

Horizon 2020 Program (2014-2020)
FET-Open Novel ideas for radically new technologies
FETOPEN-01-2018-2019-2020



Architecting More than Moore – Wireless Plasticity for
Massive Heterogeneous Computer Architectures [†]

D5.1: Die Level Simulator and Tutorial

Contractual Date of Delivery	30/09/2020
Actual Date of Delivery	30/09/2020
Deliverable Security Class	Public
Editor	Alexandre Levisse
Contributors	EPFL (Leader)
Quality Assurance	S. Abadal (UPC), G. Karunaratne (IBM)

[†] This project is supported by the European Commission under the Horizon 2020 Program with Grant agreement no: 863337

Document Revisions & Quality Assurance

Deliverable Number	D5.1
Deliverable Responsible	EPFL
Work Package	WP5
Main Editor	Alexandre Levisse

Internal Reviewers

1. Geethan Karunaratne (IBM)
2. Sergi Abadal (UPC)

Revisions

Version	Date	By	Overview
1.0	02/09/2020	Alexandre Levisse	First Draft
1.1	14/09/2020	Alexandre Levisse	Draft including David Atienza's comments
1.2	16/09/2020	Alexandre Levisse	Draft including Marina Zapater's and Giovanni Ansaloni's comments
1.3	24/09/2020	Alexandre Levisse	Draft addressing Reviews from G. Karunaratne and S. Abadal
1.4	30/09/2020	Alexandre Levisse	Final version considering all the inputs and reviews from EPFL contributors and internal reviewers

Legal Disclaimer

The information in this document is provided “as is”, and no guarantee or warranty is given that the information is fit for any particular purpose. The above referenced consortium members shall have no liability to third parties for damages of any kind including without limitation direct, special, indirect, or consequential damages that may result from the use of these materials subject to any liability which is mandatory due to applicable law. © 2020 by WiPLASH Consortium.

Executive Summary

The main objective of D5.1 is to describe the development of the simulation framework that will be used for exploration in the WiPLASH project. This simulation framework, targeting at this point die-level architectures, covers several aspects. (1) the support for low-power RISC-V Instruction Set Architecture (ISA) in Gem5 simulator full system simulation mode. (2) the integration of machine-learning specific accelerators leveraging emerging resistive memories inside the processing unit. (3) the introduction of different flavors of wireless interconnect within the computing architecture. Ultimately, these elements will be considered as a base for the follow-up tasks in WiPLASH.

Abbreviations and Acronyms

NoC	Network-on-Chip
WL	Wireless
ISA	Instruction Set Architecture
RTL	Register Transfer Level
FS	Full System
HBM	High Bandwidth Memory
ABI	Application Binary Interface
SBI	Supervisor Binary Interface
AI	Artificial Intelligence
RRAM	Resistive Random Access Memories
PCM	Phase Change Memories
CM	Computational Memories
ADC	Analog to Digital Converter
EX	EXecution
CPU	Central Processing Unit
LSTM	Long Short Term Memory
CNN	Convolutional Neural Networks
RNN	Recurrent Neural Network
MAC	Multiply Accumulate
ROI	Region Of Interest
DRAM	Dynamic Random Access Memory
LLC	Last Level Cache
WFM	Wait For Memory
IO	Input-Output
LLCMPKI	Last Level Cache Misses per Kilo Instruction
IPS	Instructions Per Nanoseconds
THz	TeraHertz
MPSoC	Multiple Processor System on Chip
BM	Broadcast Memories
FIFO	First In First Out

The *WiPLASH* consortium is composed by:

UPC	Coordinator	Spain
IBM	Beneficiary	Switzerland
UNIBO	Beneficiary	Italy
EPFL	Beneficiary	Switzerland
AMO	Beneficiary	Germany
UoS	Beneficiary	Germany
RWTH	Beneficiary	Germany



IBM **Research** | Zurich



Table of Contents

DOCUMENT REVISIONS & QUALITY ASSURANCE.....	2
EXECUTIVE SUMMARY	4
ABBREVIATIONS AND ACRONYMS.....	5
TABLE OF CONTENTS.....	7
LIST OF FIGURES	9
LIST OF TABLES	10
1 INTRODUCTION	11
2 RISC-V INSTRUCTION SET ARCHITECTURE SUPPORT IN GEM5.....	12
2.1 STATE OF THE ART AND LIMITATIONS	12
2.2 SUMMARY OF THE MAIN CONTRIBUTIONS OF EPFL IN GXR5	13
3 MATRIX-VECTOR-MULTIPLICATION RRAM ACCELERATION IN GEM5	14
3.1 STATE OF THE ART.....	14
3.2 SYSTEM-LEVEL ARCHITECTURE.....	14
3.2.1 <i>System Integration of the IMC module</i>	14
3.2.2 <i>New IMC-Core Interface Instructions</i>	15
3.2.3 <i>Using the CM Interface instructions in a Recurrent Neural Network test case</i>	16
3.3 EXPERIMENTAL SETUP AND RESULTS.....	17
3.3.1 <i>Experimental setup</i>	17
3.3.2 <i>Performance and energy results</i>	19
3.4 CONCLUSION AND FUTURE WORKS	21
4 WIRELESS INTERCONNECT INSIDE THE SIMULATION FRAMEWORK	22
4.1 WIRED VERSUS WIRELESS INTERCONNECT.....	22
4.2 PREVIOUS WORKS ON WIRELESS INTERCONNECT	23
4.3 APPLICATION AND ARCHITECTURE TARGETS FOR WIRELESS INTERCONNECT	23
4.4 MODELING PARAMETERS AND IMPLEMENTATION OF WL INTERCONNECT IN GEM5.....	25
4.5 CONCLUSIONS AND PERSPECTIVES	26
5 CONCLUSIONS AND PERSPECTIVES	27
BIBLIOGRAPHY	28

List of Figures

Figure 1: A classic 5-stage in-order pipelined CPU, separated by stage registers, with instruction fetch, instruction decode, execute, memory, and write back stages. The CM core exists alongside the arithmetic logic unit (ALU) within the EX stage of the CPU.15

Figure 2: Weight storage in the CM core for an LSTM composed of a single cell and a dense layer. The last row of the array is reserved for the biases.16

Figure 3: All performance and energy results for the system architecture exploration study. Rows (1), (2), and (3) show the results for system configurations using the 0.8GHz, 1.3GHz, and 2.3GHz core frequencies, respectively. Columns (a), (b), (c), and (d) show the simulated time (s), instructions per ns (IPnS), memory intensity (LLCMPKI), and energy (J) results, respectively. The experiment number refers to the numbered LSTM networks (larger experiment number represents larger LSTM).....20

Figure 4: Full system average power distribution for the architectural exploration study. (a), (b), and (c) show the average power distribution for system configurations tested with 0.8GHz, 1.3GHz, and 2.3GHz cores, respectively.20

Figure 5: (a) illustration of 2.5D chiplet on interposer computing unit. Each chiplet contains a 2D NoC for local communication. Chiplet-to-chiplet communication is managed by a NoC on interposer (as described in [25]) with limited bandwidth and higher risks of congestion. (b) illustration of a wireless NoC in which any computing unit could communicate with any computing unit regardless of their physical organization.23

Figure 6: (a) floorplan of recent 14nm Samsung Exynos-type SoC [29] and 22nm IBM Power-8 processor [30]. (b) Floorplan of a OpenPiton tile-based 25cores architecture using a wired 5x5 2D mesh NoC topology [31].24

Figure 7 : (a) example of the core-to-core communication simulation benchmark implemented in WiPLASH. (b) example of the cluster-to-cluster simulation benchmark.26

List of Tables

Table 1: CM instructions definition.....	15
Table 2: Gem5-X Full System Mode Parameters.....	18
Table 3 System Energy and power Numbers	18
Table 4 CM core parameters	19

1 Introduction

The Multiscale simulation work package 5 (WP5) targets at demonstrating the gains provided by architectures embedding wireless interconnect from the work package 3 (WP3), hardware accelerators and computing architectures developed in the work package 4 (WP4), while running data-intensive applications. Such gains can only be achieved and evaluated through co-optimization methodologies considering accelerator design, architecture design and optimization, wireless transceivers design and management, power and thermal modeling, run-time strategies, and application mapping strategies.

Deliverable 5.1 wraps-up the works that have been done during Task 5.1 (T5.1) and the beginning of Task 5.2 (T5.2) by EPFL. Its goal is to release a beta version of the simulator that will be used in WiPLASH.

The main contributions of this deliverable are:

- Introduction of the RISC-V Instruction Set Architecture (ISA) in the gem5 simulator.
- Introduction of an In-Memory Computing architecture developed in WP4 in gem5 and evaluation of performance gains.
- Introduction of WireLess (WL) interconnection networks as part of the gem5 simulator.

2 RISC-V Instruction Set Architecture Support in Gem5

This section presents the status of the developments of RISC-V support in architectural simulators in the community and introduces the different contributions that have been achieved in EPFL toward the support of RISC-V ISA. This section only discusses the motivation and positioning of the work, technical details can be found as an appendix of this document in the gXR5 technical documentation.

The source for the gXR5 simulator is available for the consortium in the WiPLASH git repository (<https://github.com/orgs/wiplashproject/>) and will be made available for the community after publication in a peer-reviewed venue. Reviewers can request access to the WiPLASH project repository to the project coordinator.

2.1 State of the art and limitations

While RISC-V enjoys good support from emulators such as QEMU [1], functional simulators such as SPIKE [2] and RTL simulation [3], support from architectural simulators is not developed particularly for Linux capable systems. RTL simulations enable extremely precise simulations as every logic gate is being simulated, however, extracting statistics out of simulations of programs running on the hardware may lead to months of simulations even with the most powerful servers in the market forbidding such a solution. On the other hand, emulators and functional simulators only model the execution of instructions to validate the functional results and thereby enable quick load and run time on top of an operating system. However, such a simulation framework does not allow the designer to extract architectural metrics, such as the Instructions per Cycle, the specific units of the pipeline used, the intensity on the cache (misses, hits), and the memory subsystem, and to perform a detailed performance assessment.

The middle ground between RTL and functional simulators, and the focus of this work, are system-level simulators. While not as fast to load and run programs as functional simulators, system-level simulators can represent major hardware components and interconnects as high-level software models with timing information, leading to the attainment of functionally accurate results as well as reasonable hardware performance statistics in significantly less time than RTL simulators. With extensions such as McPat, power, and energy data can also be asserted by the generated performance statistics [4]. While not being able to boast the same level of precision offered by RTL simulators, the result is the ability to rapidly prototype and redesign hardware with reasonable insight into performance ramifications, thus decreasing the time to market of a hardware product. Within WiPLASH, we focus on two system-level simulators, namely the PULP platform virtual platform [3] and the gem5 simulator [5]. The first being more adapted to low-power architectures and specific to RISC-V-based PULP processors, while gem5 provides a larger set of supported ISA (e.g. ARM, X86), enables full software stack simulation, thereby targeting it for higher performance architectural exploration testcases.

In this context, in 2019, the Embedded System Laboratory from EPFL published gem5-X [6], an extension to gem5 which introduces architectural extensions such as in-cache computing and support for High Bandwidth Memories (HBM). It also integrates a methodology for optimizing power and performance in manycore systems. gem5-X was validated against real hardware and provides an error below 4%. However, one major limitation of gem5-X is that it only supports x86 and ARM systems. The goal of this section is to describe the development of a new version of gem5-X, called gXR5,

which extends gem5-X by enabling full-system support for RISC-V single-core architectures running on top of a Linux kernel.

2.2 Summary of the main contributions of EPFL in gXR5

The base target execution environment for this work is a combination software and hardware execution environment with both an Application Binary Interface (ABI) and a Supervisor Binary Interface (SBI). In other words, **our main contribution with gem5-eXtensions for RISC-V (gXR5) is creating the first Linux-capable Full-System (FS) mode system in gem5**. The work leading up to this is as follows:

- We extend the FS mode configuration in gem5 for Linux-capable RISC-V systems.
- We implement instructions from the RISC-V privileged specification and extend or verify instructions implemented in prior work.
- We implement the missing Zifencei extension from the unprivileged specification.
- We extend CSR implementations from prior work.
- We develop a RISC-V compliant MMU capable of processing virtual memory, checking PMP/PMA, and interfacing a page table walker.
- We implement RISC-V ISA devices, including a PLIC (platform-level interrupt controller) and CLINT (core-local interrupter).
- We develop and configure a gem5-compatible bootloader, Linux kernel, device tree, and buildroot image for storage.

With these contributions combined, we are able to demonstrate running programs on top of the Linux kernel, on top of a disk image.

3 Matrix-Vector-Multiplication RRAM Acceleration in Gem5

This section describes the integration of an in-memory computing accelerator developed by IBM in the gem5-X simulator from EPFL and presents performance and energy results extracted from full system simulations of an AI workload. As this work has been carried out in parallel with the development of the gXR5 simulator (i.e., work presented in section 2), the gem5-X simulator configured to run regular ARM architectures has been used for this work. However, it should be noted that an analogous methodology could be considered to enable RRAM acceleration in gXR5.

3.1 State of the art

While Resistive Random Access Memories (RRAM) have been an exciting topic for the entire scientific community for almost 10 years now, besides being considered as a “simple” non-volatile memory (for conventional non-volatile memories replacement), many research institutes try to take advantage of its electrical properties to perform computation with it. RRAM are non-volatile memories and their memory state is held by a resistance value. The current flowing through a RRAM device can be determined from the voltage across its terminals V divided by the device electrical resistance R . By organizing X number of these devices in a one-dimensional array and applying X voltages values (V_x) across X devices holding X different resistive values (R_x), the current measured on the common grounded node is $i = \sum_x \left(V_x * \frac{1}{R_x} \right)$. Such operation is basically a dot-product operation between the different elements of V_x and $1/R_x$ and can be then chained to achieve a convolution. As dot-products and convolutions are widely used in AI workloads, this highlights the usefulness of this type of accelerator for AI workloads.

One of the most advanced and industrial-ready implementations of this concept have been developed by IBM [7] and uses Phase Change Memory (PCM), a RRAM technology in which the non-volatile resistance state is controlled by the structure of a chalcogenide material. Depending on whether the material is crystalline or amorphous, its resistance varies. Programming phases are achieved by finely tuned current pulses (time and current) through the PCM cells, in order to precisely control the volume being amorphized, i.e., the resistance state. By this technique, 8-bit values can be stored in each of the PCM cells reliably.

The Computational Memory (CM) core considered in the rest of this section contains, besides the memory array itself, peripheral circuits to achieve programming operations, row selection, current reading and translation (Analog to Digital Converters – ADC), self timed pulse generations, voltage references and registers for input and outputs.

3.2 System-Level Architecture

3.2.1 System Integration of the IMC module

We implement the CM core within the execution (EX) stage of a classic 5-stage in-order pipelined processor, as shown in Figure 1. This concept is similar to, and takes inspiration from, [8], but the inputs and outputs are not binary and therefore we must adapt the interface to process 8-bit inputs and outputs. While certain operations can cause the CPU to stall in case of a long computation (such as when the CM core performs a MAC operation), the benefit of this configuration, on top of the fact that the CPU no longer needs to go out to lower levels of the memory hierarchy for critical data,

is that the latency of memory transactions between the CM core and rest of the CPU is on the order of nanoseconds. We consider conservative (worst-case) estimates for energy and time [7] [9] [10] [11] for the CM core.

