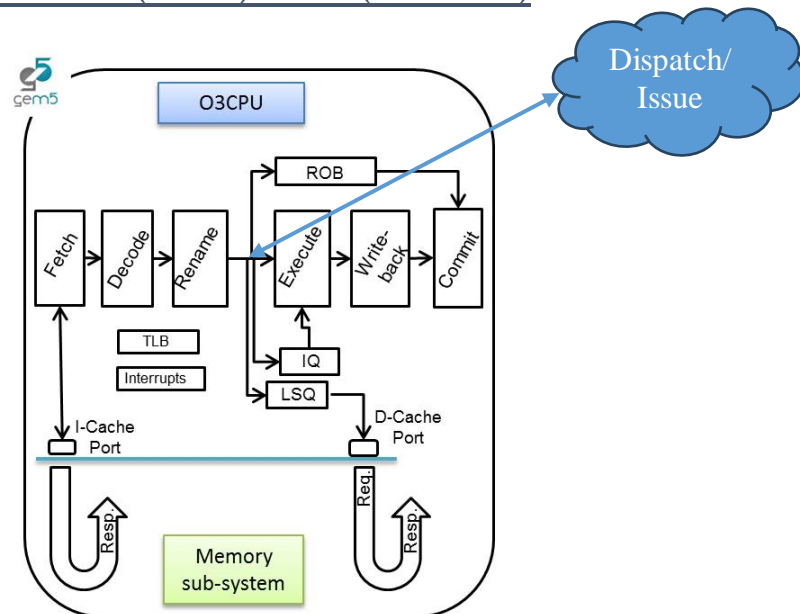


Expected delivery of **lab_4.zip** must include:

- each configuration of the custom architecture (riscv_o3_custom.py) that you modify.
- This document with all the field compiled and in PDF form.

Introduction and Background

Simulating an Out-of-Order (OoO) CPU (O3CPU)



In this laboratory, you will be able to configure an OoO CPU by using a script called `riscv_o3_custom.py`. In a few words, the script configures an Out-of-Order (O3) processor based on the *DerivO3CPU*, a superscalar processor with a reduced number of features.

Pipeline

The processor pipeline stages can be summarized as:

- **Fetch stage:** instructions are fetched from the instruction cache. The `fetchWidth` parameter sets the number of fetched instructions. This stage does branch prediction and branch target prediction.
- **Decode stage:** This stage decodes instructions and handles the execution of unconditional branches. The `decodeWidth` parameter sets the maximum number of instructions processed per clock cycle.
- **Rename stage:** As suggested by the name, registers are renamed, and the instruction is pushed to the IEW (Issue/Execute/Write Back) stage. It checks that the *Instruction Queue (IQ)*/*Load and Store Queue (LSQ)* can hold the new instruction. The maximum number of instructions processed per clock cycle is set by the `renameWidth` parameter.

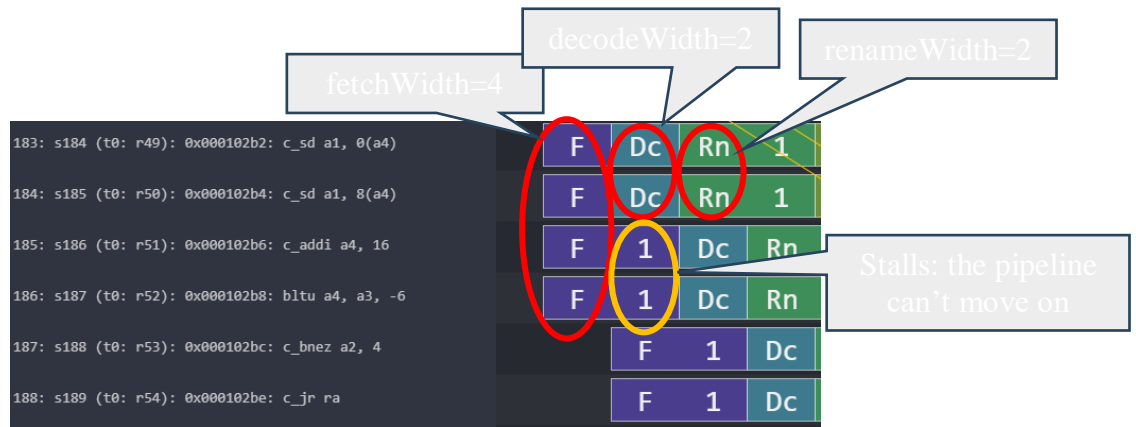


Figure 1: Understanding configurable OoO CPU parameters.

