAP (Adaptive Compute Acceleration Platform).

These devices integrate, on a single chip, two scalar processors, programmable logic (PL), which can be used for data management, specific computations, or as a custom hierarchical memory, and AI Engines—VLIW (Very Long Instruction Word) SIMD (Single Instruction Multiple Data) computational units.

These AI Engines are arranged in a two-dimensional array, where each unit corresponds to one tile, located in a dedicated region of the device. Each tile operates at 1 GHz and, depending on data precision, can simultaneously perform from 8 single-precision floating-point (SPFP) MACs to 128 INT8 MACs.

The number of tiles can range from 30 to 300, and data transfer among them is managed through AXI4 protocols, a dedicated cascade mechanism, or shared memory—allowing direct communication between tiles without occupying the NoC (Network-on-Chip) bandwidth unnecessarily.

The NoC interconnects all chip components, enabling flexible partitioning of computation. In addition, a dedicated connection between the AI Engines and the PL is available.

The design process for exploiting this technology involves several steps, greatly simplified by AMD Vitis Model Composer, a high-level graphical development environment based on Simulink:

- Mapping the algorithm across the chip components (preferably using AI Engines for highly parallelizable computations);
- Programming AI Engine kernels in C/C++;
- Describing the PL hardware (Model Composer offers a library of common building blocks);
- Connecting the PL and AIE subsystems via AXI interfaces (again, simplified through dedicated Model Composer blocks);
- Synthesizing and analyzing the design to meet performance and resource constraints.

The selected case study was the Actuator Control Electronics (ACE): the actuator control loop of an FCC, composed of 34 second-order IIR filters connected through a complex data flow graph (DFG).

Unfortunately, due to data dependencies and the need for single-sample processing, this application is not ideal for a massively parallel architecture like the AIE array, which is capable of delivering several TFLOPs.

Nonetheless, its implementation is feasible—albeit underutilized—especially considering the flexibility of selectively activating only the required tiles, thus conserving power and resources.

Initial experiments also revealed that the number and order of filters processed in parallel had minimal impact on performance when using AI Engines as the primary computational units.

Moreover, multiple kernels can be deployed within a single tile, and switching from fixed-point to floating-point arithmetic may offer significant benefits in terms of dynamic

range and precision—without incurring a noticeable latency penalty.

From a performance standpoint, as shown in Chapter 4, the entire ACE can be traversed in less than  $2: \mu s$ , which is more than 50 times faster than the original requirement.

Regarding reliability, studies in [17] and [14] highlight the effectiveness of built-in mitigation techniques and the XilSEM firmware in preventing uncorrectable errors, ensuring high dependability in safety-critical applications such as avionics and automotive.

For these reasons, this technology presents a compelling option for enhancing hardware flexibility—e.g., increasing filter order or count. However, it must be noted that other applications with much higher intrinsic parallelism—such as radar signal processing or convolutional neural networks (CNNs)—stand to benefit far more from the computational acceleration offered by AI Engines.

Future research in such domains is certainly recommended.

## Appendix A

## Math demonstrations

## A.1 Two IIR order 2 filter compressed in a single IIR order 4

Supposing to have two IIR filters of order 2 in series, the overall transfer function is given by the product of the standalone transfer functions of each filter.

Given

$$F_1(z) = \frac{B_1 + B_2 \cdot z^{-1} + B_3 \cdot z^{-2}}{1 + A_2 \cdot z^{-1} + A_3 \cdot z^{-2}}$$
(A.1)

and

$$F_2(z) = \frac{C_1 + C_2 \cdot z^{-1} + C_3 \cdot z^{-2}}{1 + D_2 \cdot z^{-1} + D_3 \cdot z^{-2}}$$
(A.2)

so the overall H(z) will be:

$$H(z) = F_{1}(z) \cdot F_{2}(z) = \tag{A.3}$$

$$= \frac{B_{1} + B_{2} \cdot z^{-1} + B_{3} \cdot z^{-2}}{1 + A_{2} \cdot z^{-1} + A_{3} \cdot z^{-2}} \cdot \frac{C_{1} + C_{2} \cdot z^{-1} + C_{3} \cdot z^{-2}}{1 + D_{2} \cdot z^{-1} + D_{3} \cdot z^{-2}} = \tag{A.4}$$

$$= \frac{B_{1}C_{1} + (B_{1}C_{2} + B_{2}C_{1}) \cdot z^{-1} + (B_{1}C_{3} + B_{3}C_{1} + B_{2}C_{2}) \cdot z^{-2} + (B_{2}C_{3} + B_{3}C_{2}) \cdot z^{-3} + B_{3}C_{3} \cdot z^{-4}}{1 + (A_{2} + D_{2}) \cdot z^{-1} + (A_{3} + D_{3} + A_{2}D_{2}) \cdot z^{-2} + (A_{2}D_{3} + A_{3}D_{2}) \cdot z^{-3} + A_{3}D_{3} \cdot z^{-4}} \tag{A.5}$$

that can be considered as an IIR filter of order 4.

## Appendix B

# Model composer generated code

Codes in this Section are relative to the ACE project.

#### B.1 ACE host.cpp

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <stdint.h>
4 #include <unistd.h>
5 #include "platform.h"
6 #include "xparameters.h"
7 #include "xil io.h"
8 #include "xil_cache.h"
9 #include "PL_Subsystem_new_currentfb.h"
#include "PL_Subsystem_new_ddv_ratio.h"
11 #include "PL_Subsystem_new_ddv_sum.h"
12 #include "PL_Subsystem_new_flutter_exc.h"
13 #include "PL_Subsystem_new_ram_ratio.h"
14 #include "PL_Subsystem_new_ram_sum.h"
15 #include "PL Subsystem new setpoint.h"
16 #include "PL_Subsystem_new_desired_ram_cmd.h"
#include "PL Subsystem new ram sum out.h"
#include "PL_Subsystem_new_ram_pos_fb.h"
19 #include "PL_Subsystem_new_ddv_sum_out.h"
20 #include "PL_Subsystem_new_ddv_position_feedback.h"
21 #include "PL_Subsystem_new_currrentfeedback_amp.h"
22 #include "PL_Subsystem_new_to_pwm.h"
23 #include "xmc_ps_main.h"
24 #include "../../ip/AIE_Subsystem/src/AIE_Subsystem.cpp"
void InitData(uint32_t** out, int size)
30 int i;
```

```
*out = (uint32_t*) malloc(sizeof(uint32_t) * size);
31
32
      if(*out == NULL) {
33
         printf("Memory Allocation failed."
34
                 "Please reduce the Sample Time in your Model Composer design"
35
                 "to reduce the number of samples in theoutput(s) of the design
      .\n");
         \operatorname{exit}(-1);
37
38
39
     for(i = 0; i < size; i++) {
40
         (*out)[i] = 0xABCDEF00;
41
42
43
44
45
  int RunTest(uint64 t mm2s base,
47
               uint64_t s2mm_base,
48
                uint32_t* in0, int input0_size,
49
               uint32_t* in1, int input1_size,
               uint32_t* in2, int input2_size,
               uint32_t* in3, int input3_size,
               uint32_t* in4, int input4_size,
53
               uint32_t* in5, int input5_size,
               uint32 t* in6, int input6 size,
               uint32_t* golden0,
56
               uint32_t* golden1,
               uint32_t* golden2,
58
               uint32_t* golden3,
               uint32_t* golden4,
               uint32_t* golden5,
                uint32\_t* golden6,
62
                uint32_t* out0, int output0_size,
63
                uint32_t* out1, int output1_size,
64
                uint32_t* out2, int output2_size,
               uint32_t* out3, int output3_size,
66
               uint32_t* out4, int output4_size,
67
               uint32_t* out5, int output5_size,
68
69
               uint32_t* out6, int output6_size)
70
     int i;
71
     int errCount = 0;
72
     int totalErrCount = 0;
     uint64\_t mem\_in0Addr = (uint64\_t)in0;
74
      printf("in0=\%11x \setminus n", mem_in0Addr);
75
     uint64\_t mem\_in1Addr = (uint64\_t)in1;
76
      printf("in1=\%llx \setminus n", mem_in1Addr);
77
     uint64 t mem in2Addr = (uint64 t)in2;
78
      printf("in2=\%llx \setminus n", mem_in2Addr);
79
     uint64\_t mem\_in3Addr = (uint64\_t)in3;
80
      printf("in3=\%llx \n", mem_in3Addr);
81
     uint64_t mem_in4Addr = (uint64_t)in4;
82
      printf("in4=\%llx \n", mem_in4Addr);
83
     uint64\_t mem\_in5Addr = (uint64\_t)in5;
84
85
      printf("in5=\%llx \n", mem_in5Addr);
     uint64_t mem_in6Addr = (uint64_t)in6;
86
```

```
printf("in6=\%llx \n", mem_in6Addr);
87
      uint64\_t mem\_out0Addr = (uint64\_t)out0;
88
      printf("out0=\%llx \n", mem_out0Addr);
      uint64\_t mem\_out1Addr = (uint64\_t)out1;
90
      printf("out1=%llx\n", mem_out1Addr);
91
      uint64_t mem_out2Addr = (uint64_t)out2;
92
      printf("out2=%llx\n", mem_out2Addr);
      uint64_t mem_out3Addr = (uint64_t)out3;
94
      printf("out3=\%llx \n", mem_out3Addr);
95
      uint64\_t mem\_out4Addr = (uint64\_t)out4;
96
      printf("out4=%llx\n", mem_out4Addr);
97
      uint64 t mem out5Addr = (uint64 t)out5;
98
      printf("out5=\%llx \n", mem_out5Addr);
99
      uint64\_t mem\_out6Addr = (uint64\_t)out6;
      printf("out6=%llx\n", mem_out6Addr);
      printf("Starting test w/ cu\n");
103
104
      Xil Out32 (mm2s base + M INO OFFSET, (uint32 t) mem inOAddr);
      Xil\_Out32(mm2s\_base + M\_IN0\_OFFSET + 4, 0);
106
      Xil_Out32(mm2s_base + M_IN0_SIZE_OFFSET, input0_size);
      Xil_Out32(mm2s_base + M_IN1_OFFSET, (uint32_t) mem_in1Addr);
109
      Xil\_Out32(mm2s\_base + M\_IN1\_OFFSET + 4, 0);
      Xil_Out32(mm2s_base + M_IN1_SIZE_OFFSET, input1_size);
112
      Xil Out32 (mm2s base + M IN2 OFFSET, (uint32 t) mem in2Addr);
113
      Xil\_Out32(mm2s\_base + M\_IN2\_OFFSET + 4, 0);
114
      Xil_Out32(mm2s_base + M_IN2_SIZE_OFFSET, input2_size);
      Xil\_Out32(mm2s\_base + M\_IN3\_OFFSET, (uint32\_t) mem\_in3Addr);
      Xil\_Out32(mm2s\_base + M\_IN3\_OFFSET + 4, 0);
118
      Xil_Out32(mm2s_base + M_IN3_SIZE_OFFSET, input3_size);
119
120
      Xil Out32 (mm2s base + M IN4 OFFSET, (uint32 t) mem in4Addr);
      Xil\_Out32(mm2s\_base + M\_IN4\_OFFSET + 4, 0);
      Xil Out32 (mm2s base + M IN4 SIZE OFFSET, input4 size);
124
      Xil_Out32(mm2s_base + M_IN5_OFFSET, (uint32_t) mem_in5Addr);
      Xil\_Out32(mm2s\_base + M\_IN5\_OFFSET + 4, 0);
126
      Xil_Out32(mm2s_base + M_IN5_SIZE_OFFSET, input5_size);
127
128
      Xil_Out32(mm2s_base + M_IN6_OFFSET, (uint32_t) mem_in6Addr);
      Xil\_Out32(mm2s\_base + M\_IN6\_OFFSET + 4, 0);
130
      Xil_Out32(mm2s_base + M_IN6_SIZE_OFFSET, input6_size);
132
      Xil_Out32(s2mm_base + S_OUT0_OFFSET, (uint32_t) mem_out0Addr);
133
      Xil\_Out32(s2mm\_base + S\_OUT0\_OFFSET + 4, 0);
134
      Xil_Out32(s2mm_base + S_OUT0_SIZE_OFFSET, output0_size);
135
136
      Xil Out32(s2mm base + S OUT1 OFFSET, (uint32 t) mem out1Addr);
      Xil\_Out32(s2mm\_base + S\_OUT1\_OFFSET + 4, 0);
      Xil_Out32(s2mm_base + S_OUT1_SIZE_OFFSET, output1_size);
139
140
      Xil_Out32(s2mm_base + S_OUT2_OFFSET, (uint32_t) mem_out2Addr);
141
      Xil\_Out32(s2mm\_base + S\_OUT2\_OFFSET + 4, 0);
142
143
      Xil_Out32(s2mm_base + S_OUT2_SIZE_OFFSET, output2_size);
144
```