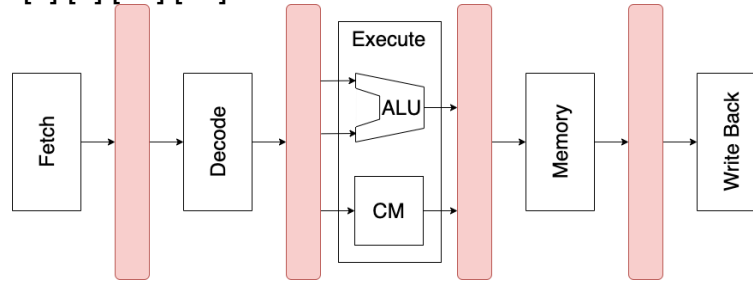


Figure 1: A classic 5-stage in-order pipelined CPU, separated by stage registers, with instruction fetch, instruction decode, execute, memory, and write back stages. The CM core exists alongside the arithmetic logic unit (ALU) within the EX stage of the CPU.

3.2.2 New IMC-Core Interface Instructions

To use the CM core, we define four additional instructions for the ARMv8 (64-bit) ISA: CM_QUEUE, CM_DEQUEUE, CM_PROCESS, and CM_GETSCALE. All instructions are multiple-register operations, referring to three input registers (Rm, Ra, Rn) and one output register (Rd). These instructions are exemplified in Table 1. Certain registers remain unused or reserved for future work, depending on the instruction (marked with 'X' in the table). The instruction bit width is 32, with the register fields having a bit width of 5. We next describe the four instructions briefly.

Table 1: CM instructions definition

Op	OpCode	Rm	Ra	Rn	Rd
CM_QUEUE	0010 0001 000	Rm	X	Rn OR 1	X
CM_DEQUEUE	0010 0001 000	Rm	X	Rn AND !1	Rd
CM_PROCESS	0000 0001 000	X	X	X	Rd
CM_GETSCALE	0100 0001 000	X	X	X	Rd

The CM_QUEUE instruction takes register Rn, packed with 8-bit input values, and places them into the input memory of the CM core. Register Rm holds queuing metadata, including the index of the input memory from where to place the packed 8-bit inputs and the number of valid 8-bit inputs in register Rn. The least significant bit of register Rm is hard-coded to 1 to signify that this is a queuing operation.

The CM_DEQUEUE instruction uses register Rn to refer to which portion of the output memory values shall be fetched from, and pack up to 8 sequential 8-bit output values into register Rd. Just like with CM_QUEUE, register Rm holds dequeuing metadata, including the number of elements being dequeued. CM_DEQUEUE shares its opcode with CM_QUEUE, so the least significant bit of Rm hard-coded to 0 to signify that this is a dequeuing operation.

The CM_PROCESS instruction performs the parallel MAC operation using the input vector held in the input memory of the CM core and its stored weights and biases. Once the MAC operation is complete, the CM core will store the output (after ADC conversion) in the CM core output memory. Register Rd can be used to indicate operation success, and any of the other unused registers can be used to refer to a specific CM core in a multi-CM core-enabled system.

The CM_GETSCALE instruction reads a pre-written 32-bit value from the output memory of the CM core to the register Rd. This value is used to scale the result of the

CM_PROCESS (the 8-bit ADC output), into the proper range for the subsequent calculations.

The op latency of the instructions CM_QUEUE, CM_DEQUEUE, CM_PROCESS, and CM_GETSCALE, are constant, based on hardware characterization [7] [9] [10] [11] and considered to be 1ns, 1ns, 100ns, and 1ns, respectively.

3.2.3 Using the CM Interface instructions in a Recurrent Neural Network test case

This subsection describes how one may use the custom CM instructions in the context of running an RNN LSTM network inference with a single cell layer and a dense layer. In this context, we neglect CM core programming operation as we consider that weights values are already programmed. We assume there to be two CM cores in the system, set up as pictured in Figure 2 [12] [13] [14]. We also assume a vectorization level of V , V referring to the maximum number of elements we can "queue" or "dequeue" into the CM core using a single custom instruction. We consider a vectorization level that can range from 1 to 8 64-bit general-purpose registers for both the input and output of the CM core.

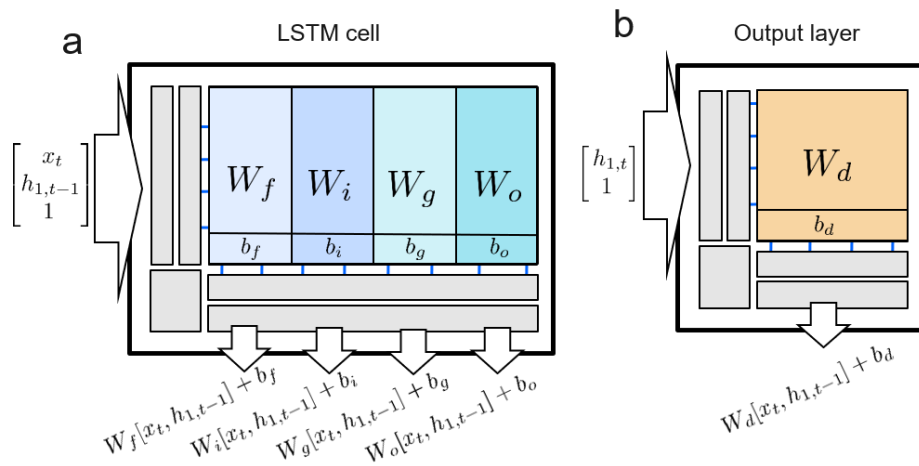


Figure 2: Weight storage in the CM core for an LSTM composed of a single cell and a dense layer. The last row of the array is reserved for the biases.

The rest of this section describes how to use the CM core when it contains the hidden layers of the LSTM RNN we have considered as a test case. The first operation to achieve is to concatenate x , $h_{1,t-1}$, and 1 (the input for the biases) into a single input vector for the CM core. Once concatenated, we save the absolute maximum of the array and cast the input vector to `int8_t`, using the absolute maximum to scale the input. Then, we split the input vector according to the vectorization level.

To prepare the queue instruction call, for every Vectorization 8-bit values in the input vector, we pack those items into a single variable, such that the first value resides in bits 0 to 7, the second value in bits 8 to 15, and so on. Additionally, we set up the metadata register to indicate the number of elements that will be queued into the CM core and at which starting index.

Referring to the queuing and metadata registers, we call the CM_QUEUE instruction. This will place the 8-bit input values stored in the queuing register into the input memory of the CM core. We then repeat this process for all values for the input vector. Once completed, we call CM_PROCESS to perform the MAC operation. This will

perform the dot products of the input vector with W_f , W_i , W_g , and W_o , in parallel. The output of CM_PROCESS will be stored in the CM core output memory.

To retrieve the output vectors from the CM core, we use the CM_DEQUEUE instruction. This requires setting up a metadata register at a similar value as the one used for the CM_QUEUE instruction, where it is specified how many values are being dequeued from which starting index in the output memory array. For every vectorization output values, we call CM_DEQUEUE. The result of this instruction call will be up to 8-bit output values packed into a single register (first value in bits 0 to 7, the second value in bits 8 to 15, and so on).

Once all of the values are dequeued, they will need to be scaled up with the proper scaling factor, to both accommodate for the input scaling (for the DACs) and output scaling (for the ADCs). To fetch the latter, call the CM_GETSCALE instruction, which will return a 32-bit floating-point number. We then perform the scaling by multiplying the output with the aforementioned scaling factors.

At this point, the outputs of the dot products are ready for digital operations to be performed in the CPU (sigmoid and tanh activation functions, and element-wise operations) to finish the LSTM cell portion of the inference. The same process described above then occurs with the dense layer input, MAC, and output, but with a different input vector and the instruction metadata set up to refer to the secondary CM. A softmax is applied at the output of this layer.

3.3 Experimental Setup and Results

3.3.1 Experimental setup

Simulation platform:

To evaluate the CM core under investigation in this project, we implement it within the gem5-X architectural simulator [6] previously introduced in this document in Section 2. Gem5-X is validated using the ARM Juno platform, and thus offers performance statistics representative of a real ARM system. The CM core is then implemented in gem5-X within the ARMv8 ISA templates and custom instructions. It is important to note that, as RISC-V support introduced in Section 2 is under development, thereby, in this part of the work, we focused on ARMv8 ISA knowing that the proposed methodology would be analogous in a RISC-V-based system. From the gem5-X simulator perspective, CM core custom instruction latencies are implemented as op latencies, as a function of the CPU frequency. Instructions that are dependant on the processor frequency are thereby straightforwardly covered by defining their latency as a multiple of the CPU frequency. For constant time instructions (i.e., CM core additional instructions), on the other hand, we tune, for each CPU frequency the multiplier to keep the instruction time constant. In this way, we ensure that the CM core instruction match with the latencies described in 3.2.2.

All of the experiments presented here run in gem5-X FS mode, which simulates the full computer architecture stack, including user-space programs running on top of Linux 5.5 and a Ubuntu LTS 16.04 disk image. The simulated full system includes models for the CPU, caches, memory, interconnects, interrupt controller, and IO. Additionally, gem5 includes support for checking-pointing and statistics reset/dump, in order to isolate the region of interest (ROI) of the LSTM test bench [15]. We reset statistics before an inference starts (after loading the system, operating system, and application binary) and dump them when the inference finishes (before the application exits).

For every experiment, we run two different LSTM test benches programmed in C++ using the C++17 standard and compiled using the aarch64-linux-gnu-g++ compiler. One LSTM test bench represents a 'typical' LSTM implementation that does not use computational memory. The other LSTM test bench is implemented using the CM core via the CM core interface instructions. The instructions can be compiled by placing them in wrappers for the built-in `__asm` method.

Experimental Parameters and system configuration:

Table 2 presents the system parameters considered in this work. We select three core frequencies and three different cache sizes commonly available on ARM-based chips.

Table 2: Gem5-X Full System Mode Parameters

CPU Core Model	Minor (In-Order) CPU
ISA	ARMv8 (AArch64)
CPU Core Frequencies	0.8Ghz, 1.3Ghz, 2.3Ghz
L1 Data/Instruction Cache Size	64kB
L2 (Last Level) Cache Size	128kB, 256kB, 512kB
DDR4 Model	8GB @ 2400Mhz

We select a variety of LSTM parameters to look at LSTM network sizes ranging from very small to moderately large. The parameters are the size of the input x , the size of a single matrix from the hidden layer h , the size of the output, and dense layer y , and finally the number of inference time steps. The number of time ticks is set to a constant value (10) because preliminary experiments showed a very low standard deviation (0.0014) concerning the performance statistics for each time tick.

Full System Power Model:

Our power model is shown in Table 3. Our core and cache power model is based on a 28nm bulk system with an ARM A57 core [6], while our DRAM power model is based on [16]. The core and cache power model is comprised of active and WFM (wait for memory) CPU core energy per cycle, in addition to energy and power for the last-level cache (LLC). The full core power model includes a wide range of frequencies between 0.1GHz and 2.3GHz [6]. To calculate the total energy of the system using the power model, we use the gem5-X statistics. The generated statistics include total CPU cycles, simulated time, total LLC read/write hits, LLC accesses, total DRAM read and write accesses, etc. The full system energy is then calculated as the sum of the energies for the core, LLC, and DRAM components.

Table 3 System Energy and power Numbers

Metric	0.8Ghz Core	1.2Ghz Core	2.3Ghz Core
WFM Core Energy (pJ/Cycle)	46.04	125.52	1234.11
Active Core Energy (pJ/Cycle)	60.92	166.06	1632.71
Mem Controller + IO power (W)	3.03	5.06	5.82
LLC leakage (mW/256kB)	271.62	604.43	874.08
LLC read energy (pJ/Byte)	1.81	2.76	8.144
LLC Write Energy (pJ/Byte)	1.63	2.49	7.32
DRAM Energy (pJ/Access)	120		

For the version of the LSTM test bench that uses the CM core, we also add the energy of the CM core to the total energy calculated for the full system, using the parameters

from Table 4. We consider that an analog MAC operation inside the CM core takes 100ns based on characterization provided by IBM in work package 4. Further reading on this topic can be done here [7].

Table 4 CM core parameters

CM Size	M x N
CM_PROCESS instruction time	100ns
All other CM instruction time	1ns
Output Bandwidth	N x 8bits / 100ns
MAC Operation Energy	50fJ
Performance	2 x M x N / 100ns OPS
CM idle power	2.55mW

3.3.2 Performance and energy results

Figure 3 shows the performance (time), Instructions per nanosecond (IPS), memory intensity metric (Last Level Cache Misses per Kilo Instruction - LLCMPKI), and energy for all the considered hardware (core frequencies from 0.8 to 2.3Ghz and L2 cache size from 128 to 512kB) and LSTM input size (from 9kB to 9.8MB), with and without the CM core.

While the CPU power is made worse (by about 3% for the full system and 10% for the CPU alone) with the addition of the CM core, the total energy it uses remains far lower because of the achieved performance gains. The reason for this significant speedup can be seen in the memory intensity metric (LLCMPKI): per inference, the working set of the digital testbench includes the entirety of the weights matrices and biases, as well as numerous temporary accumulators that have little locality from one inference to the next. In contrast, because the weights and biases are never fetched or stored from memory in the analog version of the testbench, its working set includes only the vectors one queues and dequeues from the CM core (once per LSTM cell and once per dense layer cell). As a result, the system using the CM core utilizes the local caches more (as evidenced from the higher LLC energy in the analog test bench) and it does not have to wait for memory as often.

Relative to the digital test bench, the analog test bench gives us a time speedup of 20.2x, an IPS speedup of 30.6x, and an energy improvement of 19.0x. The mismatch in time and IPS speedups come from very granular control of the CM core in C++ leads to worse compiler optimizations (e.g., the addition of more instructions), and thus the analog test bench runs approximately 20% more instructions than the digital test bench.

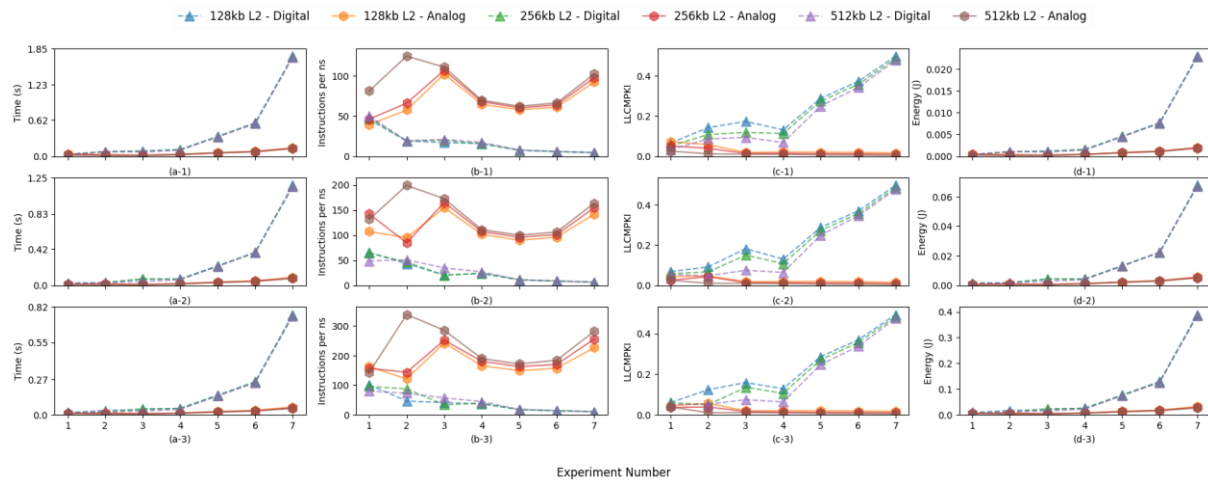


Figure 3: All performance and energy results for the system architecture exploration study. Rows (1), (2), and (3) show the results for system configurations using the 0.8GHz, 1.3GHz, and 2.3GHz core frequencies, respectively. Columns (a), (b), (c), and (d) show the simulated time (s), instructions per ns (IPnS), memory intensity (LLCMPKI), and energy (J) results, respectively. The experiment number refers to the numbered LSTM networks (larger experiment number represents larger LSTM).