- **Dispatch stage:** instructions whose renamed operands are available are dispatched to functional units (FU). For loads and stores, they are dispatched to the Load/Store Queue (LSQ). The maximum number of instructions processed per clock cycle is set by the dispatchWidth parameter.
- **Issue stage:** The simulated processor has a single instruction queue from which all instructions are issued. Ordinarily, instructions are taken in-order from this queue. An instruction is issued if it does not have any dependency.
- **Execute stage:** the functional unit (FU) processes their instruction. Each functional unit can be configured with a different latency. Conditional branch mispredictions are identified here. The maximum number of instructions processed per clock cycle depends on the different functional units configured and their latencies.
- **Writeback stage:** it sends the result of the instruction to the reorder buffer (ROB). The maximum number of instructions processed per clock cycle is set by the wbWidth parameter.
- **Commit stage:** it processes the reorder buffer, freeing up reorder buffer entries. The maximum number of instructions processed per clock cycle is set by the commitWidth parameter. Commit is done in order.

In the event of a **branch misprediction**, trap, or other speculative execution event, "squashing" can occur at all stages of this pipeline. When a pending instruction is squashed, it is removed from the instruction queues, reorder buffers, requests to the instruction cache, etc.

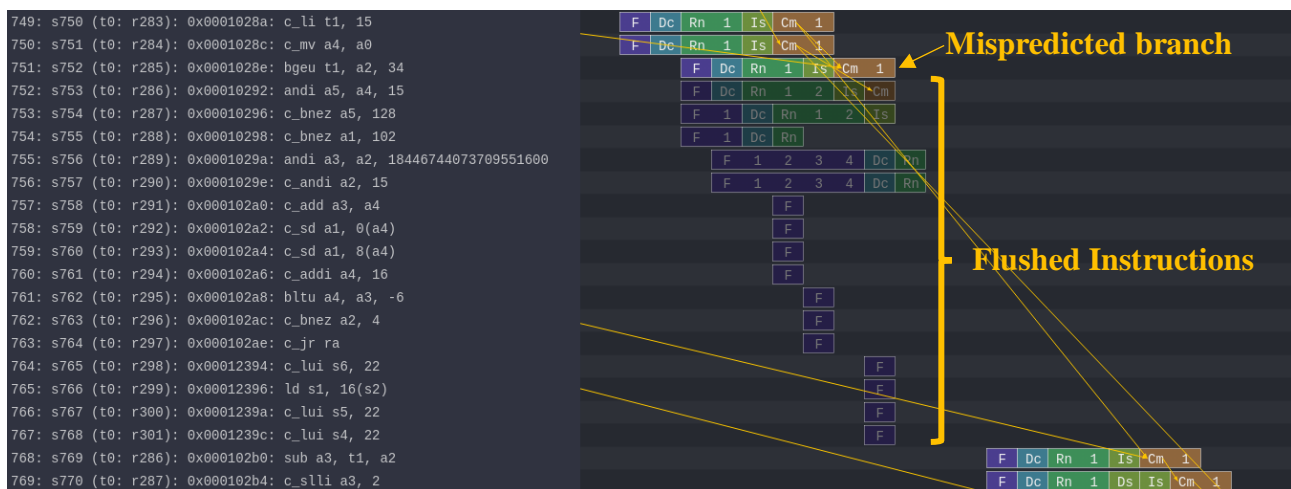


Figure 2: Example of a branch *misprediction* (transparent rows)

Pipeline Resources

Additionally, it has the following structures:

- Branch predictor (BP)
 - Allows for selection between several branch predictors, including a local predictor, a global predictor, and a tournament predictor. Also has a branch target buffer (BTB) and a return address stack (RAS).
- Reorder buffer (ROB)
 - Holds instructions that have reached the back end. Handles squashing instructions and keep instructions in program order.
- Instruction queue (IQ)
 - Handles dependencies between instructions and scheduling ready instructions. Uses the **memory dependence predictor** to tell when memory operations are ready.
- Load-store queue (LSQ)
 - Holds loads and stores that have reached the back end. It hooks up to the d-cache and initiates accesses to the memory system once memory operations have been issued and executed. Also handles forwarding from stores to loads, replaying memory operations if the memory system is blocked, and detecting memory ordering violations.
- Functional units (FU)
 - Provides timing for instruction execution. Used to determine the latency of an instruction executing, as well as what instructions can issue each cycle.
 - **Floating point units, floating point registers,** and respective instructions are supported.

560: s561 (t0: r160): 0x00010106: fmv_w_x fa5, zero	F	Dc	Rn	1	Is	1	2	3	Cm	1	
561: s562 (t0: r161): 0x0001010a: c_addi16sp sp, -64	F	Dc	Rn	1	Is	Cm	1	2	3	4	
562: s563 (t0: r162): 0x0001010c: c_fsdsp fs0, 0(sp)	F	1	Dc	Rn	1	Is	Mc	1	2	3	4
563: s564 (t0: r163): 0x0001010e: c_fsdsp fs1, 0(sp)	F	1	Dc	Rn	1	2	3	Is	Mc	1	2

Figure 3: Pipeline example of FP instructions and FP registers

Laboratory: hands-on

All the needed resources are at a GitHub repository:

https://github.com/cad-polito-it/ase_riscv_gem5_sim

To create your simulation environment:

For HTTPS clone:

```
~/my_gem5Dir$ git clone https://github.com/cad-polito-it/ase_riscv_gem5_sim.git
```

For SSH:

```
~/my_gem5Dir$ git clone git@github.com:cad-polito-it/ase_riscv_gem5_sim.git
```

The environment is configured to be executed on the **LABINF MACHINES**.

Follow the HOWTO instructions available on the GitHub Repository for simulating a program.


Exercise 1:

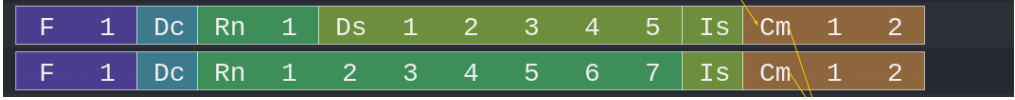
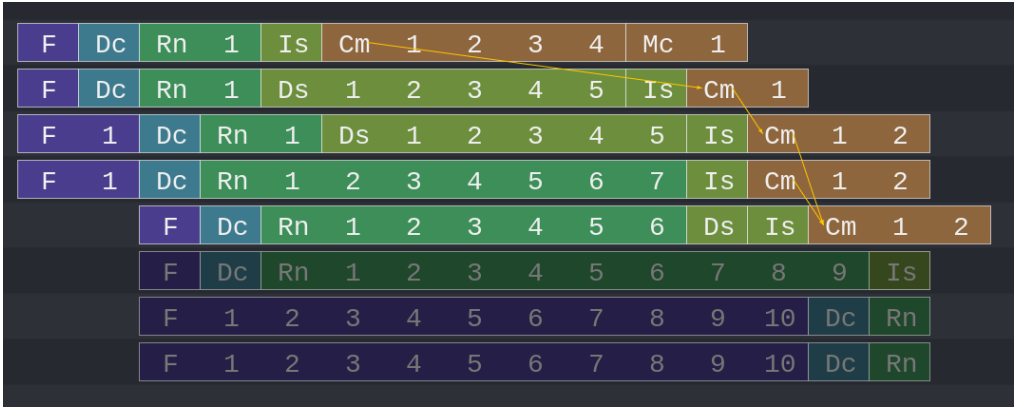
Simulate the benchmark *my_c_benchmark* (*main.c*) by using the gem5 simulator to obtain the *trace.out* file. Then, you can visualize the pipeline (i.e., load the *trace.out* file on Konata).

Based on the CPU architecture described in *riscv_o3_custom.py*, visualize the Konata's pipeline to find out the conditions:

1. Out-of-order execution (issue), in-order commit (commit)
2. Two commits in the same clock cycle
3. Flush of the pipeline.