```
Xil_Out32(s2mm_base + S_OUT3_OFFSET, (uint32_t) mem_out3Addr);
145
      Xil\_Out32(s2mm\_base + S\_OUT3\_OFFSET + 4, 0);
146
      Xil_Out32(s2mm_base + S_OUT3_SIZE_OFFSET, output3_size);
147
148
      Xil_Out32(s2mm_base + S_OUT4_OFFSET, (uint32_t) mem_out4Addr);
149
      Xil\_Out32(s2mm\_base + S\_OUT4\_OFFSET + 4, 0);
150
      Xil_Out32(s2mm_base + S_OUT4_SIZE_OFFSET, output4_size);
151
      Xil_Out32(s2mm_base + S_OUT5_OFFSET, (uint32_t) mem_out5Addr);
153
      Xil\_Out32(s2mm\_base + S\_OUT5\_OFFSET + 4, 0);
154
      Xil_Out32(s2mm_base + S_OUT5_SIZE_OFFSET, output5_size);
155
      Xil_Out32(s2mm_base + S_OUT6_OFFSET, (uint32_t) mem_out6Addr);
157
      Xil\_Out32(s2mm\_base + S\_OUT6\_OFFSET + 4, 0);
      Xil_Out32(s2mm_base + S_OUT6_SIZE_OFFSET, output6_size);
159
      Xil_Out32(mm2s_base + M_PULSE_OFFSET, 2);
161
      Xil Out32 (mm2s base + M CTRL OFFSET, 1);
162
      Xil Out32 (s2mm base + S CTRL OFFSET, 1);
163
164
      printf("AI Engine graph init\n");
165
      mygraph.init();
167
      printf("AI Engine graph run\n");
168
      mygraph.run();
169
170
      int iteration = 0;
171
      int num out 0 samples = 0;
      int num_out1_samples = 0;
      int num_out2\_samples = 0;
174
      int num_out3\_samples = 0;
      int num\_out4\_samples = 0;
      int num_out5\_samples = 0;
177
178
      int num\_out6\_samples = 0;
      \mathbf{while}(1) {
179
         uint32\_t v =
180
             Xil In32 (s2mm base + S CTRL OFFSET);
         if (iteration \% 1000 = 0) {
182
183
             if (out0 [0] != 0xABCDEF00) {
184
                while ((num_out0_samples < output0_size) &&
185
                       (out0[num\_out0\_samples] != 0xABCDEF00)) {
186
                   int incr = (output0_size < (num_out0_samples + 50)) ? 1 :</pre>
187
      50;
                   num_out0_samples += incr;
                }
189
                printf("%d out of %d samples have been processed for
190
      new_desired_ram_cmd...\n", num_out0_samples, output0_size);
             }
191
192
             if (out1 [0] != 0xABCDEF00) {
193
                while((num_out1_samples < output1_size) &&</pre>
194
                       (out1[num\_out1\_samples] != 0xABCDEF00)) {
195
                   int incr = (output1 size < (num out1 samples + 50)) ? 1 :
196
      50;
                   num_out1_samples += incr;
197
198
```

```
printf("%d out of %d samples have been processed for
199
      new_ram_sum_out...\n", num_out1_samples, output1_size);
             }
201
             if (out2 [0] != 0xABCDEF00) {
202
                while((num_out2_samples < output2_size) &&</pre>
203
                        (out2[num\_out2\_samples] != 0xABCDEF00)) {
                    int incr = (output2 \text{ size} < (num out2 \text{ samples} + 50)) ? 1:
205
       50;
                   num_out2_samples += incr;
206
                }
207
                printf("%d out of %d samples have been processed for
208
      new_ram_pos_fb...\n", num_out2_samples, output2_size);
209
210
             if (out3 [0] != 0xABCDEF00) {
211
                while((num_out3_samples < output3_size) &&</pre>
212
                       (out3[num_out3_samples] != 0xABCDEF00)) {
213
                    int incr = (output3 size < (num out3 samples + 50)) ? 1 :
214
       50:
                    num_out3_samples += incr;
215
                printf("%d out of %d samples have been processed for
217
      new_ddv_sum_out...\n", num_out3_samples, output3_size);
             }
218
219
             if (out4 [0] != 0xABCDEF00) {
220
                while ((num_out4_samples < output4_size) &&
221
                       (out4[num_out4_samples] != 0xABCDEF00)) {
                    int incr = (output4_size < (num_out4_samples + 50)) ? 1 :</pre>
223
       50;
                    num_out4_samples += incr;
224
                }
226
                printf("%d out of %d samples have been processed for
       \frac{\text{new\_ddv\_position\_feedback}... \setminus \text{n"}}{\text{n num\_out4\_samples}}, \text{ output4\_size});
             }
227
             if (out5 [0] != 0xABCDEF00) {
229
                while ((num out5 samples < output5 size) &&
230
                       (out5[num_out5_samples] != 0xABCDEF00)) {
231
                    int incr = (output5_size < (num_out5_samples + 50)) ? 1 :</pre>
       50;
                   num_out5_samples += incr;
233
                }
234
                printf("%d out of %d samples have been processed for
       new_currrentfeedback_amp...\n", num_out5_samples, output5_size);
             }
236
237
             if (out6 [0] != 0xABCDEF00) {
                while ((num out6 samples < output6 size) &&
239
                       (out6[num out6 samples] != 0xABCDEF00)) {
240
                    int incr = (output6_size < (num_out6_samples + 50)) ? 1 :</pre>
       50;
                    num out6 samples += incr;
242
                }
243
                printf("%d out of %d samples have been processed for new_to_pwm
             , num_out6_samples, output6_size);
245
```

```
246
         ++iteration;
247
         if (v & 6) {
            break;
249
250
      }
251
253
254
      255
      ********* \ n \ n " ) ;
256
      errCount = 0;
257
      for(i = 0; i < output0_size; ++i) {
258
         if (out0[i] != golden0[i]) {
259
            printf("Mismatch found for output signal new_desired_ram_cmd in
260
      sample %d."
                    "Hardware output: %x,"
261
                   "Model Composer output: %x\n",
262
                   i, out0[i], golden0[i]);
263
            ++errCount;
264
         }
266
      if (errCount == 0)  {
267
         printf("***** Model Composer and Hardware outputs match for all %d
268
      samples
                "for output signal new desired ram cmd\n", output0 size);
269
      \} else \{
         printf("***** Mismatch(es) found between Model Composer and Hardware
271
       outputs
                 "for output signal new_desired_ram_cmd\n");
272
273
      totalErrCount += errCount;
274
275
      errCount = 0;
276
      for(i = 0; i < output1\_size; ++i) {
277
         if (out1[i] != golden1[i]) {
            printf("Mismatch found for output signal new_ram_sum_out in sample
279
       %d."
                    "Hardware output: %x,"
280
                   "Model Composer output: %x\n",
281
                   i, out1[i], golden1[i]);
282
            ++errCount;
283
         }
284
      if (errCount == 0)  {
286
         printf("***** Model Composer and Hardware outputs match for all %d
287
      samples
                "for output signal new_ram_sum_out\n", output1_size);
      } else {
289
         printf("***** Mismatch(es) found between Model Composer and Hardware
290
       outputs
                 "for output signal new_ram_sum_out\n");
291
292
      totalErrCount += errCount;
293
294
295
      errCount = 0;
      for(i = 0; i < output2\_size; ++i) {
296
```

```
if (out2[i] != golden2[i]) {
297
             printf("Mismatch found for output signal new_ram_pos_fb in sample
298
      %d. "
                    "Hardware output: %x,"
299
                    "Model Composer output:%x\n",
300
                    i, out2[i], golden2[i]);
301
            ++errCount;
302
         }
303
304
      if (errCount == 0)  {
305
         printf("***** Model Composer and Hardware outputs match for all %d
306
      samples
                 "for output signal new_ram_pos_fb\n", output2_size);
307
      \} else \{
308
         printf("***** Mismatch(es) found between Model Composer and Hardware
309
       outputs
                 "for output signal new_ram_pos_fb\n");
310
311
      totalErrCount += errCount;
312
313
      errCount = 0;
314
      for(i = 0; i < output3_size; ++i) {
         if (out3 [i] != golden3 [i]) {
316
             printf("Mismatch found for output signal new_ddv_sum_out in sample
317
       %d."
                    "Hardware output: %x,"
318
                    "Model Composer output: %x\n",
319
                    i, out3[i], golden3[i]);
320
            ++errCount;
321
         }
322
      if (errCount == 0)  {
324
         printf("***** Model Composer and Hardware outputs match for all %d
325
      samples
                 "for output signal new ddv sum out\n", output3 size);
      } else {}
327
         printf("***** Mismatch(es) found between Model Composer and Hardware
       outputs
                 "for output signal new ddv sum out\n");
329
330
      totalErrCount += errCount;
331
332
      errCount = 0;
333
      for(i = 0; i < output4\_size; ++i) {
334
         if(out4[i] != golden4[i]) {
             printf("Mismatch found for output signal new_ddv_position_feedback
336
       in sample %d.'
                    "Hardware output: %x,"
337
                    "Model Composer output:%x\n",
338
                    i, out4[i], golden4[i]);
339
            ++errCount;
340
         }
341
      if (errCount == 0)  {
343
         printf("***** Model Composer and Hardware outputs match for all %d
344
      samples
                 "for output signal new_ddv_position_feedback\n", output4_size)
345
```

```
\} else \{
346
         printf("***** Mismatch(es) found between Model Composer and Hardware
347
                 "for output signal new_ddv_position_feedback\n");
348
349
      totalErrCount += errCount;
350
351
      errCount = 0;
352
      for(i = 0; i < output5_size; ++i) {
353
         if (out5[i] != golden5[i]) {
354
             printf("Mismatch found for output signal new_currrentfeedback_amp
355
      in sample %d.
                    "Hardware output: \%x,"
356
                    "Model Composer output:%x \ n",
357
                    i, out5[i], golden5[i]);
358
            ++errCount;
359
         }
360
      if (errCount == 0)  {
362
         printf("***** Model Composer and Hardware outputs match for all %d
363
      samples
                 "for output signal new_currrentfeedback_amp\n", output5_size);
      \} else \{
365
         printf("***** Mismatch(es) found between Model Composer and Hardware
366
       outputs
                 "for output signal new_currrentfeedback_amp\n");
368
      totalErrCount += errCount;
369
371
      errCount = 0;
      for(i = 0; i < output6\_size; ++i) {
372
         if (out6[i] != golden6[i]) {
373
             printf("Mismatch found for output signal new_to_pwm in sample %d."
374
                     "Hardware output: %x,'
375
                    "Model Composer output:%x\n",
376
                    i, out6[i], golden6[i]);
377
            ++errCount;
         }
379
380
      if (errCount == 0)  {
381
         printf("***** Model Composer and Hardware outputs match for all %d
      samples
                 "for output signal new_to_pwm\n", output6_size);
383
      } else {}
384
         printf("***** Mismatch(es) found between Model Composer and Hardware
       outputs
                 "for output signal new_to_pwm\n");
386
387
      totalErrCount += errCount;
388
389
      return totalErrCount;
390
391
392
393
      */
  int main()
394
   uint32 t* PL Subsystem new desired ram cmd;
```