The trends in performance and energy across all system configurations (core frequencies and LLC sizes) remains consistent. While there are energy and performance improvements across all LSTM experiments and system configurations, they only jump significantly once the LSTM is large enough. Experiments 1 through 3, the smallest LSTM networks, have a very small working set and so the largest improvements in performance and energy observed come from the system configurations with the larger LLCs. In this case, 4x improvement in performance and IPS, and 4.4x in energy consumption is achieved. From experiment 4 onwards, however, the trends converge regardless of LLC size and CPU core frequency, showing very significant performance and energy gains. In the largest experiment, the analog over digital test bench time and IPS speedup was smallest in the 0.8GHz core with 128kB LLC (12x time speedup, 20x IPS speedup) and largest in the 2.3GHz core with 512kB LLC (16x time speedup, 28x IPS speedup).

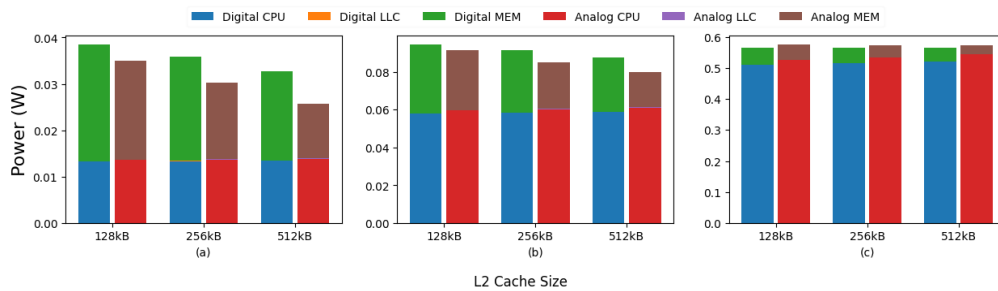


Figure 4: Full system average power distribution for the architectural exploration study. (a), (b), and (c) show the average power distribution for system configurations tested with 0.8GHz, 1.3GHz, and 2.3GHz cores, respectively.

Figure 4 shows the evolution of the average power distribution in different hardware configurations for both versions of the architecture with and without the CM core. The reason why the slowest CPU core shows the greatest energy improvement is because the CPU power is naturally very low and the DRAM power contributes a significant portion of the total power. By decreasing the memory intensity of the LSTM application using the CM core, we cut our reliance on the DRAM and therefore decrease its power contribution by 10-25% depending on the LLC size.

3.4 Conclusion and future works

By implementing an analog CM core and interface instructions in the CPU to queue inputs, process MAC operations, and dequeue outputs, we decrease the size of the effective working set of the LSTM network inference to exclude the weights and biases, thus decreasing the memory intensity of an inference and reaping significant time and energy improvements over CPU-only architectures. Using the CM core, we explored the architecture design space and showed significant speedup (20.2x) and energy gains (19x) across a variety of system configurations and LSTM network sizes.

This section demonstrated that IBM CM core can be successfully integrated inside a complete system and forms a strong baseline for future exploration works. Integrated AI accelerators are fundamental elements of future systems and, for the CM core, will be extremely useful to leverage wireless interconnect as it drastically cuts the communication needs of the system. On the other hand, since accelerators such as the CM core cannot store entire AI models, their integration, as described in this deliverable, will foster the exploration of novel communication strategies as described later in Sections 4.3 and 4.4.

As future work, we will explore the integration of the CM core as an element of an heterogeneous multi-core architecture. Then, we will integrate it with the gXR5 version of gem5 to evaluate its performance and energy implications within a system leveraging wireless interconnects.

The sources and documentation are accessible from the main Gem5-X repository on the ESL-EPFL website:

<https://www.epfl.ch/labs/esl/research/2d-3d-system-on-chip/gem5-x/>

This work is currently being reviewed as a submission to the IEEE/ACM Design Automation and Test in Europe (DATE) conference 2021 and thereby will be released online after publication.

4 Wireless Interconnect Inside the Simulation Framework

The introduction of wireless interconnects inside computing architectures opens many questions from the energy, area, and performance perspectives. Conventional state-of-the-art integrated antennas and transceiver circuits exhibit large area overhead, low bandwidths, and high energy consumption [17] [18] [19] [20]. Therefore, they are not viable candidates for replacing conventional communications buses and networks-on-chip for intra-chip and chip-to-chip communication.

The introduction of graphene-based antennas, as proposed in WiPLASH, is expected to bring in massive gains: up to 10X in execution speed for some applications (e.g. Artificial Intelligence processing). Thereby, in this part of the WiPLASH project, taking as input developments from the other work packages, we intend to set up the working baseline enabling us to explore how systems could take advantage of this new element. As this deliverable aims at providing a first structure to support for WL interconnect networks, we keep simulation for later in T5.2 and we present in this section a general overview of wired versus WL interconnect networks, previously achieved works within the consortium, potential architectures targets and finally, implemented architectures templates within the gem5-X simulator from EPFL.

4.1 Wired versus wireless interconnect

The major interest of wireless interconnects, in opposition to wire-based interconnects, resides in the plasticity of the communication network. By not being tied to physical considerations, it enables the design of less intrusive and more modular communication networks. Communication paths and bandwidths can be efficiently re-allocated depending on the needs of the application without impact on the latency at runtime. Finally, both intra-chip and inter-chip can be targeted. In WiPLASH, we intend to address both on-chip and off-chip wireless interconnection in the TeraHertz (THz) band enabled by graphene antennas.

Figure 5-a presents the organization of a standard 2D mesh NoC for a tile-based homogeneous multiprocessor architecture. In this case, core-to-core communication delay and energy costs are NoC-topology-dependent, and depend on the physical distance between the communicating elements. In a wireless NoC, as shown in Figure 5-b, communication costs become constant if the different elements are in range. In other words, communication between any processing element in range could be achieved in parallel without any congestion as long as there are enough frequency, time, or space (beam) channels available.

In modern chiplet-on-interposer architectures, an NoC is usually integrated into the interposer to enable communication between the different elements of the MPSoC [21] [22]. While such architectures are extremely efficient from a process cost and yield perspective, due to long distances between the different elements, NoC performances are drastically reduced, forcing designers to look for non-conventional technology or architectural solutions such as silicon bridges [23], silicon photonics [24] or asynchronous NoCs [25]. In this context, an integrated wireless interconnect is expected to play a game-changer role thanks to all its aforementioned advantages.

From an architectural simulation viewpoint, as we explain in Section 4.4, many architectural concepts can be reused and parametrized to introduce wireless interconnect inside the system.

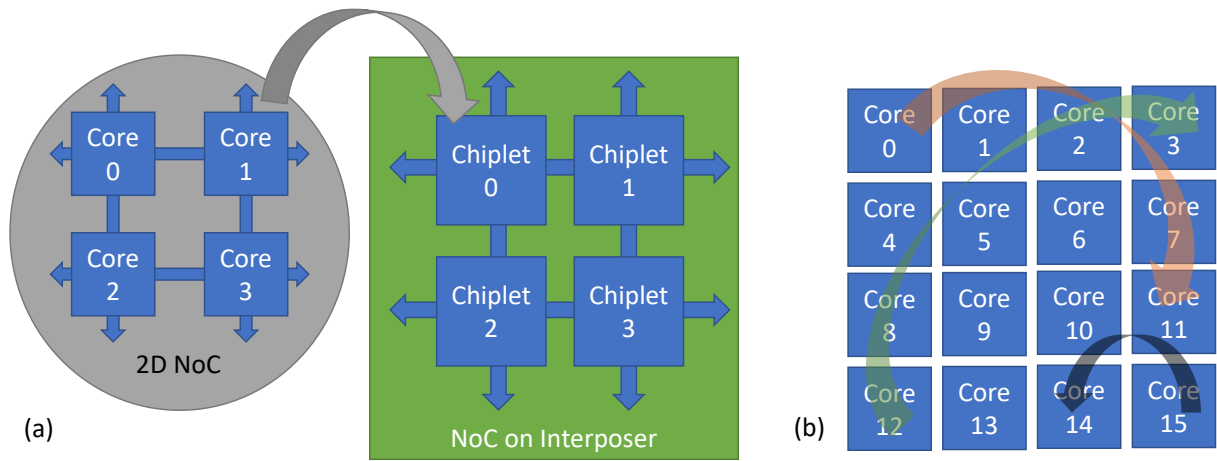


Figure 5: (a) illustration of 2.5D chiplet on interposer computing unit. Each chiplet contains a 2D NoC for local communication. Chiplet-to-chiplet communication is managed by a NoC on interposer (as described in [25]) with limited bandwidth and higher risks of congestion. (b) illustration of a wireless NoC in which any computing unit could communicate with any computing unit regardless of their physical organization.

4.2 Previous Works on Wireless Interconnect

Before the start of this project, UPC has been exploring the track of wireless interconnect-based NoCs. In this subsection, we summarize previous architectural exploration works that have been led by consortium members.

WiSync [26]: One of the main interests of wireless interconnect is the ability to broadcast. In other words, once a transceiver starts streaming, all the receivers in range can receive at the same time the data without any additional cost. In a massive many-core architecture (100+), in complement with a standard wired NoC, each core is tied to a Broadcast Memory (BM) that is then synchronized through the wireless stream interconnect to all the other cores in a really short time (less than 10 clock cycles). Even with a moderate bandwidth (20Gb/s), WiSync achieves an average 1.41x performance improvement compared to accessing synchronization variables through the regular memory hierarchy.

Replica [27]: In Replica, the authors demonstrated that shared BM could be used not only for the synchronization of locks and barriers, but of any useful data. Thereby, in the works, the authors exploited their application knowledge to put latency-critical data in the BM. In this way, they demonstrated up to 1.86x performance improvement versus a conventional architecture.

Beyond these works, that target specific applications and architectures, in WiPLASH, we intend to go one step further and address more general test cases targeting AI and Deep Learning applications, as well as more traditional benchmarks in HPC and cloud computing systems (e.g., SPEC Benchmarks and CloudSuite).

4.3 Application and Architecture targets for Wireless Interconnect

Today, artificial intelligence workloads are becoming pervasive. These applications are required, on the one hand, to run on ultra-low-power autonomous devices, and, on the other hand, on high-performance computing platforms. When running AI workloads, both target power budgets are dominated by data movement and storage associated energy costs. AI workloads and particularly artificial Neural Networks such as CNNs or RNNs require a huge amount of data to operate [28]. Most of the computation being

multiply and accumulate, this part can be accelerated like previously described through specific accelerators. However, the model sizes are huge and cannot fit within the local memory of a core nor in the shared memory, highlighting the well-known memory wall problem that we tackle in WiPLASH.

Two main architecture targets for wireless integration can be identified. On-chip communication and chip-to-chip communication. From a system simulation perspective, wireless interconnects could be modeled in many ways and the system organization (on-chip or chip-to-chip communication) can be abstracted to a reduced set of parameters (e.g. latencies, bandwidths, energy per bit) enabling the design exploration to be somehow agnostic from physical constraints and thereby ease the interaction with physical design considerations.

In this context, two major simulation methodologies can be considered, for either (1) a low-power system featuring deterministic programmer-controlled memory management such as PULP-based architectures, or (2) higher performance systems with a non-deterministic cache hierarchy memory management and running with an Operating System (OS). While wireless interconnects could be integrated in both architectures, in this work package we focus on the integration of wireless interconnects in relatively complex high-performance multicore architectures managed by an operating system. Such architectures could be considered in smartphone or server platforms.

On-chip WL communication integration on conventional smartphones or server-like SoC architectures is not straightforward. As architectures are already highly optimized and spatially organized to minimize wired interconnect energy-delay product and maximize BWs, there is not much space for as-is integration of WL transceivers (Figure 6-a presents the floorplans of several recent chips in which physical organization and blocks placement is highly optimized). On the other hand, emerging architectures such as homogeneous and heterogeneous manycore, or manycore accelerators for AI exhibit much more constraints on the communication network. One reason is that such an organization does not enable anymore the minimization of the path to large shared memories, forcing the system to rely on NoCs which may drastically lose their efficiency if a large amount of data is required to be communicated simultaneously. Figure 6-b exemplifies manycore architectures and highlights the need for efficient communication protocols to handle them.

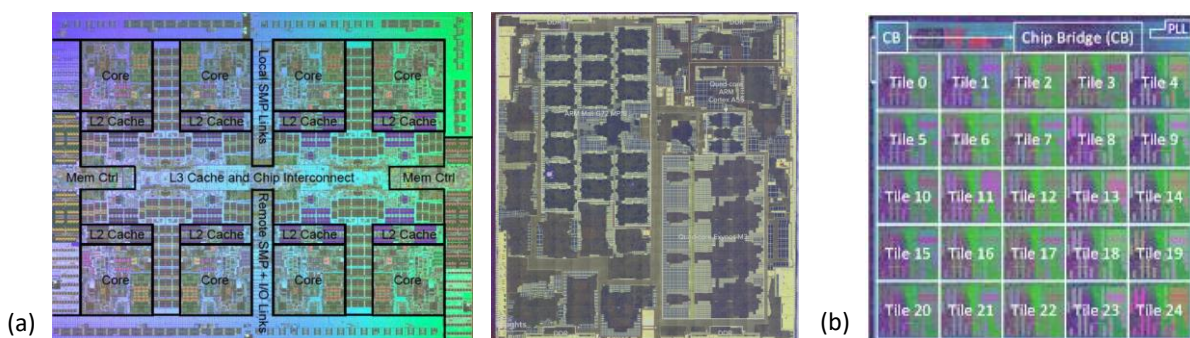


Figure 6: (a) floorplan of recent 14nm Samsung Exynos-type SoC [29] and 22nm IBM Power-8 processor [30]. (b) Floorplan of a OpenPiton tile-based 25cores architecture using a wired 5x5 2D mesh NoC topology [31].

Chip-to-chip WL communication integration appears as an appealing solution for 2.5D chiplet integration on interposer as it would enable the communication path between any element of the chip to bypass the TSVs and the supporting interposer chip. In this context, communication between processing units that may be part of

different chiplets would be drastically boosted. Additionally, communication between processing units and DRAM chiplets would also strongly benefit from both high bandwidths, reconfigurable channels, and streaming capabilities.

4.4 Modeling parameters and Implementation of WL interconnect in gem5

To address the aforementioned considerations and ideas, we have been envisioning preliminary test benches that we currently explore along Task 5.2. While this task addresses specifically on-chip interconnect, we first consider simple topologies containing a few processing units. This section presents several possible implementations of WL interconnect networks and three benchmarks. While only the two first versions are already implemented within the gem5 simulator, we describe the third one in this document as it is the follow-up experiment.

Processor-DRAM link through WL interconnect. The naïve approach while simulating a system without considering NoC constraints is to simply consider that a wireless link replaces the communication between the processor and its main DRAM memory. By tuning the latency and available bandwidth as part of a new DRAM model, it is possible to emulate a wireless link. In this case, we plan to evaluate the performance and energy gains enabled by such modification in the latency and bandwidth, and explore trade-offs involving physical design considerations (related to antennas and transceivers energy, area and position for e.g.).

Core-to-core communication benchmark: As mentioned before, CNN models may be too large to fit in the local memory of processing elements. This mismatch leads to reduced performances as data must be sent in and out along the code execution. We propose to explore the scenario described in Figure 7-a. Two processing elements containing each one CPU, one AI accelerator (we propose to consider the IMC accelerator developed in the field of WiPLASH – cf Section 3), and a memory hierarchy communicating together through a wireless link. The wireless link is modeled by a First-in-First-Out (FIFO) memory with a tunable latency within gem5. We will explore the communication intensity between these two elements, and its impact on performances and energy consumption, targeting small benchmarks at first (split matrix multiplication) and more complex ones later on (including entire CNN layers). By tuning the latency and width of the FIFO, the size of the different memory elements, types, and frequency of cores, we will perform a full system evaluation of such an approach compared to a regular wired interconnect under different architecture assumption (e.g. varying the distances among cores). As for the previous experiment, we will explore trade-offs considering physical limitations of WL interconnect.

The two aforementioned implementations are currently integrated within the gem5 simulator is available for the consortium in the WiPLASH git repository (<https://github.com/orgs/wiplashproject/>) and will be made available for the community after publication in a peer-reviewed venue. Reviewers can request access to the WiPLASH project repository to the project coordinator.

Cluster-to-cluster communication benchmark: As for the previous benchmark, we will consider two clusters of several cores working together. The proposed simulation benchmark is presented in Figure 7-b and will involve the same exploration methodology as the core-to-core benchmark. By accurately modeling the behavior of a cluster of cores, it will be possible to scale it to massive many-core architectures that at present cannot be simulated in a sustainable amount of time [6].

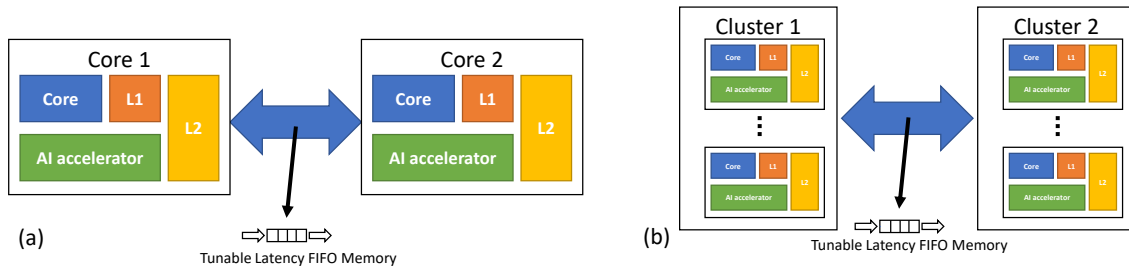


Figure 7 : (a) example of the core-to-core communication simulation benchmark implemented in WiPLASH. (b) example of the cluster-to-cluster simulation benchmark.

4.5 Conclusions and Perspectives

In this section of the deliverable, we presented an overview of wireless interconnect as part of future MPSoC architectures. In that sense, we first introduced the state of the art of WL interconnect networks compared to wired NoCs. We presented previous works proposed by other consortium members (UPM) and introduced different WL network topologies as well as their implementation within the simulation framework proposed by EPFL. This section only focused on implementation, as at the moment we only validated the functionality of the proposed topologies (variable latency in CPU-DRAM link, and core-to-core communication) and implemented them within the gem5-X EPFL simulator. Next steps include the implementation of cluster-to-cluster simulation and the evaluation of WL interconnect as part of the architecture. These experiment will foster and strengthen collaboration with partners from WP3 as we will be able to include in the simulation physical consideration on wireless transceivers bandwidth, latencies and energy consumption.

5 Conclusions and Perspectives

This section concludes the work achieved towards D5.2 along T5.1 (from M1 to M12) and T5.2 (from M9 to M12), and introduces future works perspectives for the ongoing (T5.2) and coming tasks (T5.3 from M15). The main achievements of this deliverable are the setup of a simulation infrastructure for future WL-enabled MPSoC running AI workloads and are summarized as follow :

- Extension of the gem5 simulator to support for RISC-V ISA.
- Extension of the gem5 simulator to support the IMC accelerator from WP4.
- Extension of the gem5 simulator to support for basic WL interconnect networks.

According to the objectives defined for the 12 first months, more than expected has been achieved. While the global simulation infrastructure has been settled-up through the extension of the gem5 simulator to cover for RISC-V ISA, AI IMC accelerator and basis of WL interconnect networks, we also extracted application gains on a full AI application running on top of a full software stack and submitted our findings from section 3 to the IEEE Design Automation and Test in Europe (DATE) conference 2021.

The current progress in the project makes us ready to target the next project objectives (i.e., D5.2 at M24). In that sense, we have in our hands all the tools to explore WL-enabled die-level architectures. While in a short term perspective, we will explore the impact of WL interconnect in the different benchmarks proposed section 4.4, in a long term perspective, we will explore techniques such as task mapping strategies and on-the-fly hardware reconfiguration, for performance, energy efficiency optimization or thermal management.

Bibliography

- [1] F. Bellard, "QEMU, a fast and portable dynamic translator," *USENIX Annual Technical Conference*, 2005.
- [2] A. Waterman and Y. Lee, "Spike, a RISC-V ISA Simulator," <https://github.com/riscv/riscv-isa-sim>.
- [3] "Plup platform website," [Online]. Available: <https://pulp-platform.org/>.
- [4] S. Li, J. H. Ahn, R. Strong, J. Brockman, D. Tullsen and a. N. Jouppi, "McPAT: an integrated power, area, and timing modeling framework for multicore and manycore architectures," *Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture*, pp. 469-480, 2009.
- [5] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. R. Hower, T. Krishna, S. Sardashti, R. Sen, K. Sewell, M. Shoaib, N. Vaish, M. D. Hill and D. A. Wood, "The gem5 simulator," *ACM SIGARCH Computer Architecture*, 2011.
- [6] Y. M. Qureshi, W. A. Simon, M. Zapater, D. Atienza and K. Olcoz, "Gem5-X: A Gem5-Based System Level Simulation Framework to Optimize Many-Core Platforms," *Spring Simulation Conference (SpringSim)*, 2019.
- [7] A. Sebastian, I. Boybat, M. Dazzi, I. Giannopoulos, V. Jonnalagadda, V. Joshi, G. Karunaratne, B. Kersting, R. Khaddam-Aljameh, S. R. Nandakumar, A. Petropoulos, C. Piveteau, T. Antonakopoulos, B. Rajendran, M. Gallo and E. Eleftheriou, "Computational memory-based inference and training of deep neural networks," *Symposium on VLSI Technology*, 2019.
- [8] J. Vieira, E. Giacomini and Y. Qureshi et al., "A Product Engine for Energy-Efficient Execution of Binary Neural Networks Using Resistive Memories," *IFIP/IEEE 27th International Conference on Very Large Scale Integration (VLSI-SoC)*, 2019.
- [9] V. Joshi, M. Le Gallo and S. Haefeli et al., "Accurate deep neural network inference using computational phase-change memory," *Nature Communications*, 2020.
- [10] S. Nandakumar, M. Le Gallo and C. Piveteau et al., "Mixed-precision deep learning based on computational memory," *Frontiers in Neuroscience*, 2020.
- [11] L. Kull, D. Luu, C. Menolfi and M. Braendli et al., "A 10b 1.5 GS/s pipelined-SAR ADC with background second-stage common-mode regulation and offset calibration in 14nm CMOS FinFET," *IEEE International Solid-State Circuit Conference (ISSCC)*, 2017.
- [12] T. Gokmen, M. J. Rasch and W. Haensch, "Training LSTM networks with resistive cross-point devices," *Frontiers in Neuroscience*, 2018.
- [13] C. Li, Z. Wang and M. Rao et al., "Long short-term memory networks in memristor crossbar array," *Nature Machine Intelligence*, 2019.

- [14] S. Nandakumar, M. L. Gallo and C. Piveteau et al., "Mixed-precision deep learning based on computational memories," *Frontiers in Neuroscience*, 2020.
- [15] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, 1997.
- [16] S. Lee, H. Cho and S. e. a. Y. H. al., "Leveraging power-performance relationship of energy-efficient modern dram devices," *IEEE Access*, 2018.
- [17] J. Kim, K. Choi and G. Loh, "Exploiting New Interconnect Technologies in On-Chip Communication," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems (JETCAS)*, 2012.
- [18] S. Laha, S. Kaya, D. W. Matolak, W. Rayess, D. DiTomaso and A. Kodi, "A New Frontier in Ultralow Power Wireless Links: Network-on-Chip and Chip-to-Chip Interconnects," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, 2015.
- [19] S. Deb, A. Ganguly, P. P. Pande, B. Belzer and D. Heo, "Wireless NoC as interconnection backbone for multicore chips: Promises and challenges," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems (JETCAS)*, 2012.
- [20] X. Yu, J. Baylon, P. Wettin, D. Heo, P. P. Pande and S. Mirabbasi, "Architecture and design of multichannel millimeter-wave wireless NoC," *IEEE Design and Test*, 2014.
- [21] S. Bharadwaj, J. Yin, B. Beckmann and T. Krishna, "Kite: A Family of Heterogeneous Interposer Topologies Enabled via Accurate Interconnect Modeling," *IEEE/ACM Design Automation Conference (DAC)*, 2020.
- [22] A. Kannan, N. E. Jerger and G. H. Loh, "Enabling Interposer-based Disintegration of Multi-core Processors," *IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2015.
- [23] D. Greenhill, R. Ho, D. Lewis and H. Schmit et al., "A 14nm 1GHz FPGA with 2.5 D transceiver integration," *IEEE International Solid-State Circuit Conference (ISSCC)*, 2017.
- [24] P. Fotouhi, S. Werner, J. Lowe-Power and S. J. B. Yoo, "Enabling scalable chiplet-based uniform memory architectures with silicon photonics," *ACM International Symposium on Memory Systems (MEMSYS)*, 2019.
- [25] P. Vivet, E. Guthmuller, Y. Thonnart and G. Pillonnet et al., "A 220GOPS 96-Core Processor with 6 Chiplets 3D-Stacked on an Active Interposer Offering 0.6ns/mm Latency, 3Tb/s/mm² Inter-Chiplet Interconnects and 156mW/mm² @ 82%-Peak-Efficiency DC-DC Converters," *IEEE International Solid- State Circuits Conference - (ISSCC)*, 2020.
- [26] S. Abadal, A. Cabellos-Aparicio, E. Alarcon and J. Torrellas, "WiSync: An Architecture for Fast Synchronization through On-Chip Wireless Communication," *ACM SIGPLAN*, 2016.
- [27] V. Fernando, A. Franques, S. Abadal, S. Misailovic and J. Torrellas, "Replica: A Wireless Manycore for Communication-Intensive and Approximate Data,"

International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2019.

- [28] S. Bianco, R. Cadene, L. Celona and P. Napoletano, "Benchmark Analysis of Representative Deep Neural Network Architectures," *IEEE Access*, 2018.
- [29] "TechInsights website - smartphone teardown," [Online]. Available: <https://www.techinsights.com/blog/samsung-galaxy-s9-teardown>.
- [30] J. Stuecheli, "HotChips 2013 - IBM presentation," [Online]. Available: https://www.hotchips.org/wp-content/uploads/hc_archives/hc25/HC25.20-Processors1-epub/HC25.26.210-POWER-Studecheli-IBM.pdf.
- [31] M. McKeown, A. Lavrov, M. Shahradeh, P. Jackson, Y. Fu, J. Balkind, T. Nguyen, K. Lim, Y. Zhou and D. Wentzlaff, "Power and energy characterization of an open source 25-cores manycore processor," *IEEE International Symposium on High Performance Computer Architecture (HPCA)*, 2018.



gem-5 eXtensions for RISC-V: Full System Manual

^{*} EMBEDDED SYSTEMS LABORATORY,
SWISS FEDERAL INSTITUTE OF TECHNOLOGY, LAUSANNE (EPFL)

[‡] REDS INSTITUTE
SCHOOL OF ENGINEERING AND MANAGEMENT VAUD (HEIG-VD),
UNIVERSITY OF APPLIED SCIENCES WESTERN SWITZERLAND (HES-SO)

JOSHUA KLEIN^{*}, YASIR QURESHI^{*}, MARINA ZAPATER^{*‡}, AND DAVID ATIENZA^{*}

June 2020



Contents

1	Executive Summary	4
1.1	Abstract	4
1.2	Release Information	4
1.3	Collaboration and Contact Information	4
1.4	Licenses	4
2	Introduction	5
2.1	Motivation	5
2.2	Background	5
2.2.1	The RISC-V Instruction Set Architecture	6
2.2.2	The RISC-V Unprivileged ISA Specification	6
2.2.3	The RISC-V Privileged ISA Specification	6
2.2.4	Introduction to gem5 and gem5-X	7
2.3	Prior Work	7
2.4	Contributions	8
3	Running gem5 Full System Mode with RISC-V and Linux	9
3.1	Necessary Files	9
3.1.1	Device Tree	9
3.1.2	Bootloader	9
3.1.3	Linux Kernel	10
3.1.4	Disk Image	10
3.2	Recommended Script Options	11
3.3	Quick-Start Guide	11
3.3.1	Prerequisites	11
3.3.2	Building the gem5 Binary	11
3.3.3	Setting Up Your Experiment	12
3.3.4	Running Your Experiment	12
4	Full System Configuration Overview	14
4.1	Full System Model and Child Tree	14
4.2	Full System Memory Map	15
4.3	SimpleBoard Platform	15
5	Implementation of Privileged RISC-V ISA Specification	17
5.1	Control and State Registers (CSRs)	17
5.1.1	CSRs versus Other Registers	17
5.1.2	CSRs in gem5-X	17
5.1.3	Interfacing and Accessing CSRs	17
5.2	RISC-V Instructions	18
5.2.1	Instructions in gem5	18
5.2.2	Privileged Specification Instructions	18
5.2.3	Unprivileged Specification Instructions	19
5.3	Fault Handling	19
5.3.1	RISC-V Terminology	19
5.3.2	RISC-V Fault Handling Algorithm	19



5.3.3	Exceptions in gem5	20
5.3.4	Interrupts in gem5	20
5.4	Virtual Memory	20
6	Virtual Memory Subsystem	22
6.1	TLB Implementation	22
6.2	Physical Address Translation	22
6.3	Virtual Address Translation	22
6.4	Page Table Walker Implementation	22
7	SimpleBoard Platform	23
7.1	SimpleBoard Implementation	23
7.2	SimpleBoard Devices and Parameters	23
8	Platform-Level Interrupt Controller (PLIC)	24
8.1	PLIC Implementation	24
8.2	PLIC Registers and Memory Layout	24
8.2.1	PLIC Source Registers	24
8.2.2	PLIC Interrupt Pending Array	24
8.2.3	PLIC Interrupt Enable Registers	25
8.2.4	PLIC Priority Thresholds	25
8.2.5	PLIC Claim/Complete Registers	26
8.3	PLIC Claim/Complete Operations	26
8.4	PLIC postInt and clearInt Methods	26
9	Core-Local Interrupter	27
9.1	CLINT Implementation	27
9.2	CLINT Registers and Memory Layout	27
9.2.1	CLINT Timer Register and the mtime CSR	27
9.2.2	CLINT Time Compare Register	28
9.2.3	CLINT Software Interrupt Pending Register	28
9.3	CLINT Timer Interrupts	28
10	Miscellaneous SimpleBoard SoC Devices	29
10.1	RISC-V PCI Host	29
10.2	UART	29
11	Future Work	30
11.1	Formal Verification and Validation	30
11.2	MPSoC Support	30
12	References and Acknowledgements	31
12.1	Acknowledgements	31
12.2	Links Reference	31
13	Glossary	32
13.1	Acronyms	32
13.2	Terms	33



Bibliography

34



1 Executive Summary

1.1 Abstract

While RISC-V has enjoyed both strong functional simulation support via ISA simulators such as QEMU (2) and Spike (3) and RTL simulation support via Chisel and other HDL simulators (4), it has little support in the realm of full system-level simulators, especially for simulation of Linux-capable systems. This presents a bottleneck in the RISC-V hardware development process because it is difficult to quickly and reliably prototype and verify the performance of hardware designs for complex high-level applications, such as deep learning. To resolve this bottleneck, we present in this technical manual, *gem5-eXtensions for RISC-V*, or **GXRS**: a functioning Linux-capable full system simulator built into the gem5 system-level architectural simulator and gem5-X (5)(6). We extend prior work by implementing the RISC-V privileged specification in gem5 (7). Our contributions include implementing privileged specification instructions and control and state registers (CSRs), support for user and supervisor privilege modes, a RISC-V compliant MMU capable of processing virtual memory, and ISA devices and interrupters, in addition to creating and configuring a gem5-compatible bootloader, device tree, Linux kernel, and disk image file system. With the privileged specification implemented and external components configured, we are able to demonstrate functionally correct execution of basic programs and benchmarks fetched from the disk image and executed on top of the Linux kernel.