```
uint32_t* PL_Subsystem_new_ram_sum_out;
397
      uint32_t * PL_Subsystem_new_ram_pos_fb;
398
      uint32_t* PL_Subsystem_new_ddv_sum_out;
399
      uint32_t* PL_Subsystem_new_ddv_position_feedback;
400
      uint32_t* PL_Subsystem_new_currrentfeedback_amp;
401
      uint32_t* PL_Subsystem_new_to_pwm;
402
      int totalErrCount;
403
404
      Xil_DCacheDisable();
405
406
      init_platform();
407
      sleep (1);
408
      printf("Beginning test\n");
409
      InitData(&PL_Subsystem_new_desired_ram_cmd, OUTPUT_0_SIZE);
411
      InitData(&PL_Subsystem_new_ram_sum_out, OUTPUT_1_SIZE);
412
      InitData(&PL_Subsystem_new_ram_pos_fb , OUTPUT_2_SIZE) ;
413
      InitData(&PL_Subsystem_new_ddv_sum_out, OUTPUT_3_SIZE);
414
      InitData(&PL Subsystem new ddv position feedback, OUTPUT 4 SIZE);
415
      InitData(&PL Subsystem new currrentfeedback amp, OUTPUT 5 SIZE);
416
      InitData(&PL_Subsystem_new_to_pwm, OUTPUT_6_SIZE);
417
      totalErrCount = RunTest(MM2S_BASE, S2MM_BASE,
419
                          PL_Subsystem_new_currentfb_indata, INPUT_0_SIZE,
420
                          PL_Subsystem_new_ddv_ratio_indata, INPUT_1_SIZE,
421
                          PL_Subsystem_new_ddv_sum_indata, INPUT_2_SIZE,
                          PL Subsystem new flutter exc indata, INPUT 3 SIZE,
423
                          PL_Subsystem_new_ram_ratio_indata, INPUT_4_SIZE,
424
                          PL_Subsystem_new_ram_sum_indata, INPUT_5_SIZE,
                          PL_Subsystem_new_setpoint_indata, INPUT_6_SIZE,
                          PL\_Subsystem\_new\_desired\_ram\_cmd\_goldendata\;,
427
                          PL_Subsystem_new_ram_sum_out_goldendata,
428
                          PL_Subsystem_new_ram_pos_fb_goldendata,
429
430
                          PL_Subsystem_new_ddv_sum_out_goldendata,
                          PL_Subsystem_new_ddv_position_feedback_goldendata,
431
                          PL_Subsystem_new_currrentfeedback_amp_goldendata,
432
                          PL Subsystem new to pwm goldendata,
                          PL_Subsystem_new_desired_ram_cmd, OUTPUT_0_SIZE,
434
                          PL_Subsystem_new_ram_sum_out, OUTPUT_1_SIZE,
435
                          PL_Subsystem_new_ram_pos_fb , OUTPUT_2_SIZE,
436
                          PL_Subsystem_new_ddv_sum_out, OUTPUT_3_SIZE.
437
                          PL_Subsystem_new_ddv_position_feedback, OUTPUT_4_SIZE
438
                          PL_Subsystem_new_currrentfeedback_amp, OUTPUT_5_SIZE,
439
                          PL_Subsystem_new_to_pwm, OUTPUT_6_SIZE);
441
      if (totalErrCount == 0) {
442
443
         printf("\n****** TEST PASSED\n");
445
         printf("ERROR: TEST FAILED! Error count: %d \n",
446
                  totalErrCount);
448
449
      printf("\n******* VMC_TEST_DONE\n");
450
      cleanup_platform();
451
452
      return totalErrCount;
453
```

```
454 }
```

This code runs on the processing system. The main function calls the RunTest function, which initializes and starts the DMAs and then initializes and starts the AIE graph. If a mismatch in the hardware output is detected, an error message is printed on the serial terminal; if not, a TEST PASSED message is displayed.

### B.2 ACE AIE\_subsystem.h

```
XMC AIE SUBSYSTEM H
<sup>2</sup> #define __XMC_AIE_SUBSYSTEM_H
4 #include <adf.h>
5 #include "kernels/IIR_30.hpp"
6 #include "kernels/IIR_2.hpp'
7 #include "kernels/IIR_1.hpp"
8 #include "kernels/IIR_input_block_1.hpp"
9 #include "kernels/IIR_input_block_2.hpp"
#include "kernels/IIR_input_block_3.hpp"
11 #include "kernels/IIR_CONTROLLER_1.hpp"
12 #include "kernels/IIR_CONTROLLER_2.hpp"
13 #include "kernels/IIR_CONTROLLER_3.hpp"
14 #include "kernels/IIR_25.hpp"
15 #include "kernels/IIR_20.hpp"
16 #include "kernels/IIR_19.hpp"
17 #include "kernels/IIR_18.hpp"
  class AIE_Subsystem_base : public adf::graph {
19
  public:
20
     adf::kernel IIR 30 0;
21
     adf::kernel IIR 2 0;
22
     adf::kernel IIR_1_0;
23
     adf::kernel IIR_input_block_1_0;
     adf::kernel IIR_input_block_2_0;
     adf::kernel IIR_input_block_3_0;
26
     adf::kernel IIR_controller_1;
27
     adf::kernel IIR_controller_2;
28
     adf::kernel IIR_controller_3;
     adf::kernel IIR_25_0;
30
     adf::kernel IIR 20 0;
31
     adf::kernel IIR_19_0;
     adf::kernel IIR_18_0;
34
  public:
35
     adf::input_port TF1_in, TF2_in, TF19_in, TF20_in, input_block_1_in,
     controller_1_in, TF25_in, input_block_2_in, controller_2_in, TF18_in,
     TF30_in, input_block_3_in, controller_3_in;
     adf::output\_port\ TF1\_out,\ TF2\_out,\ TF19\_out,\ TF\_20\_out,
     input_block_1_out, controller_1_out, TF_25_out, input_block_2_out,
      controller\_2\_out\;,\;\; TF18\_out\;,\;\; TF30\_out\;,\;\; input\_block\_3\_out\;,
      controller_3_out;
38
     AIE_Subsystem_base() {
39
         // create kernel IIR_30_0
40
        IIR_30_0 = adf::kernel::create(IIR_30);
41
         adf::source(IIR\_30\_0) = "kernels/IIR\_30.cpp";
42
```

```
43
         // create kernel IIR_2_0
44
         IIR_2_0 = adf :: kernel :: create(IIR_2);
         adf::source(IIR_2_0) = "kernels/IIR_2.cpp";
46
47
         // create kernel IIR_1_0
         IIR_1_0 = adf::kernel::create(IIR_1);
         adf::source(IIR_1_0) = "kernels/IIR_1.cpp";
         // create kernel IIR_input_block_1_0
         IIR_input_block_1_0 = adf::kernel::create(IIR_input_block_1);
         adf::source(IIR_input_block_1_0) = "kernels/IIR_input_block_1.cpp";
54
         // create kernel IIR_input_block_2_0
         IIR_input_block_2_0 = adf::kernel::create(IIR_input_block_2);
         adf::source(IIR_input_block_2_0) = "kernels/IIR_input_block_2.cpp";
58
59
         // create kernel IIR_input_block_3_0
         IIR input block 3 0 = adf::kernel::create(IIR input block 3);
61
         adf::source(IIR_input_block_3_0) = "kernels/IIR_input_block_3.cpp";
         // create kernel IIR_controller_1
         IIR_controller_1 = adf::kernel::create(IIR_CONTROLLER_1);
65
         adf::source(IIR_controller_1) = "kernels/IIR_CONTROLLER_1.cpp";
67
         // create kernel IIR_controller_2
         IIR controller 2 = adf::kernel::create(IIR CONTROLLER 2);
         adf::source(IIR_controller_2) = "kernels/IIR_CONTROLLER_2.cpp";
70
         // create kernel IIR_controller_3
         IIR controller 3 = adf::kernel::create(IIR CONTROLLER 3);
73
         adf::source(IIR_controller_3) = "kernels/IIR_CONTROLLER_3.cpp";
74
75
76
         // create kernel IIR_25_0
         IIR_25_0 = adf::kernel::create(IIR_25);
77
         adf::source(IIR_25_0) = "kernels/IIR_25.cpp";
         // create kernel IIR 20 0
80
         IIR_20_0 = adf::kernel::create(IIR_20);
81
         adf::source(IIR_20_0) = "kernels/IIR_20.cpp";
82
         // create kernel IIR_19_0
84
         IIR_19_0 = adf::kernel::create(IIR_19);
85
         adf::source(IIR_19_0) = "kernels/IIR_19.cpp";
         // create kernel IIR_18_0
88
         IIR_18_0 = adf::kernel::create(IIR_18);
89
         adf::source(IIR_18_0) = "kernels/IIR_18.cpp";
90
         // create kernel constraints IIR 30 0
92
         adf::runtime < ratio > (IIR_30_0) = 0.9;
93
         // create kernel constraints IIR_2_0
         adf::runtime < ratio > (IIR_2_0) = 0.9;
96
97
         // create kernel constraints IIR_1_0
98
99
         adf::runtime < ratio > (IIR_1_0) = 0.9;
100
```

```
// create kernel constraints IIR_input_block_1_0
         adf::runtime<ratio>(IIR_input_block_1_0) = 0.9;
102
         // create kernel constraints IIR_input_block_2_0
104
         adf::runtime<ratio>(IIR_input_block_2_0) = 0.9;
105
106
         // create kernel constraints IIR_input_block_3_0
         adf::runtime<ratio>(IIR_input_block_3_0) = 0.9;
         // create kernel constraints IIR_controller_1
         adf::runtime<ratio>(IIR_controller_1) = 0.9;
111
112
         // create kernel constraints IIR_controller_2
113
         adf::runtime<ratio>(IIR_controller_2) = 0.9;
         // create kernel constraints IIR_controller_3
         adf::runtime<ratio>(IIR_controller_3) = 0.9;
117
118
         // create kernel constraints IIR 25 0
119
         adf::runtime < ratio > (IIR_25_0) = 0.9;
         // create kernel constraints IIR_20_0
         adf::runtime < ratio > (IIR_20_0) = 0.9;
123
124
         // create kernel constraints IIR_19_0
         adf::runtime < ratio > (IIR_19_0) = 0.9;
127
         // create kernel constraints IIR 18 0
         adf::runtime < ratio > (IIR_18_0) = 0.9;
130
         // create nets to specify connections
         adf::connect < adf::stream > net0 (TF1_in, IIR_1_0.in[0]);
         adf::connect < adf::stream > net1 \ (TF2_in, IIR_2_0.in |0|);
133
134
         adf::connect < adf::stream > net2 (TF19_in, IIR_19_0.in |0|);
         adf::connect < adf::stream > net3 (TF20_in, IIR_20_0.in[0]);
         adf::connect< adf::stream > net4 (input_block_1_in,
136
      IIR_input_block_1_0.in [0]);
         adf::connect< adf::stream > net5 (controller_1_in, IIR_controller_1.
137
      in [0]);
         adf::connect < adf::stream > net6 (TF25_in, IIR_25_0.in[0]);
138
         adf::connect< adf::stream > net7 (input_block_2_in,
139
      IIR\_input\_block\_2\_0.in[0];
         adf::connect< adf::stream > net8 (controller_2_in, IIR_controller_2.
140
      in [0]);
         adf::connect < adf::stream > net9 (TF18_in, IIR_18_0.in [0]);
         adf::connect < adf::stream > net10 (TF30_in, IIR_30_0.in[0]);
142
         adf::connect< adf::stream > net11 (input_block_3_in,
143
      IIR\_input\_block\_3\_0.in[0];
         adf::connect< adf::stream > net12 (controller_3_in, IIR_controller_3.
         adf::connect< adf::stream > net13 (IIR_30_0.out[0], TF30_out);
145
         adf::connect < adf::stream > net14 (IIR_2_0.out[0],
                                                               TF2\_out);
                                             (IIR_1_0.out [0], TF1_out);
         adf::connect< adf::stream > net15
147
         adf::connect < adf::stream > net16 (IIR_input_block_1_0.out|0|,
148
      input_block_1_out);
         adf::connect< adf::stream > net17 (IIR_input_block_2_0.out[0],
149
      input_block_2_out);
```

```
adf::connect< adf::stream > net18 (IIR_input_block_3_0.out[0],
150
      input_block_3_out);
         adf::connect< adf::stream > net19 (IIR_controller_1.out[0])
      controller_1_out);
         adf::connect< adf::stream > net20 (IIR_controller_2.out[0],
152
      controller_2_out);
         adf::connect< adf::stream > net21 (IIR_controller_3.out[0],
      controller 3 out);
         adf::connect < adf::stream > net22 (IIR_25_0.out[0], TF_25_out);
154
         adf::connect < adf::stream > net23 (IIR_20_0.out[0], TF_20_out);
         adf::connect < adf::stream > net24 (IIR_19_0.out[0], TF19_out);
         adf::connect< adf::stream > net25 (IIR_18_0.out[0], TF18_out);
157
      }
158
159
160
  class AIE Subsystem : public adf::graph {
161
  public:
162
      AIE_Subsystem_base mygraph;
163
164
  public:
165
      adf::input_plio TF1_in, TF2_in, TF19_in, TF20_in, input_block_1_in,
      controller_1_in, TF25_in, input_block_2_in, controller_2_in, TF18_in,
      TF30_in, input_block_3_in, controller_3_in;
      adf::output_plio TF1_out, TF2_out, TF19_out, TF_20_out,
167
      input_block_1_out, controller_1_out, TF_25_out, input_block_2_out,
      controller_2_out, TF18_out, TF30_out, input_block_3_out,
      controller 3 out;
168
      AIE_Subsystem() {
         TF1_in = adf::input_plio::create("TF1_in",
170
               adf::plio_64_bits,
                './data/input/TF1_in.txt");
173
174
         TF2_in = adf::input_plio::create("TF2_in",
               adf::plio_64_bits,
               "./data/input/TF2_in.txt");
         TF19 in = adf::input plio::create("TF19 in",
178
               adf::plio_64_bits,
179
                "./data/input/TF19_in.txt");
180
181
         TF20_in = adf::input_plio::create("TF20_in",
182
               adf::plio_64_bits,
183
                " ./data/input/TF20_in.txt " );
         input_block_1_in = adf::input_plio::create("input_block_1_in",
186
               adf::plio_64_bits,
187
                "./data/input/input_block_1_in.txt");
188
189
         controller 1 in = adf::input plio::create("controller 1 in",
190
               adf::plio_64_bits,
191
                "./data/input/controller_1_in.txt");
         TF25_in = adf::input_plio::create("TF25_in",
194
               adf::plio_64_bits,
195
                "./data/input/TF25_in.txt");
196
197
         input_block_2_in = adf::input_plio::create("input_block_2_in",
198
```

```
adf::plio_64_bits,
199
                "./data/input/input_block_2_in.txt");
200
         controller_2_in = adf::input_plio::create("controller_2_in",
202
                adf::plio_64_bits,
203
                "./data/input/controller_2_in.txt");
204
         TF18_in = adf::input_plio::create("TF18_in",
206
                adf::plio_64_bits,
207
                "./data/input/TF18_in.txt");
208
209
         TF30_in = adf::input_plio::create("TF30_in",
210
                adf::plio_64_bits,
211
                "./data/input/TF30_in.txt");
212
213
         input_block_3_in = adf::input_plio::create("input_block_3_in",
214
                adf::plio_64_bits,
215
                "./data/input/input_block_3_in.txt");
216
217
         controller_3_in = adf::input_plio::create("controller_3_in",
218
                adf::plio_64_bits,
219
                "./data/input/controller_3_in.txt");
221
         TF1_out = adf::output_plio::create("TF1_out",
222
                adf::plio_64_bits,
223
                "TF1_out.txt");
224
225
         TF2_out = adf::output_plio::create("TF2_out",
226
                adf::plio_64_bits,
                "TF2_out.txt");
         TF19_out = adf::output_plio::create("TF19_out",
230
                adf::plio_64_bits,
231
232
                "TF19_out.txt");
233
         TF_20_out = adf::output_plio::create("TF_20_out",
234
                adf::plio_64_bits,
                "TF_20_out.txt");
236
         input_block_1_out = adf::output_plio::create("input_block_1_out",
238
                adf::plio_64_bits,
239
                "input_block_1_out.txt");
240
241
         controller_1_out = adf::output_plio::create("controller_1_out",
242
                adf::plio_64_bits,
                "controller_1_out.txt");
244
245
         TF_25_out = adf::output_plio::create("TF_25_out",
246
                adf::plio_64_bits,
                "TF 25 out.txt");
248
         input_block_2_out = adf::output_plio::create("input_block_2_out",
                adf::plio_64_bits,
                "input_block_2_out.txt");
252
253
         controller_2_out = adf::output_plio::create("controller_2_out",
254
255
                adf::plio_64_bits,
                "controller_2_out.txt");
256
```

```
257
         TF18_out = adf::output_plio::create("TF18_out",
258
               adf::plio_64_bits,
                "TF18_out.txt");
260
261
         TF30_out = adf::output_plio::create("TF30_out",
               adf::plio_64_bits,
                "TF30 out.txt");
264
265
         input_block_3_out = adf::output_plio::create("input_block_3_out",
               adf::plio_64_bits,
267
                "input block 3 out.txt");
268
269
         controller_3_out = adf::output_plio::create("controller_3_out",
               adf::plio_64_bits,
271
                "controller_3_out.txt");
272
273
         adf::connect <> (TF1_in.out [0], mygraph.TF1_in);
         adf::connect < > (TF2 in.out [0], mygraph.TF2 in);
         adf::connect <> (TF19_in.out[0], mygraph.TF19_in);
                          \left(\,TF20\_in\,.\,out\,[\,0\,]\,\,,\,\,\,mygraph\,.\,TF20\_in\,\right)\,;
         adf::connect<>
                          (input_block_1_in.out [0], mygraph.input_block_1_in);
         adf::connect<>
                          (controller_1_{in.out}[0], mygraph.controller_1_{in});
         adf::connect<>
279
         adf::connect<>
                          (TF25_{in.out}[0], mygraph.TF25_{in});
280
                          (input_block_2_in.out[0], mygraph.input_block_2_in);
         adf::connect<>
281
         adf::connect<>
                          (controller_2_in.out[0], mygraph.controller_2_in);
                         (TF18_in.out [0], mygraph.TF18_in);
         adf::connect<>
283
         adf::connect<>
                          (TF30_in.out [0], mygraph.TF30_in);
                          (input_block_3_in.out[0], mygraph.input_block_3_in);
         adf::connect<>
         adf::connect<>
                          (controller_3_in.out[0], mygraph.controller_3_in);
         adf::connect<>
                          (mygraph.TF1\_out, TF1\_out.in[0]);
287
         adf::connect<>
                          (mygraph.TF2\_out, TF2\_out.in[0]);
288
                          (mygraph.TF19_out, TF19_out.in [0]);
         adf::connect<>
289
290
         adf::connect<>
                          (mygraph.TF_20_out, TF_20_out.in [0]);
         adf::connect<>
                          (mygraph.input_block_1_out, input_block_1_out.in[0]);
291
                          (mygraph.controller_1_out, controller_1_out.in[0]);
         adf::connect<>
292
                          (mygraph.TF_25_out, TF_25_out.in[0]);
         adf::connect<>
         adf::connect<>
                          (mygraph.input block 2 out, input block 2 out.in[0]);
294
         adf::connect<>
                          (mygraph.controller_2_out, controller_2_out.in[0]);
295
         adf::connect<>
                          (mygraph.TF18\_out, TF18\_out.in[0]);
296
         adf::connect<>
                          (mygraph.TF30_out, TF30_out.in[0]);
297
         adf::connect < > (mygraph.input_block_3_out, input_block_3_out.in[0]);
         adf::connect<> (mygraph.controller_3_out, controller_3_out.in[0]);
299
300
301
303 #endif // XMC_AIE_SUBSYSTEM_H
```

This source code contains the definition and application of constraints of the AIE graph. All the inputs and outputs of the graph are also connected to the specific .txt file containing input data and golden output, so this can be used to simulate the AIE array.

## B.3 ACE AIE\_subsystem.cpp

```
<sup>1</sup> #include "AIE_Subsystem.h"
```

```
// instantiate ADF graph
AIE_Subsystem mygraph;

// initialize and run the dataflow graph
#if defined(__AIESIM__) || defined(__X86SIM__)

int main(void) {
    mygraph.init();
    mygraph.run();
    mygraph.end();
    return 0;
}

#endif
```

This code initializes, runs and ends the AIE graph execution.

#### B.4 ACE s2mm.cpp

```
1 #include <hls_stream.h>
2 #include <ap_int.h>
3 #include <ap_axi_sdata.h>
4 #include <stdint.h>
     template <int BITWIDTH>
8
     void Push(hls::stream<ap_axiu<BITWIDTH, 0, 0, 0> >& s, uint32_t* target,
9
      int size){
      constexpr int dwidth = ((BITWIDTH/32) + (((BITWIDTH\%32)==0)? 0: 1))
      for (int i = 0; i < size; i+=dwidth/32) {
 #pragma HLS PIPELINE II=1
        ap_axiu < BITWIDTH, 0, 0, 0 > v = s.read();
14
        target[i]
                 = uint32_t(v.data);
        if constexpr (BITWIDTH == 64)
          target[i+1] = uint32_t(v.data >> 32);
        else if constexpr (BITWIDTH == 128) {
18
          target[i+1] = uint32_t(v.data >> 32);
19
          target[i+2] = uint32_t(v.data >> 64);
20
          target[i+3] = uint32\_t(v.data >> 96);
        \} else if constexpr (BITWIDTH > 32) {
          23
            target[i+j] = uint32\_t(v.data >> shift);
24
        }
26
27
29
30
    void PushAll(
31
      hls::stream<ap_axiu<32, 0, 0, 0>>& in0, uint32_t* out0, uint32_t size0
32
      hls::stream<ap_axiu<32, 0, 0, 0>>& in1, uint32_t* out1, uint32_t size1
33
```