1.2 Release Information

Version	Date	Changes
v1.0	June 2020	Initial release.

1.3 Collaboration and Contact Information

The maintainers of this project can be contacted via email at {joshua.klein, yasir.queshi, david.atienza}@epfl.ch and marina.zapater@heig-vd.ch.

Because the scope of this project is very large, we are always interested in potential collaboration efforts to develop new features and bring **GXRS** to gem5 master. For inquiries, source code, and additional information, please contact one of the aforementioned emails.

1.4 Licenses

GXRS is released under the GNU GPL v2.0 license. Please refer to the LICENSE file in the main repository for more details.

2 Introduction

2.1 Motivation

Due to its status as an open-source and free instruction set, the RISC-V ISA has long been a popular candidate for implementing new hardware systems-on-chip (SoCs) and accelerators by striving to be the Linux of hardware (1). Even though it was only introduced in 2014, the open-source nature of RISC-V has led to varying levels of support across the computer architecture stack. Towards the software end, there are multiple Linux ports for various distributions (11) (12), as well as Linux kernel support upstream made possible by the RISC-V GNU toolchain (13). The hardware end is more limited however: there are multiple different kinds of simulators currently available for RISC-V and each represents a significant trade-off.

On one end of the spectrum, there are RTL simulators that typically interpret hardware description languages (HDLs) like Chisel and Verilog. They can provide extremely accurate simulations of hardware interactions, leading to precise performance results with respect to speed, latencies, bandwidth, power, and energy. However, even on the most powerful machines, RTL simulators can take on the order of weeks or even months to run and generate statistics for high-level, complex applications. This can significantly increase the time-to-market of a hardware product.

On the other end of the spectrum, there are functional ISA simulators such as QEMU (2) and spike (3) for RISC-V. These simulators only model the execution of instructions to verify the functional results of a program. While it is not possible to attain precise performance results of the underlying hardware with these simulators, one can very quickly load and run a program on top of an operating system using these simulators. What may take on the order of weeks or months in an RTL simulator can easily run on a functional simulator in minutes.

The middle ground between RTL and functional simulators, and the focus of this work, are system-level simulators. While not as fast to load and run programs as functional simulators, system-level simulators can represent major hardware components and interconnects as high-level software models with timing information, leading to the attainment of functionally accurate results as well as reasonable hardware performance statistics in significantly less time than RTL simulators. With extensions such as McPat, power and energy data can also be asserted by the generated performance statistics (14). While not being able to boast the same level of precision offered by RTL simulators, the end result is the ability to rapidly prototype and redesign hardware with reasonable insight into performance ramifications, thus decreasing the time to market of a hardware product.

Unfortunately, the premier system-level simulator in academia, gem5 (5), only has limited support for RISC-V. Prior work has implemented the unprivileged instruction set (15) as well as limited bare metal full system support for RISC-V (10). The RISC-V privileged ISA specification has not been implemented in gem5 however, and so it is impossible to leverage the benefits of this system-level simulator for complex applications running on top of a Linux system. Therefore the goal of **GXR5** is to implement a Linux-capable full system simulator to allow for rapid design-space exploration of new system architectures for RISC-V.

2.2 Background

In this section we introduce the basic terms and ideas of the RISC-V instruction set, compare the unprivileged and privileged specifications, describe the target RISC-V execution stack, and introduce gem5 and gem5-X.



2.2.1 The RISC-V Instruction Set Architecture

Introduced in 2014, RISC-V is a free and open-source ISA that is built to have a minimalist base instruction set that is highly extensible and can still meet the demands of modern computer systems. The ISA comes in 32-bit, 64-bit, and 128-bit formats, which are denoted as RV32, RV64, and RV128, respectively, and all instructions are 32-bit (with the exception of shorter instructions introduced with the compressed ISA extension). The specification is split into multiple documents, including the base unprivileged specification, privileged specification, external debug specification, trace specification, and compliance framework (1) (9).

In order for a RISC-V system to run programs on top of the Linux kernel, it needs to, at minimum, implement the G (general purpose) and C (compressed) extensions from the unprivileged ISA specification, as well as the privileged ISA specification, which are explained below with respect to this work.

2.2.2 The RISC-V Unprivileged ISA Specification

The RISC-V unprivileged specification (in revision v2.1 as of writing this manual) specifies the existence, operation, formats, and bit codes for the base instruction set as well as numerous extensions. In addition to instruction listings, the unprivileged ISA also includes directives for interrupt subroutines (exceptions, traps, and interrupts), counters, and registers, as well as definitions for numerous RISC-V terms including different execution environments (EEs) and hardware threads (harts).

A Linux-capable RISC-V system with the G and C extensions is referred to as a RV32GC, RV64GC, or RV128GC system depending on the word size. The compressed extension defines 16-bit instruction layouts and the general purpose extension is a composite extension comprised of the *IMAFDZicsrZifencei* extensions, outlined in table 1.

Extension	Version	Description
RV32I	2.1	Base 32-bit integer extension.
RV64I	2.1	Base 64-bit integer extension.
M	2.0	Multiply/divide extension.
A	2.1	Atomic operations extension.
F	2.2	Single-precision floating point extension.
D	2.2	Double-precision floating point extension.
C	2.0	Compressed instruction formats extension.
Zicsr	2.0	CSR interface instructions extension.
Zifencei	2.0	Instruction-fetch fence instruction extension.

Table 1: RISC-V Unprivileged Specification GC extension versions as of writing this manual. All extensions presented above are ratified.

2.2.3 The RISC-V Privileged ISA Specification

The RISC-V privileged specification (in revision 1.11 as of writing this manual) specifies different privilege modes of operation (user, supervisor, hypervisor, and machine), and it is split into Supervisor-Level and Machine-Level ISAs (7).

The Machine-Level ISA contains the definitions and layouts of real CSRs and the privileged specification instructions accessible in M-mode (machine mode), as well as a description of physical memory protection and attributes (PMP/PMAs). The CSRs, unlike normal argument and temporary registers housed in a CPU's register file, may be memory-mapped and contain a lot of specific information pertaining to a variety of system functions, including the current status of the system, interrupt information, system information, and counters. The new instructions introduced in the Machine-Level ISA include environment call and breakpoint instructions, trap-return instructions, and a wait for interrupt instruction.

The Supervisor-Level ISA contains the definition and layouts of supervisor CSRs, most of which shadow the existing Machine-Level ISA CSRs. It also introduces a supervisor fence instruction, virtual memory management, and a paging algorithm.

2.2.4 Introduction to gem5 and gem5-X

Introduced initially in 2011, gem5 is a modular computer architecture simulator that, unlike RTL simulators, enables system-level design space exploration by simulating high-level event-driven software models for processors, peripheral devices, and memory. It comes with numerous CPU, RAM, and device models right out of the box and includes varying levels of support for numerous ISAs including x86, ARM, and RISC-V.

gem5's primary running modes are Syscall Emulation (SE) and Full System (FS). SE mode enables program simulation with simple, emulated interrupt and syscall handling. FS mode on the other hand enables program simulation on top of a full system stack, including memory hierarchy, operating system, disk image, and full interrupt service routines by in interrupt controllers (5).

gem5-X, published in April 2019 by the Embedded Systems Laboratory at EPFL, extends gem5 by introducing architectural extensions such as in-cache computing and 3D-stacked HBM models, as well as a methodology for optimizing the power and performance of many-core systems (6). It also implements several enhancements in gem5, listed below:

- ARM-64 Full System support by way of an Ubuntu 16.04 LTS disk image with kernel v4.3.
- Profiling capabilities within the simulation using the gperf profiler.
- Enhanced checkpointing by way of Region-Of-Interest (ROI) marking in applications.
- 9P over Virtio for fast modification of files without modifying the root file system, enabled by default and built into the kernel.

2.3 Prior Work

Varying levels of support for RISC-V have been introduced to the main gem5 repository over the years, but unfortunately there has been no full published effort bringing Linux-capable FS mode systems for RISC-V in gem5.

The first published effort bringing RISC-V to gem5 was published in CARRV (workshop on Computer Architecture Research with RISC-V) 2017 by Alec Roelke and Mircea R. Stan (15). In their work, "RISC5: Implementing the RISC-V ISA in gem5", they implemented RISC-V unprivileged GC extensions and verified their results against a single-core system running in gem5 SE mode. This was extended by Tuan Ta, Lin Cheng, and Christopher Batten at CARRV 2018 with their work, "Simulating Multi-Core RISC-V Systems in gem5", which brought multi-core SE mode simulation to gem5 (16).

The work closest to bringing Linux-capable FS mode to gem5 for RISC-V systems is a Master's Thesis written by Robert Scheffel of Technische Universität Dresden (TUD) (10). In this work, Scheffel implemented a RISC-V bare-metal-capable FS mode system in gem5. In this case, Scheffel could run binaries stored on a simulated flash device loaded in RAM without an operating system or MMU. This required implementing interrupts and exceptions, in addition to a few peripheral devices implicitly required by the RISC-V ISA.

Finally, Nils Asmussen of the Barkhausen Institut contributed a RISC-V Sv39 MMU to the main gem5 repository for use with microkernels (24). This MMU model is based off of the x86 MMU model in gem5 but does not support Linux.

2.4 Contributions

The base target execution environment for this work is a combination software and hardware execution environment with both an application binary interface (ABI) and supervisor binary interface (SBI). In other words, **our main contribution with gem5-eXtensions for RISC-V is creating the first Linux-capable FS mode system in gem5**. The work leading up to this is as follows:

1. We extend the FS mode configuration in gem5 for Linux-capable RISC-V systems.
2. We implement instructions from the RISC-V privileged specification and extend or verify instructions implemented in prior work.
3. We implement the missing Zifencei extension from the unprivileged specification.
4. We extend CSR implementations from prior work.
5. We develop a RISC-V compliant MMU capable of processing virtual memory, checking PM-P/PMA, and interfacing a page table walker.
6. We implement RISC-V ISA devices, including a PLIC (platform-level interrupt controller) and CLINT (core-local interrupter).
7. We develop and configure a gem5-compatible bootloader, Linux kernel, device tree, and buildroot image for storage.

With these contributions combined, we are able to demonstrate running programs on top of the Linux kernel, on top of a disk image. In the next section we discuss and describe how to set up and use **GXR5**. In the rest of this technical manual we detail the high-level implementation of our contributions.



3 Running gem5 Full System Mode with RISC-V and Linux

In this chapter we describe how we configured and ran our RISC-V model, ending with a quick-start guide.

3.1 Necessary Files

Because our model is run in FS mode with a full Linux environment, we need several major system components not included with gem5. This includes,

- A bootloader
- A static kernel binary, e.g., vmlinux
- A file system/disk image
- A device tree binary

Additionally, all of the aforementioned components must be compatible with the RV64GC flavor of ISA.

All of the configuration files are located under `system/riscv/` in the `bootloader`, `disk`, `dt`, and `linux` folders. In the following sections we describe the configuration options and how to use them in brief. Our steps for building each of the components follow from each respective component's own quick-start guide and it is assumed you have followed their guide and set up their respective environments. Once you have all of the files configured, it is possible to test them using the RISC-V variant of soft-mmio QEMU (2). Note that a common dependency is the RISC-V GNU toolchain, which includes a cross-compiler necessary for compiling the bootloader and kernel.

Additionally, you should also set up gem5 full system mode by creating the folders `src/full_system_images/disks` and `src/full_system_images/binaries` to house your additional files. These folders should sit on your M5 path after you set up gem5, which can be done with the following:

```
1 export M5_PATH=/path/to/gXR5/full\_system\_images
```

3.1.1 Device Tree

Our device tree is custom built, initially modified from a device tree generated by QEMU for a Fedora Linux emulation (2) (11). Once you have the device tree compiler set up, you can compile the device tree structure (dts) file using the following:

```
1 dtc -I dts -O dtb gem5-simple-rv64.dts -o gem5-simple-rv64.dtb
```

The output of this command is a device tree binary (dtb) file that can then be used by copying it to your **GXR5** FS mode directory, under `full_system_images/binaries`.

3.1.2 Bootloader

The bootloader we used for testing is the Open Supervisor Binary Interface, or OpenSBI. OpenSBI is a first-stage bootloader and platform-independent M-mode firmware designed to link a library with a simple set of defined methods for a specific platform. The RISC-V manual also includes a SBI specification that OpenSBI conforms to (9) (17).



For our experiments, we tested the OpenSBI bootloader compiled for the qemu/virt platform using a custom configuration file. Copy the configs.mk, objects.mk, and platform.c files to the qemu/virt directory in your OpenSBI folder. After this, you can run the following commands to build the bootloader:

```
1 git clone https://github.com/riscv/opensbi.git
2 cd ./opensbi
3 git checkout 813f7f4c250af9f7c9546f64778e9b35bb7d7dcb
4 export CROSS_COMPILE=/path/to/riscv/compiler/here
5 export PLATFORM_RISCV_XLEN=64
6 make PLATFORM=qemu/virt O=/path/to/your/build/directory
7 make PLATFORM=qemu/virt I=/path/to/your/install/directory
```

The result of these commands, if run successfully, should include a file "fw_jump.elf". This can be used as a standalone bootloader which launches a vmlinux file during run-time. To install for **GXRS**, copy the file to full_system_images/binaries.

3.1.3 Linux Kernel

Under most circumstances, gem5 is only able to directly run static binaries, and thus we must use a static version of the Linux kernel. This refers specifically to the file generated by building the Linux kernel, vmlinux.

For our tests, we built and used Linux kernel version 5.5 directly from the Linux repository (18). You will need to git checkout this version and then rebuild Linux with our custom configuration file. Copy one of the provided Linux configuration files into your Linux directory as the file ".config". You can run the following:

```
1 git clone https://github.com/torvalds/linux.git
2 cd ./linux
3 git checkout v5.5
4 make ARCH=riscv CROSS_COMPILE=/path/to/riscv/compiler/here
5 make ARCH=riscv CROSS_COMPILE=/path/to/riscv/compiler/here all
```

The result of running these commands should be both a vmlinux and Image file containing a RV64GC-compatible kernel. Though the Image file is not used in gem5, you can use it in QEMU to verify a successful RISC-V build in lieu of the vmlinux file.

Like before, in order to use your vmlinux file you will need to copy over to your **GXRS** directory under full_system_images/binaries.