```
{\tt hls::stream} < {\tt ap\_axiu} < 32, \ 0, \ 0, \ 0>> \& \ {\tt in2} \ , \ {\tt uint32\_t*} \ {\tt out2} \ , \ {\tt uint32\_t} \ {\tt size2}
      hls::stream<ap_axiu<32, 0, 0, 0 > \infty in 3, uint 32_t * out 3, uint 32_t size 3
      hls::stream<ap_axiu<32, 0, 0, 0> >& in4, uint32_t* out4, uint32_t size4
36
      hls::stream<ap_axiu<32, 0, 0, 0>>& in5, uint32_t* out5, uint32_t size5
      hls::stream<ap_axiu<32, 0, 0, 0 >  in6, uint32_t* out6, uint32_t size6
38
      ) {
  #pragma HLS DATAFLOW
      Push < 32 > (in0, out0, size0);
40
      Push < 32 > (in1, out1, size1);
41
      Push < 32 > (in2, out2, size2);
      Push < 32 > (in3, out3, size3);
43
      Push < 32 > (in4, out4, size4);
44
      Push < 32 > (in5, out5, size5);
45
      Push < 32 > (in6, out6, size6);
46
47
48
   extern "C" {
49
    void s2mm(
      hls::stream<ap_axiu<32, 0, 0, 0> \& in0,
      uint32_t * out0, uint32_t size0,
      hls::stream < ap_axiu < 32, 0, 0, 0 > & in1,
      uint32_t * out1, uint32_t size1,
      hls::stream<ap axiu<32, 0, 0, 0 >  in2,
      uint32_t * out2, uint32_t size2,
56
      hls::stream<ap_axiu<32, 0, 0, 0>>& in3,
      uint32_t * out3, uint32_t size3,
      hls::stream < ap\_axiu < 32, 0, 0, 0 > & in4,
      uint32\_t* out4, uint32\_t size4,
      hls::stream < ap\_axiu < 32, 0, 0, 0 > \& in5,
61
      uint32\_t* out5, uint32\_t size5,
      hls::stream < ap\_axiu < 32, 0, 0, 0 > \& in6,
63
      uint32_t* out6, uint32_t size6) {
65 #pragma HLS INTERFACE ap_ctrl_hs port=return bundle=control
66 #pragma HLS INTERFACE s axilite port=return bundle=control
67 #pragma HLS INTERFACE axis port=in0
68 #pragma HLS INTERFACE m_axi port=out0 bundle=gmem0 max_widen_bitwidth=4
69 #pragma HLS INTERFACE s_axilite port=out0 bundle=control
70 #pragma HLS INTERFACE s_axilite port=size0 bundle=control
71 #pragma HLS INTERFACE axis port=in1
72 #pragma HLS INTERFACE m_axi port=out1 bundle=gmem1 max_widen_bitwidth=4
73 #pragma HLS INTERFACE s_axilite port=out1 bundle=control
74 #pragma HLS INTERFACE s_axilite port=size1 bundle=control
75 #pragma HLS INTERFACE axis port=in2
76 #pragma HLS INTERFACE m_axi port=out2 bundle=gmem2 max_widen_bitwidth=4
77 #pragma HLS INTERFACE s_axilite port=out2 bundle=control
78 #pragma HLS INTERFACE s axilite port=size2 bundle=control
79 #pragma HLS INTERFACE axis port=in3
80 #pragma HLS INTERFACE m_axi port=out3 bundle=gmem3 max_widen_bitwidth=4
81 #pragma HLS INTERFACE s_axilite port=out3 bundle=control
82 #pragma HLS INTERFACE s_axilite port=size3 bundle=control
83 #pragma HLS INTERFACE axis port=in4
84 #pragma HLS INTERFACE m_axi port=out4 bundle=gmem4 max_widen_bitwidth=4
85 #pragma HLS INTERFACE s_axilite port=out4 bundle=control
86 #pragma HLS INTERFACE s_axilite port=size4 bundle=control
```

```
87 #pragma HLS INTERFACE axis port=in5
88 #pragma HLS INTERFACE m_axi port=out5 bundle=gmem5 max_widen_bitwidth=4
89 #pragma HLS INTERFACE s_axilite port=out5 bundle=control
90 #pragma HLS INTERFACE s_axilite port=size5 bundle=control
91 #pragma HLS INTERFACE axis port=in6
92 #pragma HLS INTERFACE m_axi port=out6 bundle=gmem6 max_widen_bitwidth=4
93 #pragma HLS INTERFACE s_axilite port=out6 bundle=control
94 #pragma HLS INTERFACE s_axilite port=size6 bundle=control
       PushAll(
96
         in0, out0, size0,
97
         in1, out1, size1,
98
         in2, out2, size2,
99
         in3, out3, size3,
100
         in4, out4, size4,
         in5, out5, size5.
         in6, out6, size6);
103
104
106
```

This code contains an HLS description of a stream-to-memory mapped DMA. The s2mm kernel starts when triggered via AXI-Lite (ap\_start). It reads data from AXI-Stream inputs (in0 to in6) and writes to memory via AXI-MM. Each stream is processed in parallel using HLS DATAFLOW, storing sizeX words to outX. DMA-like behavior is implemented in hardware; activation depends on the host writing control registers.

#### B.5 ACE mm2s.cpp

```
#include <hls stream.h>
2 #include <ap int.h>
3 #include <ap axi sdata.h>
4 #include <stdint.h>
  template <int BITWIDTH>
           void Push (uint32_t* data, hls::stream <ap_axiu <BITWIDTH, 0, 0, 0 > >&
       target, uint32_t size, int pulse){
      constexpr int dwidth = ((BITWIDTH/32) + (((BITWIDTH\%32) == 0)? 0: 1))
      *32;
      for (int i = 0; i < size * pulse; i+=dwidth/32) {
 #pragma HLS PIPELINE II=1
           int index = i % size;
13
           ap_axiu < BITWIDTH, 0, 0, 0 > v;
14
           v.last = 0;
           v.keep = 0xFFFF;
16
           v.data = data[index];
           if constexpr (BITWIDTH == 64)
               v. data += (uint64 t(data[index+1]) << 32);
19
           else if constexpr (BITWIDTH == 128)
20
               v.data += (ap_int < 128 > (data[index + 1]) < < 32) +
                          (ap_{int} < 128 > (data[index + 2]) < < 64) +
                          (ap_int < 128 > (data[index + 3]) < < 96);
           else if constexpr (BITWIDTH > 32) {
24
               for (int shift = 32, j=1; shift < BITWIDTH; shift +=32, j++) {
```

```
v.data += (ap_int <BITWIDTH>(data[index+j]) << shift);
26
               }
27
           target.write(v);
29
30
31
33
    void PushAll(
34
       uint32_t*in0, hls::stream < ap_axiu < 32, 0, 0, 0 > & out0, uint32_t size0
35
      uint32 t*in1, hls::stream<ap axiu<32, 0, 0, 0> \geq out1, uint32 t size1
36
      uint32_t*in2, hls::stream < ap_axiu < 32, 0, 0, 0 > & out2, uint32_t size2
      uint32 t* in3, hls::stream<ap axiu<32, 0, 0, 0> >& out3, uint32 t size3
38
      uint32 + in4, hls::stream<ap axiu<32, 0, 0, 0> \gg out4, uint32 + size4
      uint32 t* in5, hls::stream<ap axiu<32, 0, 0, 0> >& out5, uint32 t size5
40
      uint32_t*in6, hls::stream < ap_axiu < 32, 0, 0, 0 > & out6, uint32_t size6
      int pulse) {
42
  #pragma HLS DATAFLOW
43
      Push < 32 > (in0, out0, size0, pulse);
      Push < 32 > (in1, out1, size1, pulse);
45
      Push < 32 > (in2, out2, size2, pulse);
46
      Push < 32 > (in3, out3, size3, pulse);
      Push < 32 > (in4, out4, size4, pulse);
      Push < 32 > (in5, out5, size5, pulse);
49
      Push < 32 > (in6, out6, size6, pulse);
50
51
  extern "C" {
    void mm2s(
54
       uint32 t* in0,
      hls::stream<ap axiu<32, 0, 0, 0 >  out0, uint32 t size0,
56
      uint32 t* in1,
      hls::stream<ap_axiu<32, 0, 0, 0> & out1, uint32_t size1,
58
       uint32_t*in2,
59
      hls::stream<ap_axiu<32, 0, 0, 0>>& out2, uint32_t size2,
       uint32 t* in3,
61
      \label{eq:hls::stream} $$\operatorname{ap\_axiu} < 32, \ 0, \ 0, \ 0> >  \ \text{out3}, \ \operatorname{uint32\_t} \ \operatorname{size3}, 
       uint32_t*in4,
      hls::stream < ap\_axiu < 32, 0, 0, 0 > \& out4, uint32\_t size4,
64
      uint32_t*in5,
65
      hls::stream < ap\_axiu < 32, 0, 0, 0 > \& out5, uint32\_t size5,
      uint32 t* in6,
      hls::stream<ap axiu<32, 0, 0, 0 > \infty out6, uint32 t size6,
       int pulse) {
70 #pragma HLS INTERFACE ap_ctrl_hs port=return bundle=control
71 #pragma HLS INTERFACE s_axilite port=return bundle=control
72 #pragma HLS INTERFACE m_axi port=in0 bundle=gmem0 max_widen_bitwidth=4
73 #pragma HLS INTERFACE s_axilite port=in0 bundle=control
74 #pragma HLS INTERFACE axis port=out0
75 #pragma HLS INTERFACE s_axilite port=size0 bundle=control
76 #pragma HLS INTERFACE m axi port=in1 bundle=gmem1 max widen bitwidth=4
```

```
77 #pragma HLS INTERFACE s_axilite port=in1 bundle=control
78 #pragma HLS INTERFACE axis port=out1
79 #pragma HLS INTERFACE s_axilite port=size1 bundle=control
80 #pragma HLS INTERFACE m_axi port=in2 bundle=gmem2 max_widen_bitwidth=4
81 #pragma HLS INTERFACE s_axilite port=in2 bundle=control
82 #pragma HLS INTERFACE axis port=out2
83 #pragma HLS INTERFACE s_axilite port=size2 bundle=control
84 #pragma HLS INTERFACE m_axi port=in3 bundle=gmem3 max_widen_bitwidth=4
85 #pragma HLS INTERFACE s_axilite port=in3 bundle=control
86 #pragma HLS INTERFACE axis port=out3
87 #pragma HLS INTERFACE s_axilite port=size3 bundle=control
88 #pragma HLS INTERFACE m_axi port=in4 bundle=gmem4 max_widen_bitwidth=4
89 #pragma HLS INTERFACE s_axilite port=in4 bundle=control
90 #pragma HLS INTERFACE axis port=out4
91 #pragma HLS INTERFACE s_axilite port=size4 bundle=control
92 #pragma HLS INTERFACE m_axi port=in5 bundle=gmem5 max_widen_bitwidth=4
93 #pragma HLS INTERFACE s_axilite port=in5 bundle=control
94 #pragma HLS INTERFACE axis port=out5
95 #pragma HLS INTERFACE s axilite port=size5 bundle=control
96 #pragma HLS INTERFACE m_axi port=in6 bundle=gmem6 max_widen_bitwidth=4
97 #pragma HLS INTERFACE s_axilite port=in6 bundle=control
98 #pragma HLS INTERFACE axis port=out6
99 #pragma HLS INTERFACE s_axilite port=size6 bundle=control
100 #pragma HLS INTERFACE s_axilite port=pulse bundle=control
101
       PushAll(
         in0, out0, size0,
103
         in1, out1, size1,
104
         in2, out2, size2,
         in3, out3, size3,
106
         in4, out4, size4,
         in5, out5, size5,
108
         in6, out6, size6,
109
         pulse);
113
```

The mm2s kernel reads data from memory (inX) and writes it to AXI streams (outX) across 7 channels, repeating the transfer pulse times. Each Push sends data as ap\_axiu packets on the HLS stream. If target.write() blocks, it means the receiver (e.g., DMA) is not reading fast enough: the FIFO is full and the kernel waits. This indicates that the next sample cannot be sent until the previous one is consumed.

## Appendix C

# Kernel codes

The kernel codes reported in this Section are the custom C++ functions run on AIE (or AIE-ML) used in the ACE development.

#### C.1 Single IIR filter order 2 - buffer

```
_{1} #include < adf.h>
2 #include <aie_api/aie.hpp>
4 #include <aie_api/utils.hpp>
7 x =
      [x, 0, 0, 0]
p = [x(n-1), x(n-2), y(n-1), y(n-2)]
c = [b2, b3, a2, a3, b1, g, 0, 0]
y = [y, 0, 0, 0]
p_out = p updated for next iteration
12
13
void IIR_2_custom
      adf::input_buffer<float,adf::extents<4>>& __restrict x, // input sample
      adf::input_buffer<float,adf::extents<8>>& __restrict c, // coefficients
17
      adf::input_buffer<float,adf::extents<4>>& __restrict p, // partial
      adf::output_buffer<float,adf::extents<4>>& __restrict p_out, // partial
      results, output
      adf::output_buffer<float,adf::extents<4>>& __restrict y // output
20
      // iterators for in/out buffers
      auto x_int_I=aie::begin_vector<4>(x);
23
      auto v_coeff_I=aie::begin_vector<8>(c);
      auto p int I=aie::begin vector <4>(p);
      auto p_out_I=aie::begin_vector<4>(p_out);
26
      auto y_I=aie::begin_vector<4>(y);
2.7
      // iterator -> vector
      auto x_int = *x_int_I++;
30
      auto p_int = *p_int_I++;
```

```
auto v_coeff = *v_coeff_I++;
32
       x_{int}[0] = x_{int}[0] * v_{coeff}[5]; // *g multiplication just on input
33
34
       // v_{\text{sample should have }} [x(t-1), x(t-2), y(t-1), y(t-2), x, 0, 0]
35
       aie::vector<float, 8> v_sample = aie::concat(p_int, x_int);
36
       // operations for output
38
       auto v_mul = aie::mul(v_sample, v_coeff);// multiplication
39
       float y_{int} = aie :: reduce_add(v_mul.to_vector < float > (0)); // sum of
40
       aie :: vector < float, 4 > v_result = aie :: zeros < float, 4 > ();
41
       v_result.push(y_int);// out must be a vector
42
       *y\_I \!\!+\!\!\!+\!\!\!=\!\! v\_result;
44
       // partial results update
45
       aie::vector<float, 4> p_out_v;
46
       p_{out}v[1] = p_{int}[0];
       p \text{ out } v[0] = x \text{ int } [0];
48
       p_{out}v[3] = p_{int}[2];
49
       p_out_v[2] = v_result[0];
51
       *p\_out\_I++=p\_out\_v;
52
```

This code describes with AIE APIs an IIR filter of order 2, as explained in Section 4.1.1. This kernel, with buffer inputs, would also need a manual handshake of the AXI protocol. For this reason, in the next versions, coefficients and partial products are stored in AIE memory, while the input type stream will allow automatic AXI4 handshake management.

# C.2 Single IIR filter order 2 - stream (not optimal)

```
1 #include <adf.h>
2 #include <aie_api/aie.hpp>
3 #include <aie_api/aie_adf.hpp>
4 #include <aie_api/utils.hpp>
6 /*
7 \times = [X]
       [x(n-1), x(n-2), y(n-1), y(n-2)]
       \left[ \, b2 \, , \ b3 \, , \ -a2 \, , \ -a3 \, , \ b1 \ g \, , \ 0 \, , \ 0 \, \right]
y = [y]
11 */
12
  void IIR_2_custom
14
       input_stream<float>* x, // input sample
       output_stream<float>* y // output
16
  ) {
17
       static aie::vector<float, 4> p = aie::zeros<float, 4>(); // clear
18
      partial results
       const aie::vector<float, 8> coeff(0.6868f, 0.f, 0.9973f, 0.f, 0.6868f,
19
      1.f, 0.f, 0.f);
20
      // input reading
```

```
float x_{int} = readincr(x);
22
23
      // to vector
      x_{int} = x_{int}*coeff[5]; // *g multiplication just on input sample
26
      aie::vector<float, 4> x_int_v = aie::zeros<float,4>();
      x_int_v.push(x_int);// out must be a vector
29
      // \text{ v\_sample should have } [x(t-1), x(t-2), y(t-1), y(t-2), x, 0, 0, 0]
30
       aie::vector<float, 8> v_sample = aie::concat(p, x_int_v);
31
      // operations for output
33
      auto v_mul = aie::mul(v_sample, coeff);// multiplication
34
      float y_{int} = aie :: reduce_add(v_mul.to_vector < float > (0)); // sum of
      v_mul
      writeincr(y, y_int);
36
37
      // partial results update
      p[1] = p[0];
39
      p[0] = x_{int};
40
      p[3] = p[2];
41
      p[2] = y_{int};
42
43
```

This is a streaming version of the single IIR filter of order 2, but calling scalar multiplication and sum functions as '\*' and '+' would result in the call of an appropriate function that would waste too many clock cycles (AIE does not contain a scalar FP unit). This operation is transformed into a vectorial one in the next versions of the code.

#### C.3 Single IIR filter order 2 - stream

```
1 #include <adf.h>
2 #include <aie_api/aie.hpp>
3 #include <aie_api/aie_adf.hpp>
4 #include <aie_api/utils.hpp>
6 /*
7 \times = [x]
p = [x(n-1), x(n-2), y(n-1), y(n-2)]
a = [b2, b3, -a2, -a3, b1, g, 0, 0]
y = [y]
11
12
void IIR_2_custom
14
      input\_stream < float > * x, // input sample
16
      output_stream<float>* y // output
  ) {
17
      static aie::vector<float, 4> p = aie::zeros<float, 4>(); // clear
      partial results
      const aie::vector<float, 8> coeff(0.6868f, 0.f, 0.9973f, 0.f, 0.6868f,
19
      1.f, 0.f, 0.f);
20
      // input reading
      float x_int = readincr(x);
22
23
```

```
// -> vector
24
       aie::vector<float, 4> x_int_v(x_int, 0.0f, 0.0f, 0.0f);
25
       x_{int}v = aie :: mul(x_{int}v, coeff[5]) . to_vector < float > (0);
26
27
       // v_{\text{sample should have }} [x(t-1), x(t-2), y(t-1), y(t-2), x, 0, 0, 0]
28
       aie::vector<float, 8> v_sample = aie::concat(p, x_int_v);
29
      // operations for output
31
      auto v_mul = aie::mul(v_sample, coeff);// multiplication
32
       float y_{int} = aie :: reduce_add(v_{mul}.to_vector < float > (0)); // sum of
       writeincr(y, y_int);
34
35
       // partial results update
      p[1] = p[0];
      p[0] = x_{int};
38
      p[3] = p[2];
39
      p[2] = y_{int};
40
41
```

#### 64 bit version

```
1 #include <adf.h>
2 #include <aie_api/aie.hpp>
3 #include <aie_api/aie_adf.hpp>
4 #include <aie_api/utils.hpp>
6 /*
7 \ x =
      [x]
p = [x(n-1), x(n-2), y(n-1), y(n-2)]
a = [b2, b3, -a2, -a3, b1, g, 0, 0]
y = y
11 */
12
  void IIR_1
13
14
      input_stream<uint64>* x, // input sample
      output_stream<uint64>* y // output
17 ) {
      static aie::vector<float, 4> p = aie::zeros<float, 4>(); // clear
      partial results
      const aie::vector<float, 8> coeff(0.0f, 0.f, 0.0f, 0.f, 1.f, 1.f, 0.f,
19
      0.f);
21
      // input reading
      uint64 x_real = readincr(x);
23
      float* floats = (float*)&x_real;
24
      float x_int = floats[0];
25
26
      // -> vector
      aie::vector<float, 4> x int v(x int, 0.0f, 0.0f, 0.0f);
28
      x_{int}v = aie :: mul(x_{int}v, coeff[5]) .to_vector < float > (0);
29
30
      // v_{\text{sample should have }} [x(t-1), x(t-2), y(t-1), y(t-2), x, 0, 0]
31
      aie::vector<float, 8> v_sample = aie::concat(p, x_int_v);
33
      // operations for output
34
```

```
auto v_mul = aie::mul(v_sample, coeff);// multiplication
35
       float y_{int} = aie :: reduce\_add(v_{int}.to\_vector < float > (0)); // sum of
36
      v_{mul}
       float raw[2] = \{y_int, y_int\};
37
       uint64 \ y_real = *((uint64_t*)raw);
38
       writeincr(y, y_real);
39
       // partial results update
41
      p[1] = p[0];
42
      p[0] = x_int_v[0];
      p[3] = p[2];
44
      p[2] = y_{int};
45
46
```

This is the final version of the IIR filter of order 2. Input and output are 64-bit in order to allow a single sample streaming (min is 64-bit from PL to AIE and vice versa).

# C.4 Single IIR filter order 4 - stream (improved)

```
1 #include <adf.h>
2 #include <aie_api/aie.hpp>
3 #include <aie_api/aie_adf.hpp>
4 #include <aie_api/utils.hpp>
7 X =
      [x]
      [x(n-1), y(n-1), x(n-2), y(n-2), x(n-3), y(n-3), x(n-4), y(n-4)]
8 p =
g = [b2, -a2, b3, -a3, b4, -a4, b5, -a5, b1, g, 0, 0, 0, 0, 0]
y = [y]
11
  */
12
  void IIR_4_custom
14
      input_stream < float >* x, // input sample
      output_stream<float>* y // output
16
  ) {
17
      static aie::vector<float, 8> p = aie::zeros<float, 8>(); // clear
18
      partial results
      const aie::vector< float, 16> coeff(0.0567984f, 1.6190948f, 0.0820072f,
      -1.1351098f, 0.0567984f, 0.3802160f, 0.0157913f, -0.0547274f, 0.0157913f
      , 1.0f, 0.f, 0.f, 0.f, 0.f, 0.f, 0.f);
20
      // iterators for in/out buffers
21
      float x_{int} = readincr(x);
23
      // input reading
      aie::vector<float, 8> x int v(x) int, 0.0f, 0.0f, 0.0f, 0.0f, 0.0f, 0.0f
26
      , 0.0 f);
      x_{int}v = aie :: mul(x_{int}v, coeff[5]) .to_vector < float > (0);
2.7
28
      // v_{\text{sample should have }} [x(t-1), x(t-2), y(t-1), y(t-2), x, 0, 0]
29
      aie :: vector < float, 16 > v_sample = aie :: concat(p, x_int_v);
30
```

```
// operations for output
auto v_mul = aie::mul(v_sample,coeff);// multiplication
float y_int = aie::reduce_add(v_mul.to_vector<float>(0));// sum of
v_mul
writeincr(y, y_int);

// partial results update
p.push(y_int);
p.push(x_int);
```

This is an order-4 version of the previous code.

#### C.5 Controller

```
1 #include <adf.h>
 2 #include <aie_api/aie.hpp>
 3 #include <aie_api/aie_adf.hpp>
 4 #include <aie_api/utils.hpp>
 5
 6 /*
 7 x =
               X
               [x(n-1), y(n-1), x(n-2), y(n-2)]
 p = 0
 a = [b2, -a2, b3, -a3, b1, g, 0, 0]
y = [y]
    */
11
13 // constants' section
static const bool LIMITER_ENABLE = true;
static const float LIMITER LOWER LIMIT = -30.0 \,\mathrm{f};
static const float LIMITER UPPER LIMIT = 30.0 f;
17 static const bool INTEGRATOR_RST = false;
    0.0 \, f, 0.0 \, f); //13
     static const aie::vector<float, 8> c_B(0.0f, 0.0f, 0.0f, 0.0f, 350.0f, 1.0f
              , 0.0 f, 0.0 f); //14
     static const aie::vector<float, 8> c_C(0.0f, 0.0f, 0.0f, 0.0f, 1.701f, 1.0f
              , 0.0 f, 0.0 f); //15
static const aie::vector<float, 8 > c_D(0.0f, 0.0f, 0.0f, 0.0f, 1.0f, 1.0f,
              0.0 f, 0.0 f); //16
static const aie::vector<float, 8 > c_E(0.0f, 0.0f, 
              0.0 f, 0.0 f); //17
23
    void IIR_CONTROLLER
24
25
               input\_stream < float > * x, // input sample
               output_stream<float>* y // output
27
     ) {
28
29
               // partial results
30
               static aie::vector<float, 4> p_A = aie::zeros<float, 4>();
31
               static aie::vector<float, 4> p_B = aie::zeros<float, 4>();
32
               static aie::vector<float, 4> p_C = aie::zeros<float, 4>();
33
               static aie::vector<float, 4> p_D = aie::zeros<float, 4>();
34
               static aie::vector<float, 4> p_E = aie::zeros<float, 4>();
35
36
```