3.1.4 Disk Image

The disk image we used is a minimal rootfs created by buildroot (19). Similarly with building Linux, in order to build the root file system you will need to copy over the configuration file and make. **Note that if you previously did not set up buildroot with RISC-V, buildroot will need to download and set up the entire RISC-V GNU toolchain, so your first time creating the minimal rootfs will likely take tens of minutes** (if not longer) depending on your computer and internet connection.

So just like with Linux, you will need to copy the provided buildroot.config file into the .config file of the buildroot directory. By default, the rootfs file will be a 60MB image with only the minimal root file system. The username and password are set to root/buildroot by default, but this



can be changed in the nconfig menu for buildroot, or for auto-login, in /etc/inittab of the generated rootfs.ext2. Note that the rootfs file is technically marked as an ext4 file system, but what is generated is an ext2 file system with an ext4 link. Linux will still treat this file system as ext4 however.

So after you have completed configuration of your rootfs file, you are ready to build it. Simply call "make" in the buildroot directory and the disk image, rootfs.ext2, will be available in output/images. In order to use it, copy it over to your **GXR5** directory under full_system_images/disks.

Lastly, just like with the Image file in Linux, you can use this rootfs file with RISC-V QEMU to verify its contents. Because it is a miniaml root file system however, it will not have the proper utilities for creating test programs, so simply mount it and copy over files from your host system using the following:

```
1 sudo mount <path_to_rootfs>/rootfs.ext2 /path/to/mount/point
2 cp /your/binary /path/to/mount/point
3 sudo umount /path/to/mount/point
```

3.2 Recommended Script Options

In addition to the files included above, we tested our model using the AtomicSimpleCPU model, 2GB of RAM, and the DDR3_2133_8x8 memory type. Future work will include testing and verification of other CPU models with respect to real hardware.

3.3 Quick-Start Guide

In this brief start-up guide, we will guide you through the basic steps to running your first experiment with gXR5. This guide assumes you have already built a working bootloader, device tree, static kernel file, and disk image using the aforementioned configuration files as described in the previous sections and placed them in the proper file locations: gem5-simple-rv64.dtb, fw_jump.elf, and vmlinux in gXR5/full_system_images/binaries and rootfs.ext2 in gXR5/full_system_images/disks.

3.3.1 Prerequisites

You will need to set up the gem5 environment in order to compile and run the gem5 binary using the scon (SConstruct) builder. If running on an Ubuntu-based host system, you can use the following command to get all the required libraries:

```
1 sudo apt install build-essential git m4 scon zlib1g \
2 zlib1g-dev libprotobuf-dev protobuf-compiler libprotoc-dev \
3 libgoogle-perftools-dev python-dev python-six python \
4 libboost-all-dev
```

3.3.2 Building the gem5 Binary

Once the above is done, you will need to build a RISC-V gem5 binary. You can create multiple builds including .fast, .opt, and .debug. If you are only concerned about running experiments, it is recommended to only create gem5.fast. However, if you need to debug anything or want to generate traces, you will need to build gem5.opt or gem5.debug. Do this with the following:



```
1 cd gXR5
2 scons build /RISCV/gem5.{fast, opt, debug}
```

Additionally, if you would like to speed up the compilation process, you can use the option "-jN" on the scons build line where N is the number of threads you want to assign for compilation.

3.3.3 Setting Up Your Experiment

Before you start running gem5, you will need to set up the rootfs image to run your desired program. Copy over your desired program to the rootfs, and then modify the /etc/inittab in your rootfs to include the program you want to run. The inittab is responsible for initializing various file systems and the login screen but simple programs should be safe to run this way because the kernel has already entirely booted by the time inittab is run.

For example, if you want to run a cross-compiled binary called "rvhello" sitting in my file system's bin folder, you would run the following:

```
1 riscv64-linux-gnu-gcc rvhello.c -o rvhello
2 mkdir /mnt
3 sudo mount gXR5/full\_system\_images/disks/rootfs.ext2 /mnt
4 sudo cp rvhello /mnt/bin
5 sudo vi /mnt/etc/inittab
6 sudo umount /mnt
```

When you run "sudo vi" in the above, you will modify the inittab to have the following:

```
1 # /etc/inittab
2 #
3 # Copyright (C) 2001 Erik Andersen <andersen@codepoet.org>
4
5 ...
6
7 # Run your binary
8 ::sysinit:/bin/rvhello
9
10 ...
```

Exit vi with escape + ":wq" + enter, and then you can unmount your file system and you are ready to run your experiment.

3.3.4 Running Your Experiment

Once your program sits in your rootfs and your inittab file is modified appropriately, run **GXR5** using a script that has some variation of the following:

```
1 cd gXR5
2 ./build/RISCV/gem5.{fast, opt, debug} \
3 -d /path/to/your/output/directory \
4 ./configs/example/fs.py \
5 --disk-image=full_system_images/disks/rootfs.ext2 \
6 --kernel=full_system_images/binaries/vmlinux \
7 --os-type=linux \
```



```
8  --dtb-filename=full_system_images/binaries/gem5-simple-rv64.dtb \  
9  --cpu-type=AtomicSimpleCPU \  
10 -n 1 \  
11 --mem-size=8GB \  
12 --mem-type=DDR3_2133_8x8 \  
13 --cpu-clock=1GHz \
```

At this point you should be able to connect to your running gem5 instance in another shell with,

```
1  telnet localhost 3456
```

Upon connecting to your gem5 instance, you should be able to see the OpenSBI logo as well as the kernel dmesg, followed by whatever output is apart of the program you are running (including basic "Hello, world!" programs and other benchmarks).

4 Full System Configuration Overview

The goal with this first revision of a RISC-V FS mode-compatible model was to create the simplest possible computer system for running basic benchmarks that still aligns with performance and power models of real RISC-V processors. In this chapter we describe the FS mode configuration in detail, and, in brief, our platform.

4.1 Full System Model and Child Tree

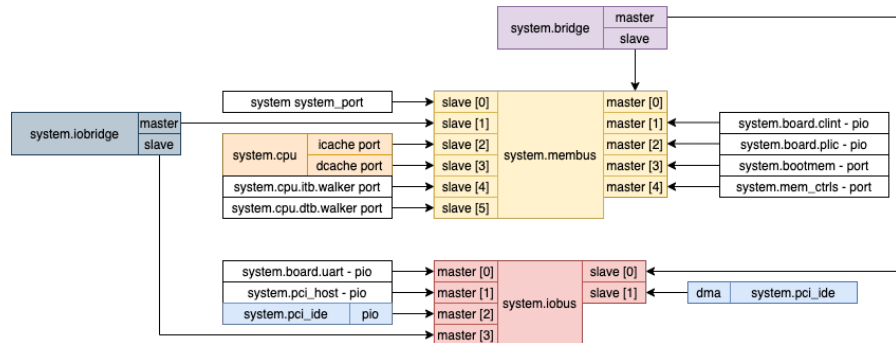


Figure 1: RISC-V Full System configuration with all models used for simulation and their interfaces. The diagram on the left An arrow from object A pointing to object B indicates that A sets its associated port to B within the gem5 configuration script.

The RISC-V Full System model consists of numerous built-in gem5 models as well as custom ISA-specific models. As shown in figure 1, the model includes a single-core CPU with instruction and data caches, TLB walkers for each of the aforementioned caches, a (currently unused) bootmem, a Platform-Level Interrupt Controller (PLIC), a Core-Local Interrupter (CLINT), a PCI host with IDE controller, and lastly a UART module. All models are connected to either a memory bus (membus) or IO bus (iobus), and these buses are interfaced through both a system bridge and IO bridge. Of these models mentioned, only the PLIC and CLINT are RISC-V-specific models required by the ISA. Additionally, the PCI host is customized for our implementation. The rest of the models are ISA-independent and are provided by gem5's built-in utilities.

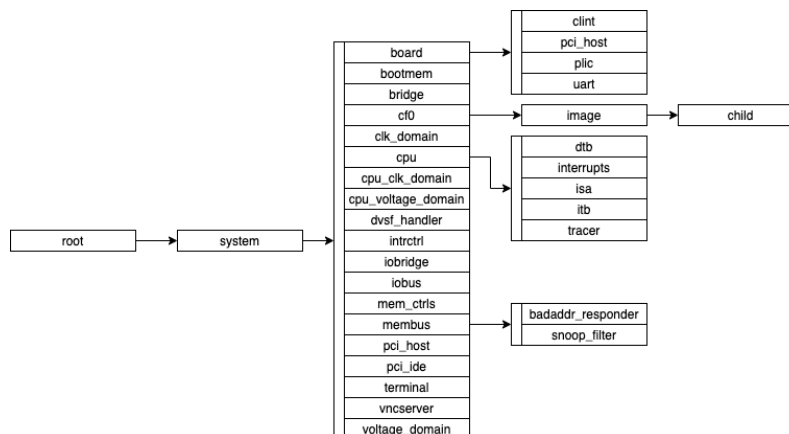


Figure 2: RISC-V Full System gem5 child tree.

The child tree is shown in figure 2 and shows the gem5 model hierarchy of the simulated system.

The full system configuration is implemented in `configs/example/fs.py` and `configs/common/FSConfig.py`. `fs.py` contains the general FS-mode setup while `FSConfig.py` contains the specific RISC-V system setup, including mapping memory and making all port connections.

4.2 Full System Memory Map

The full system memory map can be seen in table 2. It is defined in `configs/example/fs.py` and `src/dev/riscv/SimpleBoard.py`.

Device	Address Range
PLIC	0x0c000000:0x0c2fffff
UART	0x10000000:0x10000100
CLINT	0x20000000:0x2effffff
PCI	0x2f000000:0x5fffffff
RAM	0x80000000:0xffffffff
– bootloader	0x80000000:0x801fffff
– kernel	0x80200000:0x87ffffff
– device tree	0x88000000:0x8fffffff

Table 2: Rough memory map of simulated full system.

For a precise listing of all CSRs and their offsets, please refer to the ISA specification (7) (8).

4.3 SimpleBoard Platform

To interface ISA-specific devices as well as external devices, we further develop the SimpleBoard platform that was initially introduced in Scheffel's thesis (10). The relation between the CPU core and SimpleBoard SoC, as well as external devices, can be seen in figure 3.

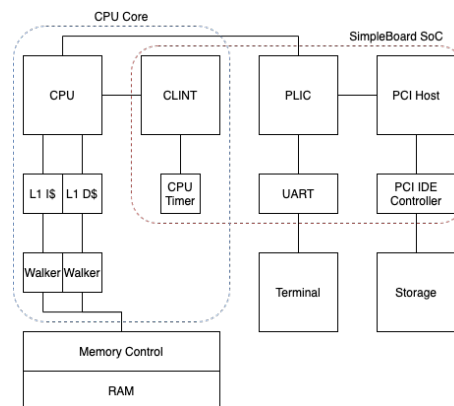


Figure 3: Target hardware system for emulation.

In order to enable gem5 FS mode, models for both a CLINT and PLIC needed to be developed and interfaced through the SimpleBoard platform. Additionally, a platform in gem5 is used to also interface certain external interrupts such as interrupts from PCI devices as well as implement



timers, connect a terminal, etc. All code for the SimpleBoard platform and associated devices can be found in `src/dev/riscv`.



5 Implementation of Privileged RISC-V ISA Specification

In this section we go into more detail about how the privileged ISA specification was implemented in **GXR5**, including the CSRs, instructions, interrupt handling, and virtual memory subsystems.

5.1 Control and State Registers (CSRs)

In this section, we discuss CSRs in RISC-V and their implementation in gem5.

5.1.1 CSRs versus Other Registers

The RISC-V ISA specification implicitly designates two kinds of registers: the aforementioned Control and State Registers and what this manual will refer to as 'normal' CPU registers. The normal registers include the common registers seen across many ISAs, including argument registers, temporary registers, stack pointer, program counter, etc. CSRs however contain very specific information with respect to the operation of RISC-V systems. Furthermore, CSRs don't have to be in the register file of the CPU and can instead be *memory-mapped* to help interface external devices. Keep in mind however that a CPU will still see and access memory-mapped CSRs as though they were in the register file.

Some of the most important CSRs include *mstatus*, which contains information about the current running status of the system, *mip* and *mie*, which show the pending and enabled interrupts in the system, and *mtime*, which is used as a global system timer. The full CSR listings can be found in the RISC-V privileged specification (7).

CSR names that start with 'm' refer to a CSR only accessible in machine (M) mode. The RISC-V privileged specification also lists supervisor (S) mode CSRs, most of which are not actually CSRs separate from their M-mode counterparts but are shadows of the same M-mode CSRs where M-mode-only bits cannot be reliably read or written from. There are some CSRs completely unique to S-mode however, such as the *satp* CSR.

5.1.2 CSRs in gem5-X

The register file for RISC-V is located in `src/arch/riscv/registers.hh`. This file defines the maps that contain both the normal register file registers as well as the CSRs. It is in this file that CSR indices, names, and offsets are stored and referenced in gem5. Additionally, the `registers.hh` file defines BitUnions for easy access to the different fields of the CSRs. Most of these CSRs were previously defined in prior work, and only some performance and timing registers needed to be added for **GXR5**.

5.1.3 Interfacing and Accessing CSRs

As the register file is usually unique for each hart, the usual way to interface the CSRs is via the `threadContext` class using the methods `getMiscReg` and `setMiscReg`. These methods are defined and implemented for RISC-V in `src/arch/riscv/isa.hh` and `src/arch/riscv/isa.cc`, respectively.

Hardware events usually call the ISA methods when a CSR needs to be checked or modified. For example, a CPU read to the *mtime* CSR must read the external system clock, and therefore the ISA interface will access the system's CLINT to read that CSR. Another example is when an



external interrupt needs to be posted by the PLIC, it will update the *mip* CSR through the ISA interface.

Note that with memory-mapped CSRs, the ISA interface only defines direct access to the register. Usually when a CSR must be accessed directly by address, it is through gem5's packet interface (typically associated with PIO devices), and so another interface is defined to connect external devices and the ISA interface in `src/arch/riscv/isa_device.hh` (and subsequently implemented in `src/arch/riscv/isa_device.cc`).

5.2 RISC-V Instructions

In this section, we discuss the new instruction implementations in RISC-V and gem5.

5.2.1 Instructions in gem5

gem5 implements instructions by using ISA template files that then generate all of the instruction classes and functionality during the SConstruct build process. The ISA templates for RISC-V are located in `src/arch/riscv/isa/`. The primary file used for implementing individual instruction functions is `src/arch/riscv/isa/decoder.isa`, which decodes machine code to assign and process instruction types. Most of the RV64GC instructions from the unprivileged ISA specification were implemented in prior work.

5.2.2 Privileged Specification Instructions

The privileged RISC-V spec adds very few instructions to the base ISA. These instructions are URET, SRET, MRET, SFENCE.VMA, WFI, HFENCE.BVMA, and HFENCE.GVMA. As of writing this manual, the Hypervisor specification has not yet been ratified, so we forego the implementation of the HFENCE instructions. The table of instruction layouts (Table 5.1 in the unprivileged specification document) is recreated in table 3 for convenience.

00000000	00010	00000	000	00000	1110011	URET
00010000	00010	00000	000	00000	1110011	SRET
00110000	00010	00000	000	00000	1110011	MRET
00010000	00101	00000	000	00000	1110011	WFI
00000000	rs2	rs1	000	00000	1110011	SFENCE.VMA
00000000	rs2	rs1	000	00000	1110011	HFENCE.BVMA
00000000	rs2	rs1	000	00000	1110011	HFENCE.GVMA

Table 3: Instructions introduced in the RISC-V Privileged ISA specification and their bit layouts.