```
// iterators for in/out buffers
37
      float x_int = readincr(x);
38
39
      //////// FILTER A
40
      // -> vector
41
      aie::vector<float, 4> x_{int}vA(x_{int}, 0.0f, 0.0f, 0.0f);
42
      x_{int}vA = aie :: mul(x_{int}vA, c_A[5]) . to_vector < float > (0);
43
44
      // v_{\text{sample should be }} [x(t-1), y(t-1), x(t-2), y(t-2), x, 0, 0, 0]
45
      aie::vector<float, 8> v_sample_A = aie::concat(p_A, x_int_vA);
46
47
      // operations for output
48
      auto v_mul = aie::mul(v_sample_A,c_A);// multiplication
49
      float y_A = aie :: reduce_add(v_mul.to_vector < float > (0) . extract < 4 > (0)); //
50
      sum of v_mul
      aie::vector<float, 4> y_A_v(y_A, 0.0f, 0.0f, 0.0f);
      y_A = aie :: add(y_A_v, v_mul.to_vector < float > (0) . extract < 4 > (1)) [0];
      // partial results update
54
      p_A. push(y_A);
      p_A. push(x_int_vA[0]);
      //////// FILTERS B,C,E
58
      // gain multiplication
59
      aie::vector<float, 4> x_step2(y_A, y_A, y_A, 0.0 f);
      aie::vector<float, 4> g_vector(c_B[5], c_C[5], c_E[5], 0.0 f);
      x_{step2} = aie :: mul(x_{step2}, g_{vector}).to_{vector} < float > (0);
63
      // product
      aie::vector<float, 16> v_sample_BCE = aie::concat(p_B, p_C, p_E,
65
     x step2);
      aie::vector<float, 16> v_coeff_BCE = aie::concat(c_B.extract <4>(0), c_C
66
      \det(4>0), c_E.extract (4>0), aie::vector (1) 4(2) 4(2) 4(2) 6.
      c_{E[4]}, 0.0f);
      auto v_mul_BCE = aie::mul(v_sample_BCE, v_coeff_BCE);// multiplication
67
68
      // reduction for the different filters
69
      aie::vector<float, 4> y B(aie::reduce add(v mul BCE.extract<4>(0).
70
      to\_vector < float > (0)), 0.0f, 0.0f, 0.0f);
      y_B = aie :: add(y_B, v_mul_BCE.to_vector < float > (0).extract < 4 > (3));
71
      aie::vector<float, 4> y_C(0.0f, aie::reduce_add(v_mul_BCE.extract<4>(1)
72
      . to_{vector} < float > (0) , 0.0 f, 0.0 f ;
      y_C = aie :: add(y_C, v_mul_BCE.to_vector < float > (0).extract < 4 > (3));
73
      aie::vector<float, 4> y_E(0.0f, 0.0f, aie::reduce_add(v_mul_BCE.extract
      <4>(2).to_vector<float>(0), 0.0 f;
      y_E = aie :: add(y_E, v_mul\_BCE.to\_vector < float > (0).extract < 4 > (3));
75
76
      // partial results update
77
      p_B. push(y_B[0]);
      p B.push(x step2[0]);
79
      p_C. push(y_C[1]);
80
      p_C. push (x_step2 [1]);
81
      p_E. push (y_E[2]);
      p_E. push(x_step2[2]);
83
84
      /////// OUT PART
85
86
      // partial sum
      auto y_BE = aie :: add(y_B, aie :: shuffle_down(y_E, 2));
87
```

```
float y_mux;
88
89
       // multiplexer
       if (INTEGRATOR_RST) {
91
          y_{mux} = 0.0 f;
92
      }else{
93
          y_mux = p_D[1];
      aie::vector<float, 4> y_mux_v(y_mux, 0.0 f, 0.0 f, 0.0 f);
96
97
      //final sum
98
      auto y_fin = aie :: add(y_BE, y_mux_v);
99
       writeincr(y, y_fin[0]);
100
      /////// LIMITER
      aie::vector<float, 4> y_lim_in(aie::add(y_mux_v, aie::shuffle_down(y_C
104
      (1))[0],0.0f,0.0f,0.0f);
      auto y lim out = y lim in;
      if (LIMITER ENABLE) {
106
          if (aie::gt(y_lim_in, LIMITER_UPPER_LIMIT).test(0)) {
              y_{lim_out}[0] = LIMITER_UPPER_LIMIT;
          } else if (aie::lt(y_lim_in, LIMITER_LOWER_LIMIT).test(0)) {
109
              y_lim_out [0] = LIMITER_LOWER_LIMIT;
      }
112
113
114
      /////// FILTER D
      auto x_D = aie :: mul(y_lim_out, c_D[5]) . to_vector < float > (0);
      aie::vector<float, 8> v_sample_D = aie::concat(p_D, x_D);
      v_{mul} = aie :: mul(v_{sample_D}, c_D); // multiplication
118
      119
       sum of v_mul
      aie::vector<float, 4> y_D_v(y_D, 0.0f, 0.0f, 0.0f);
120
      y_D = aie :: add(y_D_v, v_mul.to_vector < float > (0) . extract < 4 > (1)) [0];
      p D. push(y D);
      p D. push (x D[0]);
123
124
```

#### 64 bit version

```
1 #include <adf.h>
2 #include <aie_api/aie.hpp>
3 #include <aie_api/aie_adf.hpp>
4 #include <aie_api/utils.hpp>
6 /*
7 \times = [x]
p = [x(n-1), y(n-1), x(n-2), y(n-2)]
c = [b2, -a2, b3, -a3, b1, g, 0, 0]
y = [y]
11
  */
12
13 // constants' section
static const bool LIMITER_ENABLE = false;
static const float LIMITER LOWER LIMIT = 0.0 f;
static const float LIMITER_UPPER_LIMIT = 0.0 f;
```

```
17 static const bool INTEGRATOR RST = false;
18 static const aie::vector<float, 8 > c_A(0.0f, 0.0f, 0.0f, 0.0f, 1.0f, 1.0f,
                 0.0 f, 0.0 f); //3
19 static const aie::vector<float, 8> c_B(0.0f, 0.0f, 0.0f, 0.0f, 0.16f, 1.0f,
                   0.0 \, f, 0.0 \, f); //4
static const aie::vector<float, 8 > c_C(0.0f, 0.0f, 
                 0.0 f, 0.0 f); //5
     static const aie::vector<float, 8> c D(0.0f, 0.0f, 0.0f, 0.0f, 0.0f, 0.0f,
                 0.0 f, 0.0 f); //6
static const aie::vector<float, 8 > c_E(0.0f, 0.0f, 
                 0.0 \, f, 0.0 \, f); //7
23
      void IIR CONTROLLER 1
24
25
                  input_stream<uint64>* x, // input sample
26
                  output_stream<uint64>* y // output
27
      ) {
28
29
                  // partial results
30
                   static aie::vector<float , 4> p_A = aie::zeros<float , 4>();
31
                   static aie::vector<float, 4> p_B = aie::zeros<float, 4>();
                  static aie::vector<float , 4> p_C = aie::zeros<float , 4>();
                   static aie::vector<float , 4> p_D = aie::zeros<float , 4>();
34
                  static aie::vector<float, 4> p_E = aie::zeros<float, 4>();
35
36
                  // input section
                  uint64 \times real = readincr(x);
38
                   float* floats = (float*)&x real;
39
                  float x_int = floats[0];
40
                  ////////// FILTER A
42
                  // -> vector
43
                  aie::vector<float, 4> x_{int}vA(x_{int}, 0.0f, 0.0f, 0.0f);
44
                  x_{int}_vA = aie :: mul(x_{int}_vA, c_A[5]) . to_vector < float > (0);
45
46
                  // v_sample should be [x(t-1), y(t-1), x(t-2), y(t-2), x, 0, 0, 0]
47
                  aie::vector<float, 8> v_sample_A = aie::concat(p_A, x_int_vA);
49
                  // operations for output
50
                  auto v_mul = aie::mul(v_sample_A,c_A);// multiplication
                  float y_A = aie :: reduce_add(v_mul.to_vector < float > (0) . extract < 4 > (0)); //
                   sum of v mul
                  aie::vector<float, 4> y_A_v(y_A, 0.0f, 0.0f, 0.0f);
                 y_A = aie :: add(y_A_v, v_mul.to_vector < float > (0).extract < 4 > (1)) [0];
                  // partial results update
56
                 p_A. push (y_A);
57
                 p_A. push (x_int_vA [0]);
58
                  //////// FILTERS B,C,E
60
                  // gain multiplication
61
                  aie::vector<float, 4> x_step2(y_A, y_A, y_A, 0.0 f);
62
                  x \text{ step } 2 = \text{aie} :: \text{mul}(x \text{ step } 2, g \text{ vector}) \text{ to } \text{vector} < \text{float} > (0);
64
65
                  // product
66
                  aie::vector<float, 16> v_sample_BCE = aie::concat(p_B, p_C, p_E,
                x step2);
```

```
aie::vector<float, 16> v_coeff_BCE = aie::concat(c_B.extract <4>(0), c_C
68
       \det(4>0), c_E.extract (4>0), aie::vector (1) (2) (2) (2) (3)
       c_{E}[4], 0.0f);
       auto v_mul_BCE = aie::mul(v_sample_BCE, v_coeff_BCE);// multiplication
69
70
       // reduction for the different filters
       aie::vector<float, 4> y_B(aie::reduce_add(v_mul_BCE.extract<4>(0).
       to\_vector < float > (0)), 0.0f, 0.0f, 0.0f);
       y_B = aie :: add(y_B, v_mul_BCE.to_vector < float > (0).extract < 4 > (3));
73
       aie::vector<float, 4> y_C(0.0f, aie::reduce_add(v_mul_BCE.extract<4>(1)
74
       to_{\text{vector}} < \frac{\text{float}}{(0)}, 0.0 \, \text{f}, 0.0 \, \text{f};
       y_C = aie :: add(y_C, v_mul_BCE.to_vector < float > (0).extract < 4 > (3));
75
       aie::vector<float, 4> y_E(0.0f, 0.0f, aie::reduce_add(v_mul_BCE.extract
76
       <4>(2).to_vector<float>(0), 0.0 f);
       y_E = aie :: add(y_E, v_mul_BCE.to_vector < float > (0).extract < 4 > (3));
77
78
       // partial results update
79
       p_B. push(y_B[0]);
80
       p B.push(x step2[0]);
81
       p_C. push(y_C[1]);
82
       p_C. push (x_step2 [1]);
83
       p_E. push (y_E[2]);
       p_E. push(x_step2[2]);
85
86
       /////// OUT PART
87
       // partial sum
       auto y_BE = aie :: add(y_B, aie :: shuffle_down(y_E, 2));
89
       float y_mux;
90
       // multiplexer
       if (INTEGRATOR RST) {
           y_{mux} = 0.0 f;
94
       } else {
95
96
           y_{mux} = p_D[1];
97
       aie::vector<float, 4> y_mux_v(y_mux, 0.0 f, 0.0 f, 0.0 f);
98
       //final sum
100
       auto y_fin = aie :: add(y_BE, y_mux_v);
       float raw[2] = \{y_fin[0], y_fin[0]\};
       uint64 \ y_real = *((uint64_t*)raw);
103
       writeincr(y, y_real);
104
105
106
       /////// LIMITER
       aie::vector<float, 4> y_lim_in(aie::add(y_mux_v, aie::shuffle_down(y_C
108
       (1))[0],0.0f,0.0f,0.0f);
       auto y_lim_out = y_lim_in;
109
       if (LIMITER ENABLE) {
110
            if (aie::gt(y lim in, LIMITER UPPER LIMIT).test(0)) {
                y \lim_{n \to \infty} \operatorname{out}[0] = \operatorname{LIMITER} \operatorname{UPPER} \operatorname{LIMIT};
            } else if (aie::lt(y_lim_in, LIMITER_LOWER_LIMIT).test(0)) {
113
                y_lim_out [0] = LIMITER_LOWER_LIMIT;
114
       }
116
117
118
       /////// FILTER D
```

This is the 64-bit version of the controller block, whose behavior is described in 4.1.1.

### C.6 Input block