The URET, SRET, and MRET instructions are all trap return instructions for specific privilege modes. These instructions were initially implemented in prior work and only verified for use with Linux-capable FS mode in this work.

The WFI instruction is a "wait for interrupt" instruction. The RISC-V privileged specification states that a no-operation (NOP) is a valid implementation for WFI, so we left the instruction as such.

The SFENCE.VMA instruction is a supervisor fence that is used to enforce memory ordering. While the RISC-V privileged specification goes into very tedious and specific detail about memory ordering and constraints, the result of the specification's discussion is a relatively simple



hardware cache flush, implemented in gem5 using the overridden demapPage method found in src/arch/riscv/tlb.cc. A call to SFENCE.VMA will flush a specific page or set of pages of the L1 ITB and DTB caches by address space number (ASN), virtual address, both, or neither (resulting in a full cache flush) depending on the values of its arguments rs1 and rs2.

5.2.3 Unprivileged Specification Instructions

Only two instructions necessary for Linux-capable FS mode were missing from prior work. The instructions and their layouts are in tables 4 and 5.

fm	pred	succ	rs1	000	rd	0001111	FENCE
----	------	------	-----	-----	----	---------	-------

Table 4: Missing fence instruction from the RV32/64 Base Instruction Set.

imm[11:0]	rs1	001	rd	0001111	FENCE.I
-----------	-----	-----	----	---------	---------

Table 5: Missing fence instruction from the RV32/RV64 *Zifencei* Standard Extension.

The FENCE instruction is simply implemented as a flush instruction for the entire L1 instruction and data caches (ITB and DTB). The FENCE.I instruction is specifically for L1 ITBs, so only the ITB is flushed. The cache flush implementation in gem5 is the same as described in the previous section for the SFENCE.VMA instruction.

5.3 Fault Handling

In this section, we discuss faults, exceptions, and interrupts in RISC-V and their implementation in gem5.

5.3.1 RISC-V Terminology

The RISC-V ISA specification defines numerous terms relating to trap/faults. A synchronous fault is referred to as an *exception*, while an asynchronous fault is referred to as an *interrupt*.

Exceptions include faults due to misaligned instructions/addresses, invalid instruction/address access, illegal instruction calls, breakpoints, environment calls, and page faults. Interrupts include software, timer, and external asynchronous faults in different privilege modes.

The fault number of an exception or interrupt is stored in the cause CSR, which holds fault causes in a one-hot representation and reserves its top bit to indicate if the fault is an interrupt or an exception.

5.3.2 RISC-V Fault Handling Algorithm

Because RISC-V relies on its various binary interfaces for fault handling, a lot of the algorithm for handling various faults is delegated to software. This makes the implementation of fault handling in hardware extremely easy. The RISC-V hardware simply saves the current context (privilege mode, current PC value), escalates the privilege mode (unless delegated via the mideleg and medeleg CSRs), and sets the PC to the saved address of the fault handler. The cause value stored in the m/s/ucase CSR will tell the software fault handler which specific routine needs to be



taken. When the fault is handled, it calls the `m/s/uret` instruction to restore context and continue program execution.

5.3.3 Exceptions in `gem5`

Prior work had defined and implemented RISC-V exceptions in `src/arch/riscv/faults.hh` and `src/arch/riscv/faults.cc`, respectively. All fault types derive from a base `RiscvFault` class and most faults will use the same `invoke` method which implements the fault handling described in the previous section. The derived fault classes are mostly for implementation simplicity, but in this work we also had to ensure privilege modes were implemented for the `ecall` fault type.

Address, misalignment, instruction, and page faults of all types are typically posted by the TLB, defined in `src/arch/riscv/tlb.cc`.

5.3.4 Interrupts in `gem5`

Interrupts in RISC-V use the same fault handling algorithm as exceptions with some minor changes based on the fault code and cause CSR, and therefore use the same aforementioned fault handler defined for exceptions in `gem5`. To use the same fault handler, prior work had already defined and implemented an `InterruptFault` class in `src/arch/riscv/faults.hh` and `src/arch/riscv/faults.cc`, respectively.

Interrupts are posted from different (usually external) sources at asynchronous times. `gem5` offers a CPU interface with `postInt` and `clearInt` methods that set and clear an interrupt pending array defined in `src/arch/riscv/interrupts.hh`.

5.4 Virtual Memory

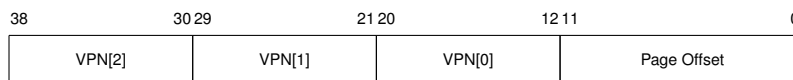


Figure 4: Sv39 virtual address.

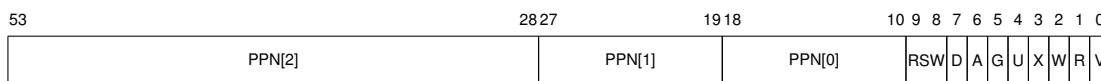


Figure 5: Sv39 page table entry.

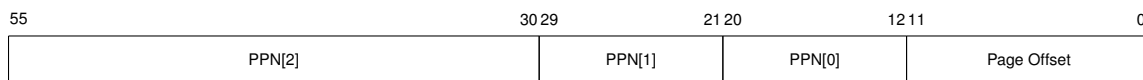


Figure 6: Sv39 physical address.

The `satp` (supervisor address translation and protection) CSR is responsible for determining the status of virtual memory. It contains three fields: the mode of translation, the address space identifier (ASID), and finally the physical page number of the root page table. When the mode is set to bare or the privilege mode is in M-mode, address translation is direct and all addresses are considered "physical". When the mode field is not bare, there are four modes of virtual address translation available.



In **GXR5**, we support only Sv39 virtual address translation. This is 39-bit virtual addressing that translates a 39-bit virtual address to a 56-bit physical address. The bit formats for Sv39 virtual address, page table entries, and physical addresses, are found in figures 4, 5, and 6. Virtual address translation is mostly implemented in `src/arch/riscv/tlb.cc`, and is discussed in more detail in the next chapter.



6 Virtual Memory Subsystem

In this section we expand upon the previous by going in-depth into our virtual memory implementation, which includes a RISC-V-compliant Sv39 MMU consisting of a TLB and page table walker.

6.1 TLB Implementation

gem5 implements a base TLB class that is used for instruction and data caches. Our gem5 TLB implementation is largely based on the TLB implementation for the ARM ISA, but with the RISC-V virtual address translation algorithm implemented. In other words, functional, atomic, and timings-based address translation calls are routed either through the `translateFs` or `translateSe` methods. In our case, we move prior work to the `translateSe` method and focus on implementing the `translateFs` method. The TLB description and implementation are located in `src/arch/riscv.hh` and `src/arch/riscv/tlb.cc`, respectively. Additionally, some page table structures are implemented and defined in `src/arch/riscv/pagetable.hh`.

The local page table of the TLB is stored in a CPP map data structure for simplicity. This data structure maps virtual address bases (virtual address without offset) to `TlbEntry` structures. The `TlbEntry` structures include the virtual address, physical address, page table entry, and ASID.

The actual address translation takes place in the `translateFs` method of the TLB. This method is the main workhorse for all address translation types, PMP/PMA checking, and cache management with respect to other TLB methods.

6.2 Physical Address Translation

When the `satp` CSR is set to bare translation mode, or the translation is occurring in machine mode, the physical address is translated directly and hence there is no need to cache the address. The address is simply checked against the PMP/PMA checker for potential access faults, before the physical address field of the requesting packet is simply set to the unchanged physical address.

6.3 Virtual Address Translation

When the `satp` CSR is not set to bare translation mode and the mode is supervisor or user, virtual address translation is required. The translation process starts with checking the TLB's map data structure. If there is a TLB hit, we simply translate the virtual address using the cached physical address. If there is a TLB miss, we proceed with the virtual address translation algorithm explained in the RISC-V privileged ISA specification, section 4.3.2. The source code is annotated with each individual step of the virtual address translation algorithm.

6.4 Page Table Walker Implementation

The page table walker definition and implementation are in `src/arch/riscv/table_walker.hh` and `src/arch/riscv/table_walker.cc`, respectively. Because a page table walker is a purely hardware component with no RISC-V specification, we use a very simple design that only acts to facilitate DMA transactions via the `walk` method.

The table walker is set up using gem5's ports interface to connect it both to the L1 caches and directly to the memory. Memory is accessed directly via the `dmaAction` method and the result is stored locally.

7 SimpleBoard Platform

The SimpleBoard platform, extended from prior work, inherits from the gem5 platform class. Platforms are typically used to define SoCs and various on-chip interconnects. Though our target SoC in future work will be the HiFive Unleashed 1 board, this SimpleBoard platform represents the simplest possible implementation of a platform for RISC-V that can interface its I/O devices for Linux-capable FS mode.

7.1 SimpleBoard Implementation

Our definition and implementation of the SimpleBoard platform can be found in `src/dev/riscv/simpleboard.hh` and `src /dev/riscv/simpleboard.cc`, respectively. Additionally, the simulation object parameters for the SimpleBoard platform, in addition to the pythonic class definitions for the other devices hosted on the SimpleBoard are located in `src/dev/riscv/SimpleBoard.py`.

In addition to hosting SoC devices for RISC-V systems, the SimpleBoard platform is also used to interface PCI and console interrupts, as well as host a generic UART device. The PCI interrupt interface is required by our implementation for disk access, however the console interrupt interface is currently unused.

The SimpleBoard SoC is instantiated and connected to the main system through the gem5 FS mode configuration script.

7.2 SimpleBoard Devices and Parameters

All of the SoC devices are defined and linked to the SimpleBoard platform in `src/dev/riscv/SimpleBoard.py`. These devices include the PLIC, CLINT, UART, and a RISC-V PCI host. A visual representation of these devices can be found in figure 7.

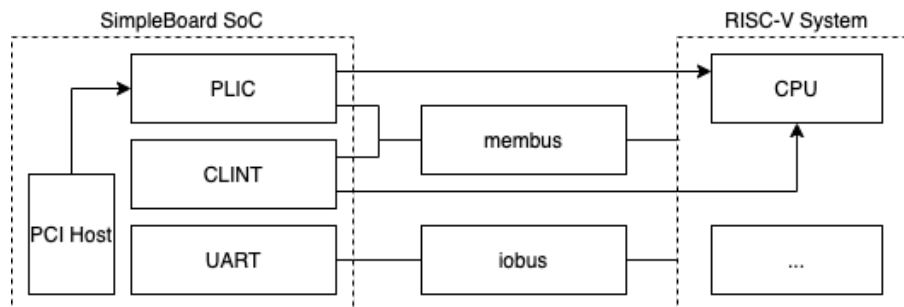


Figure 7: A rough diagram showing the interconnects between the SimpleBoard SoC devices and the RISC-V system. The lines with arrows represent interrupt lines to the PLIC and CPU while the other connections are either peripheral I/O (PIO) interconnects or interconnects via the `isa_device` interface. PIO interfaces use packets and addresses to connect these devices.



8 Platform-Level Interrupt Controller (PLIC)

In order to handle potentially simultaneous interrupts signals that come from peripheral devices, the RISC-V ISA manual specifies the existence of a PLIC. The goal of the PLIC is to receive all external interrupt signals (so-called 'global interrupts'), order them by priority, and then delegate them to an available hart for handling.

For every potential source of a global interrupt, the PLIC defines a source number and highest allowable priority mask. An interrupt source numbered 1 with a priority 7 has the highest interrupt priority. Interrupt source 0 is reserved to mean 'no interrupt'.

8.1 PLIC Implementation

In our implementation, the PLIC is a BasicPioDevice on the SimpleBoard platform, and the first of our ISA devices. The definition and implementation of the PLIC are found in `src/dev/riscv/plic.hh` and `src/dev/riscv/plac.cc`, respectively. The configuration of the registers is defined as a map in the source code and follows the same layout at the PLIC described in the FU540-C000 core manual (20). This layout is reproduced in the next section.

8.2 PLIC Registers and Memory Layout

All PLIC registers are 32-bit and are listed in table 6. Additionally, all PLIC registers are read-write registers, except for the pending array registers which are read-only. 64-bit registers are split into two 32-bit registers with a low and high field that is automatically interfaced by RISC-V software when accessing the PLIC.

Note that this PLIC was designed for the HiFive Unleashed SoC with five harts (one small CPU core and 4 large CPU cores), even though our simulated FS mode system only simulates one core (and therefore one hart) currently. We keep the registers for other harts for future work.

8.2.1 PLIC Source Registers

Each source number refers to a specific external device and designates a priority. For example, in the FU540-C000 manual, sources 1-3 refer to the L2 cache controller, source 4 refers to UART0, and so on. The values stored in the registers are 3-bit source priority values. A value of zero indicates that interrupts are disabled for the source, while a value of seven indicates that interrupts are of the highest priority for the source.

Like with the additional registers for additional harts, we preserve all of the source registers for future work. In our FS mode simulation, we technically only use the sources for UART and the PCI host.

8.2.2 PLIC Interrupt Pending Array

The PLIC interrupt pending array is a one-hot representation of pending interrupts where the bit index refers to the source number. For example, if bit 12 of the interrupt pending array is high, it means source 12 is awaiting interrupt handling.



Offset	Register
0x0c000004	Source 1 Priority Register
0x0c000008	Source 2 Priority Register
...	...
0x0c0000d8	Source 54 Priority Register
0x0c001000	Interrupt Pending Array (low)
0x0c001004	Interrupt Pending Array (high)
0x0c002000	Hart 0 M-Mode Interrupt Enable (low)
0x0c002004	Hart 0 M-Mode Interrupt Enable (high)
0x0c002080	Hart 1 M-Mode Interrupt Enable (low)
0x0c002084	Hart 1 M-Mode Interrupt Enable (high)
0x0c002100	Hart 1 S-Mode Interrupt Enable (low)
0x0c002104	Hart 1 S-Mode Interrupt Enable (high)
...	...
0x0c002380	Hart 4 M-Mode Interrupt Enable (low)
0x0c002384	Hart 4 M-Mode Interrupt Enable (high)
0x0c002400	Hart 4 S-Mode Interrupt Enable (low)
0x0c002404	Hart 4 S-Mode Interrupt Enable (high)
0x0c200000	Hart 0 M-Mode Priority Threshold
0x0c200004	Hart 0 M-Mode Claim/Complete
0x0c201000	Hart 1 M-Mode Priority Threshold
0x0c201004	Hart 1 M-Mode Claim/Complete
0x0c202000	Hart 1 S-Mode Priority Threshold
0x0c202004	Hart 1 S-Mode Claim/Complete
...	...
0x0c207000	Hart 4 M-Mode Priority Threshold
0x0c207004	Hart 4 M-Mode Claim/Complete
0x0c208000	Hart 4 S-Mode Priority Threshold
0x0c208004	Hart 4 S-Mode Claim/Complete

Table 6: PLIC register listings.

8.2.3 PLIC Interrupt Enable Registers

The PLIC interrupt enable registers masks the available interrupts using the same one-hot representation used by the interrupt pending array. If an interrupt is pending and a hart is ready to handle an external interrupt, the hart will check the interrupt pending array against its assigned interrupt enable array, based on the current hart's number and privilege mode.

8.2.4 PLIC Priority Thresholds

The PLIC priority threshold register holds a 3-bit priority value. In order for a hart of specific number and privilege mode to be able to handle an interrupt, its threshold value must be lower than the source priority value.



8.2.5 PLIC Claim/Complete Registers

The PLIC claim and complete registers are used for the PLIC claims process, described in the next section.

8.3 PLIC Claim/Complete Operations

At a high level, the PLIC is meant to perform two major operations in addition to standard reads and writes to its registers: PLIC claim, and PLIC claim complete.

When a hart is ready to handle an external (global) interrupt, it 'claims' the interrupt (making it unavailable to other harts) by sending a read request to the hart's claim register. This will read the PLIC's interrupt pending array, clear the bit associated with the highest priority pending interrupt, and return the source number to the hart. The hart can then use the source number to jump to the correct interrupt handler for the source.

When a hart is finished handling an external interrupt, it performs a 'claim complete', where it sends a write request to its claim/complete register. This register will hold the source number of the last completed external interrupt.

8.4 PLIC postInt and clearInt Methods

In our FS mode implementation, we use postInt and clearInt to send interrupts to the CPU. These methods are used primarily by the platform to forward PCI interrupts. Additionally, we only send M-mode and S-mode interrupts to the CPU; U-mode external interrupts can be delegated depending on the mideleg CSR.

The sendInt method will send an external interrupt to the CPU if the source priority exceeds the source threshold as well as set the appropriate bit in the PLIC pending array. Likewise, the clearInt method will clear an interrupt in the CPU and clear the appropriate bit in the PLIC pending array.



9 Core-Local Interrupter

The CLINT, expanded from prior work (10), houses the system's main clock. As a result, the CLINT is responsible for posting timer interrupts to its local CPU and associated harts. While there would usually only be one PLIC in a real system, there is usually one CLINT per hart.

In addition to posting timer interrupts, the CLINT is also responsible for posting software interrupts. Our implementation loosely follows that of the CLINT implementation of the HiFive Unleashed SoC, so our CLINT will post only M-mode software interrupts (20). As software interrupts are typically only used for cross-hart communication, they remain untested in **GXR5**.

9.1 CLINT Implementation

Like the PLIC, our CLINT is implemented as a BasicPioDevice on the SimpleBoard platform that is also an ISA device. The definition and implementation of the CLINT are found in `src/dev/riscv/clint.hh` and `src/dev/riscv/clint.cc`, respectively. In **GXR5**, the CLINT is set up such that there is only one instance of the object in the whole system, but to handle multi-core and multi-hart systems in future work, we define a `CpuTimer` class that contains core-local and hart-local registers. The timer in the CLINT is, by default, set to 1MHz. The registers and layout of the CLINT are shown and described in the following section.

9.2 CLINT Registers and Memory Layout

Like with the PLIC, the CLINT register layout follows that of the FU540-C000 core from the HiFive Unleashed SoC, and can be seen in table 7. There is only one `mtime` register in the CLINT, but there is one `mtimecmp` and `msip` register per hart. All registers are read-write. Unlike the PLIC, the offsets of the `mtimecmp` and `msip` registers are generated as a function of the base offset added to the width of those register.

Offset	Bit Width	Register
0x00000000	4	Hart 0 Interrupt Pending
0x00004000	8	Hart 0 Time Compare
0x0000bff8	8	Hart 0 Timer

Table 7: CLINT register listings.

For example, for hart 1 the offset of the interrupt pending (`msip`) and time compare (`mtimecmp`) CSRs would be `0x00000004` and `0x00004008`, respectively.

9.2.1 CLINT Timer Register and the `mtime` CSR

The `mtime` CSR from the CPU's register file is a memory-mapped CSR. A read/write operation to the CSR must perform the operation to the actual register location, which in our case, is the CLINT. The CLINT stores the `mtime` CSR as a simple 64-bit unsigned integer, which can be read/written to directly via address or via the ISA device interface.

In a real (non-simulated) system, the CLINT's timer would be tied to an external oscillator tuned to a constant frequency. Reading the `mtime` CSR means reading the timer register, which returns the number of cycles the oscillator has processed since system start. As `gem5` is an event-based simulator however, implementing a timer directly in this way would be incredibly inefficient and

increase simulation times significantly due to the number of added discrete events to the event queue. Therefore, the timer only changes with discrete events – primarily reads/writes to the CLINT.

When the timer needs to be updated, it is done through the `updateTime` method. This method calculates the timer's current value as a function of the default clock frequency of the timer and the current tick of the simulator. Note that it is also possible to write to the timer register, and therefore we preserve an offset value as well to incorporate into the timer update method.

9.2.2 CLINT Time Compare Register

The CLINT timer compare register, `mtimecmp`, is a register that determines when a timer interrupt must be posted. A timer interrupt is posted whenever the `mtime` CSR is greater than or equal to `mtimecmp`. Setting `mtimecmp` to `INT_MAX` effectively disables timer interrupts. The `mtimecmp` register only determines when a timer interrupt is posted to the hart it is attached to, and in a MPSoC, there would be one `mtimecmp` register per hart.

9.2.3 CLINT Software Interrupt Pending Register

The `msip` register is supposed to hold one bit that indicates if a software interrupt is pending or not. In our implementation however, we tie the `msip` register directly to the `mip` CSR, and thus a read or write operation to the `msip` register in the CLINT is directly reflected in the `mip` CSR instead.

9.3 CLINT Timer Interrupts

Timer interrupts occur whenever a CPU's `mtimecmp` register is greater than or equal to the `mtime` CSR. In our implementation, we use `gem5`'s event queue to schedule "timer alarms" that go off and post either M-mode or S-mode timer interrupts to the CPU.

Initially, the simulated system starts with `mtime` at 0 and `mtimecmp` at `INT_MAX`, so timer interrupts are disabled. `mtime` will be modified any time it is read or written to by the simulation. When `mtimecmp` is written to, `mtime` is updated and an event (timer alarm) is scheduled for when `mtime` should be greater than or equal to `mtimecmp`. This timer alarm, when run, will verify `mtimecmp` and `mtime`, and then post a timer interrupt to the CPU. The software interrupt handler will then reset `mtimecmp`, which writes `mtimecmp` and consequently starts the timer alarm process again.

If `mtimecmp` is written to again in between the aforementioned cycle, the previous timer alarm event is descheduled and a new one is scheduled in its place. Note that even setting `mtimecmp` back to `INT_MAX` would still set a timer alarm, although at a tick extremely far into the future.



10 Miscellaneous SimpleBoard SoC Devices

10.1 RISC-V PCI Host

For access to a disk image, gem5 uses a simulated PIIX4 IDE controller connected via PCI bus to a generic PCI host model. Our SimpleBoard SoC defines a `GenericRiscvPciHost` object that inherits from the base gem5 `GenericPciHost` object, and is defined and implemented in `src/dev/riscv/pci_host.hh` and `src/dev/riscv/pci_host.cc`, respectively.

The sole purpose of the custom PCI host is to map the correct source interrupt number from the PCI host to the PLIC. In our implementation, this is the base interrupt source number added to the PCI interrupt number. The base source number for the PCI host is 0x20, and thus the source numbers for all potential PCI interrupts are 0x21, 0x22, 0x23, and 0x24, although it is observed that only source number 0x21 is used.

10.2 UART

The UART module is incorporated into the RISC-V system using the `Uart8250` model that comes with the vanilla gem5 release.

11 Future Work

In this section, we discuss the trajectory of this project and its short and long-term goals. Ideally we would like to simulate most of the features of a full-fledged RISC-V SoC so that system configurations ranging from embedded systems to high-performance computing systems can be verified and analyzed. The sections below outline some of the most crucial steps towards achieving this goal.

11.1 Formal Verification and Validation

While having a functionally correct simulation of RISC-V hardware is ideal, it is only useful if it can accurately represent real-world hardware as well. The next major milestone is to tune the latencies and models (including different CPU types) against a real RISC-V system. This will show that **GXRS** can be used to accurately represent the performance of RISC-V systems, and thus can be used to conduct design-space exploration of new RISC-V systems.

11.2 MPSoC Support

Given that even the most basic embedded systems nowadays include multiple CPU cores, it is crucial that our RISC-V simulator be able to simulate multiple cores and their interactions, including parallel workloads, shared caches, etc.

12 References and Acknowledgements

12.1 Acknowledgements

This work has been partially supported by the ERC Consolidator Grant COMPUSAPIEN (GA No. 725657), the EC H2020 WiPLASH (GA No. 863337) and the RISC-AI project from HEIG-VD/HES-SO

We would also like to thank Professor Katzalin Olcoz of the Complutense University of Madrid for her advice and help in developing the virtual memory subsystem of this project, and general advice about kernel operations.

12.2 Links Reference

- **GXR5** Homepage: <https://www.epfl.ch/labs/esl/research/2d-3d-system-on-chip/gxr5>
- Embedded Systems Laboratory (ESL) at EPFL: <https://esl.epfl.ch/>
- gem5-X Homepage: <https://www.epfl.ch/labs/esl/open-source-software-projects/gem5-x/gem5-x-documentation/>
- gem5: http://gem5.org/Main_Page
- OpenSBI: <https://github.com/riscv/opensbi>
- Linux: <https://github.com/torvalds/linux>
- buildroot: <https://buildroot.org/>
- Device Tree standard: <https://www.devicetree.org/>
- The RISC-V Foundation: <https://riscv.org/>
- The RISC-V GNU Toolchain: <https://github.com/riscv/riscv-gnu-toolchain>
- QEMU: <https://www.qemu.org/>
- spike: <https://github.com/riscv/riscv-isa-sim>



13 Glossary

13.1 Acronyms

Acronym	Term	Description
CLINT	Core-Local Interrupt Controller	Interrupt controller responsible for delegating timer and software interrupts to a local hart.
CSR	Control and State Register	RISC-V term for registers responsible for preserving system information and state.
DTB	Data Table Buffer	The L1 cache for data within the CPU.
hart	RISC-V Hardware Thread	RISC-V term hardware thread within a CPU.
ITB	Instruction Table Buffer	The L1 cache for instructions within the CPU.
MMU	Memory Management Unit	Hardware component that contains a TLB and page table walker.
PCI	Peripheral Control Interface	A standard meant to facilitate control of peripheral devices. A PCI system typically includes a host connected to a bus, connected to a series of PCI-compatible devices.
PLIC	Platform-Level Interrupt Controller	Interrupt controller responsible for delegating external interrupts to RISC-V harts.
RV64GC	RISC-V 64-bit General and Compressed extensions	RISC-V 64-bit-word-width machine with General and Compressed extensions. The General and Compressed extensions are the minimum extensions required to run a Linux-capable system.
PIO	Peripheral Input/Output	Usually used in reference to a PIO device, refers to external device such as a UART model.
Sv39	39-bit Supervisor Virtual Address	One of several RISC-V virtual addressing schemes specifying a 39-bit virtual address width.



Acronym	Term	Description
TLB	Translation Lookaside Buffer	Generic term for hardware memory cache, typically apart of the MMU.

13.2 Terms

Term	Description
Exception	A synchronous fault with respect to CPU cycles.
Fault	General term in RISC-V referring to various kinds of traps (e.g., page fault).
Interrupt	An asynchronous fault with respect to CPU cycles.

Bibliography

- [1] K. Asanovic, D. A. Patterson, "Instruction Sets Should Be Free: The Case For RISC-V", in *Technical Report No. UCB/EECS-2014-146*, August 6, 2014. Accessed September 6, 2019 at <https://people.eecs.berkeley.edu/~krste/papers/EECS-2014-146.pdf>.
- [2] F. Bellard, "QEMU, a fast and portable dynamic translator.", in *USENIX Annual Technical Conference*, 2005.
- [3] A. Waterman and Y. Lee, "Spike, a RISC-V ISA Simulator", online. Accessed April 20, 2020 at <https://github.com/riscv/riscv-isa-sim>.
- [4] Jonathan Bachrach, Huy Vo, Brian Richards, Yunsup Lee, Andrew Waterman, Rimas Avižienis, John Wawrzynek, and Krste Asanović. "Chisel: constructing hardware in a scala embedded language". in *Design Automation Conference (DAC), 2012 49th ACM/EDAC/IEEE.*, 2012, pp. 1212– 1221.
- [5] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. R. Hower, T. Krishna, S. Sardashti, R. Sen, K. Sewell, M. Shoaib, N. Vaish, M. D. Hill, D. A. Wood, "The gem5 simulator", in *ACM SIGARCH Computer Architecture News*, Volume 39 Issue 2, May 2011.
- [6] Y. M. Qureshi, W. A. Simon, M. Zapater, D. Atienza, K. Olcoz, "Gem5-X: A Gem5-Based System Level Simulation Framework to Optimize Many-Core Platforms", in *2019 Spring Simulation Conference (SpringSim)*, 29 April - 2 May 2019.
- [7] A. Waterman, K. Asanovic, *et al.*, "The RISC-V Instruction Set Manual Volume II: Privileged Architecture", online, June 8 2019. Accessed September 6, 2019 at <https://riscv.org/specifications/privileged-isa/>.
- [8] A. Waterman, K. Asanovic, *et al.*, "The RISC-V Instruction Set Manual Volume I: Unprivileged ISA", online, June 8, 2019. Accessed September 6, 2019 at <https://riscv.org/specifications/>.
- [9] A. Waterman, K. Asanovic, *et al.*, "The RISC-V Instruction Set Manual", online, June 8 2019. Accessed September 6, 2019 at <https://github.com/riscv/riscv-isa-manual>.
- [10] R. Scheffel, "Simulation of RISC-V based Systems in gem5", Master's thesis, TU Dresden, August 2018.
- [11] "Architectures/RISC-V", web page. Accessed May 20, 2020 at <https://fedoraproject.org/wiki/Architectures/RISC-V>.
- [12] "Debian port information", web page. Accessed May 20, 2020 at <https://wiki.debian.org/RISC-V>.
- [13] P. Dabbelt, K. Cheng, J. Wilson, A. Waterman, *et al.*, "RISC-V GNU Compiler Toolchain", online, September 12 2014. Accessed May 20, 2020 at <https://github.com/riscv/riscv-gnu-toolchain/graphs/contributors>.
- [14] S. Li, J. H. Ahn, R. Strong, J. Brockman, D. Tullsen, and N. Jouppi, "McPAT: an integrated power, area, and timing modeling framework for multicore and manycore architectures", in *MI-CRO 42: Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture*, December 2009, pp. 469-480.

- [15] A. Roelke and M. Stan, "RISC5: Implementing the RISC-V ISA in gem5", in *First Workshop on Computer Architecture Research with RISC-V (CARRV)*. 2017.
- [16] T. Tuan, L. Cheng, and C. Batten, "Simulating multi-core RISC-V systems in gem5.", in *Workshop on Computer Architecture Research with RISC-V*. 2018.
- [17] "RISC-V Open aSource Supervisor Binary Interface (OpenSBI)", online. Accessed May 20, 2020, at <https://github.com/riscv/opensbi>.
- [18] "Linux kernel", online. Accessed May 20, 2020, at <https://github.com/torvalds/linux>.
- [19] "Buildroot: Making Embedded Linux Easy", online. Accessed May 20, 2020, at [https:// buildroot.org/](https://buildroot.org/).
- [20] "SiFive FU540-C000 Manual v1p0", online. Accessed May 20, 2020, at <https://static.dev.sifive.com/FU540-C000-v1.0.pdf>.
- [21] Sakalis, Christos, *et al.*, "Splash-3: A properly synchronized benchmark suite for contemporary research.", in *2016 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 2016.
- [22] "SPEC CPU2017", online. Accessed May 20, 2020, at <https://www.spec.org/cpu2017/>.
- [23] "riscv-tests", online. Accessed May 20, 2020, at <https://github.com/riscv/riscv-tests>.
- [24] N. Asmussen, H. Härtig, and G. Fettweis, "A Modular and Secure System Architecture for the IoT", at the *3rd gem5 Users' Workshop*. June 2020.