```
1 #include <adf.h>
2 #include <aie_api/aie.hpp>
3 #include <aie_api/aie_adf.hpp>
4 #include <aie_api/utils.hpp>
5
6 /*
7 x =
      X
      [x(n-1), y(n-1), x(n-2), y(n-2)]
p = 0
a = [b2, -a2, b3, -a3, b1, g, 0, 0]
y = y
11
13 // constants' section
static const float BIAS = 0.0 f;
static const aie::vector<float, 8> c_r1(0.0f, 0.0f, 0.0f, 0.0f, 1.0f, 1.0f,
       0.0 \, f, 0.0 \, f); //22
static const aie::vector<float, 8> c_r2(0.0f, 0.0f, 0.0f, 0.0f, 1.0f, 1.0f,
       0.0 \, f, 0.0 \, f); //21
  static const aie::vector<float, 8> c_s1(0.0f, 0.0f, 0.0f, 0.0f, 1.0f, 1.0f,
       0.0 \, f, 0.0 \, f); //24
  static const aie::vector<float, 8> c_s2(1.0f,
     0.980093260759718254426786643307423219085\,f\,\,,\,\,\,0.0\,f\,\,,\,\,\,0.0\,f\,\,,\,\,\,1.0\,f\,\,,
      0.009953369620140867582436250415867107222f, 0.0f, 0.0f); //23
19
  void IIR_input_block_1
21
22
      input_stream<uint64>* x, // ratio, sum
23
      output_stream<uint64>* y // ratio, sum
24
25
26
27
      // partial results
      static aie::vector<float , 4> p_r1 = aie::zeros<float , 4>();
      static aie::vector<float, 4> p_r2 = aie::zeros<float, 4>();
      static aie::vector<float, 4> p_s1 = aie::zeros<float, 4>();
30
      static aie::vector<float, 4> p_s2 = aie::zeros<float, 4>();
31
32
      // input reading
33
      uint64 x_real = readincr(x);
34
      float * floats = (float *)&x_real;
35
      aie::vector<float, 4> x_int_v1(floats[0], floats[1], 0.0f, 0.0f);
36
```

```
37
38
      ////////////// FILTER r1, s1
      aie::vector<float, 4> g_int_v1(c_r1[5], c_s1[5], 0.0f, 0.0f);
40
      x_{int}v1 = aie :: mul(x_{int}v1, g_{int}v1).to_vector < float > (0);
41
42
      // v_{\text{sample should be }}[x(t-1), y(t-1), x(t-2), y(t-2), x, 0, 0, 0]
      aie::vector<float, 16> v_sample_1 = aie::concat(p_r1, p_s1, x_int_v1,
44
     aie::zeros<float, 4>());
      aie:: vector < float, 16 > v_coeff_1 = aie:: concat(c_r1.extract < 4 > (0), c_s1
     . \text{ extract } < 4 > (0), \text{ aie :: vector} < \frac{\text{float}}{1}, 4 > (c_r1[4], c_s1[4], 0.0f, 0.0f), \text{ aie } 
     :: zeros < float, 4 > ());
46
      // operations for output
      auto v_mul = aie::mul(v_sample_1, v_coeff_1);// multiplication
48
      float y r1 = aie::reduce add(v mul.to vector<float>(0).extract<4>(0));
49
     // sum of v_mul
      float y_s1 = aie :: reduce\_add(v_mul.to\_vector < float > (0) . extract < 4 > (1));
     // sum of v mul
      aie::vector<float, 4> y_1_v(y_r1, y_s1, 0.0f, 0.0f);
      y_1v = aie :: add(y_1v, v_mul.to_vector < float > (0) . extract < 4 > (2));
      // partial results update
54
      p_r1.push(y_1_v[0]);
      p_r1.push(x_int_v1[0]);
56
      p_s1.push(y_1_v[1]);
      p_s1.push(x_int_v1[1]);
58
59
      //////// FILTERS r2, s2
      // -> vector
61
      aie::vector<float, 4> g_int_v2(c_r2[5], c_s2[5], 0.0f, 0.0f);
      y_1_v = aie :: mul(y_1_v, g_int_v^2) . to_vector < float > (0);
64
      // v_{\text{sample should be }} [x(t-1), y(t-1), x(t-2), y(t-2), x, 0, 0]
65
      aie::vector<float, 16> v_sample_2 = aie::concat(p_r2, p_s2, y_1_v, aie
66
     :: zeros < float, 4 > ());
      aie::vector<float, 16> v coeff 2 = aie::concat(c r2.extract<4>(0), c s2
     :: zeros < float, 4 > ());
68
      // operations for output
69
      v_mul = aie::mul(v_sample_2, v_coeff_2);// multiplication
70
      float y_r^2 = aie :: reduce\_add(v_mul.to\_vector < float > (0) . extract < 4 > (0));
71
     // sum of v_mul
      // sum of v_mul
      aie::vector<float, 4> y_2_v(y_r2, y_s2, 0.0f, 0.0f);
73
      y_2v = aie :: add(y_2v, v_mul.to_vector < float > (0) . extract < 4 > (2));
74
75
      // partial results update
76
      p r2.push(y 2 v[0]);
      p_r2.push(y_1_v[0]);
      p_s2.push(y_2_v[1]);
      p_s2.push(y_1_v[1]);
80
81
      // BIAS sum
82
      aie::vector<float, 4> bias_v(BIAS, 0.0f, 0.0f, 0.0f);
83
      y 2 v = aie :: add(y 2 v, bias v);
84
```

```
// output
float raw[2] = {y_2_v[0], y_2_v[1]};
uint64 y_real = *((uint64_t*)raw);
writeincr(y, y_real);
}
```

This is the input block, whose behavior is described in 4.1.1.

# **Bibliography**

- [1] AMD. Versal: The First Adaptive Compute Acceleration Platform (ACAP), 2020. URL https://docs.amd.com/v/u/en-US/wp505-versal-acap. Accessed: 2025-04-01.
- [2] AMD. Versal ACAP DSP Engine Architecture Manual (AM004), 2022. URL https://docs.amd.com/r/en-US/am004-versal-dsp-engine/Overview?tocId=tItzyS6 FCFNHd7m~GFfslQ. Accessed: 2025-04-01.
- [3] AMD. Versal Adaptive SoC AI Engine Architecture Manual, 2023. URL https://docs.amd.com/viewer/book-attachment/juyVkQQPlmZlTyA3UczsEA/IHjwvrFtwe3n2UxK4C0iHQ-juyVkQQPlmZlTyA3UczsEA. Accessed: 2025-04-01.
- [4] AMD. AXI4-Stream Infrastructure IP Suite v3.0, 2023. URL https://docs.amd.com/viewer/book-attachment/iDBE\_pP8gVw30\_auFL6zrQ/1SslvnphHVUfedW4\_R2 Hxg-iDBE\_pP8gVw30\_auFL6zrQ. Accessed: 2025-07-11.
- [5] AMD. Versal Adaptive SoC AIE-ML Architecture Manual (AM020), 2024. URL https://docs.amd.com/r/en-US/am020-versal-aie-ml. Accessed: 2025-04-01.
- [6] AMD. AI Engine Tools and Flows User Guide, 2024. URL https://docs.amd.com/r/en-US/ug1076-ai-engine-environment. Accessed: 2025-04-04.
- [7] AMD. AI Engine Kernel and Graph Programming Guide (UG1079), 2024. URL https://docs.amd.com/r/en-US/ug1079-ai-engine-kernel-coding?tocId=\_G~tNVucqwCOCCt0l\_v6bA. Accessed: 2025-04-01.
- [8] AMD. VCK190 Evaluation Board User Guide (UG1366), 2024. URL https://docs.amd.com/r/en-US/ug1366-vck190-eval-bd. Accessed: 2025-06-012.
- [9] AMD. Vitis Model Composer User Guide (UG1483), 2024. URL https://docs.amd.com/r/en-US/ug1483-model-composer-sys-gen-user-guide. Accessed: 2025-04-02.
- [10] AMD. Versal Adaptive SoC System and Solution Planning Methodology Guide (UG1504), 2024. URL https://docs.amd.com/r/en-US/ug1504-acap-system-solution-planning-methodology. Accessed: 2025-04-03.
- [11] AMD. Versal Adaptive SoC Programmable Network on Chip and Integrated Memory Controller v1.1, 2025. URL https://docs.amd.com/r/en-US/pg313-network-on-chip. Accessed: 2025-06-12.

- [12] AMD. Vivado Design Suite User Guide: Programming and Debugging (UG908), 2025. URL https://docs.amd.com/r/en-US/ug908-vivado-programming-debugging/Trigger-At-Startup. Accessed: 2025-06-25.
- [13] M. R. Casu and L. Macchiarulo. Adaptive latency-insensitive protocols. *IEEE Design & Test of Computers*, 24(5):442–552, 2007.
- [14] Y. P. Chen, P. Maillard, R. D. Veggalam, S. R. Madem, E. Crabill, and J. Barton. 64mev proton single-event evaluation of xilinx single event mitigation (xilsem) firmware on 7nm versal<sup>TM</sup> acap devices. *IEEE Radiation Effects Data Workshop (REDW) (in conjunction with 2022 NSREC)*, 2022.
- [15] F. Flores, O. L. Sánchez, L. J. Álvarez Ruiz de Ojeda, M. D. V. Peña, and J. M. V. Sánchez. Acceleration of a compute-intensive algorithm for power electronic converter control using versal ai engines. *Conference on Design of Circuits and Integrated Circuits (DCIS)*, 2024.
- [16] J. Lei and E. S. Quintana-Orti. Mapping Parallel Matrix Multiplication in Goto-BLAS2 to the AMD Versal ACAP for Deep Learning. *PECS* '24, 1607.
- [17] P. Maillard, Y. P. Chen, J. Barton, and M. L. Voogel. Single event latchup (sel) and single event upset (seu) evaluation of xilinx 7nm versal<sup>TM</sup> acap programmable logic (pl). *IEEE Radiation Effects Data Workshop (REDW)v*, 2021.
- [18] A. Mège and N. Wazad. Challenges in the board implementation of VERSAL ACAP on a space board. Airbus Defence & Space, Space Electronics EDHPC, 2023.
- [19] L. M. Ni and P. K. McKinley. A Survey of Wormhole Routing Techniques in Direct Networks. *Computer*, 26(2):62–76, 1993.
- [20] R. F. Stengel. Toward intelligent flight control. *IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS*, 23(6):1699–1717, 1993.
- [21] R. Sticca and A. Demichelis. Model Based Design (MBD) for hardware FPGA in Digital Flight Control Computer (DFCC). POLARIS Innvation Journal, 24:13–18, 2015.
- [22] I. Swarbrick, D. Gaitonde, S. Ahmad, B. Jayadev, J. Cuppett, A. Morshed, B. Gaide, and Y. Arbel. VERSAL NETWORK-on-CHIP (NoC). *IEEE Symposium on High-Performance Interconnects (HOTI)*, pages 13–17, 2019.
- [23] TE connectivity. LVDT tutorial. https://www.te.com/en/products/sensors/position-sensors/resources/lvdt-tutorial.html, 2017. Accessed: 2025-05-23.
- [24] Trenz Electronic. Amd versal ai edge evalboard with ve2302 device. https://shop.trenz-electronic.de/en/TE0950-03-EGBE21C-AMD-Versal-AI-Edge-Evalboard-with-VE2302-device-8-GB-DDR4-SDRAM-15-x12-cm, 2025. Accessed: 2025-04-03.
- [25] XD100. Vitis Tutorials: AI Engine Development (XD100), 2024. URL https://docs.amd.com/r/en-US/Vitis-Tutorials-AI-Engine-Development/Vitis-Tutorials-AI-Engine-Development-XD100. Accessed: 2025-26-05.