For every condition, fill the following tables.

Condition	Out-of-order execution, in-order commit
Screenshot from Konata	
Explain the reason behind the condition	Even though the instructions can be executed in a different order to maximize performance, the CPU commits the results of these instructions in the original, sequential order of the program. This happens in order to ensure consistency with the program's logical flow and, in case of exceptions or interrupts during execution, committing results in order simplifies error handling.
Briefly explain the advantages of the OoO execution in a CPU	Out-of-Order (OoO) execution allows the processor to identify and execute independent instructions simultaneously, maximizing the utilization of the functional units. This results in better performance by reducing the idle time of the processor. Moreover, it enables the CPU to rearrange the order of instructions on-the-fly, mitigating stalls caused by dependencies or resource conflicts. This leads to more efficient use of resources and faster completion of tasks.

Condition	Two or more commits in the same clock cycle
Screenshot from Konata	
Explain the reason behind the condition	When several instructions are executed out of order and their dependencies are resolved simultaneously, their results might be ready for commit at the same time. If these instructions are independent of each other and their execution completes together, they can be committed in a single clock cycle.
Briefly explain the Commit functioning	Completed instructions wait in the Reorder Buffer (ROB) until they are ready to be committed. The Commit stage validates these instructions to ensure they were executed correctly without errors or exceptions. After successfully committing an instruction, the corresponding entry in the ROB can be freed. This action makes space in the ROB for new instructions to be tracked as they progress through the pipeline.
Condition	Flush of the pipeline
Screenshot from Konata	
Explain the reason behind the condition	When the processor encounters a conditional branch, it might predict the direction the code will take. If that prediction turns out to be wrong, the pipeline is flushed to discard the instructions that were fetched and executed based on the incorrect prediction. This ensures the processor resumes execution from the correct point in the program. In addition, if an exception or error occurs during the execution of an instruction (like a divide-by-zero error), the processor needs to stop the current instructions and handle the exception. Flushing the pipeline removes any instructions that might cause incorrect results due to the exception, allowing the processor to handle the error appropriately.

Exercise 2:

Given your benchmark (*main.c* in *my_c_benchmark*), optimize the CPU architecture (i.e., modify the *riscv_o3_custom.py* file) and write down the improvements in terms of CPI and speedup.

- To optimize the CPU architecture, open the configuration file of the CPU (i.e., the *riscv_o3_custom.py*), and tune specific hardware-related parameters.

You have to change specific values in **one or more** stages of the pipeline:

- # - FETCH STAGE
 - Tune parameters such as the *fetchWidth*, *fetchBufferSize* and so on, and see the effects on your system.
- # - DECODE STAGE
- # - RENAME STAGE
 - Try changing some values, but don't touch the "Phys" ones.
- # - DISPATCH/ISSUE STAGE
- # - EXECUTE STAGE
 - Here you can optimize the Functional units of your CPU like the INT ALU, the FP ALU, the FP Multiplier/Divider and so on.
 - Tune the number of units (*count*) that you have in the system, as well as their latency (*opLat*) to see how this affects the execution of your program.
- You can create a different branch predictor. They are defined in *create_predictor.py*
- You can also try to change the parameters of the L1 Cache. Look for the "class L1Cache" in the *riscv_o3_custom.py* file. The L1 cache, also referred to as the primary cache, is the smallest and fastest level of memory. It is located directly on the processor, and it is used to store frequently accessed data by the CPU. In this way, the CPU saves time with respect to the normal access to the main memory.

HINT: To implement the best hardware optimization, and understand how to change the parameters, the best option consists in analysing the *stats.txt* file (in *ase_riscv_gem5_sim/results/my_c_benchmark*). Find information regarding the workload profiling. In other words, look for lines such as "system.cpu.commitStats0.committedInstType::IntAlu", and the following ones to understand which kind of instructions are executed the most. In this way, you can target a specific functional unit and modify its specifications.

Fill the following Tables with the CPI that you obtain with the old and the new architectures. Compute also the equivalent speedup that you obtain.

HINT: You can get the CPI and other useful information from the *stats.txt* file.

Parameters	Configuration 1	Configuration 2	Configuration 4	Configuration 5
First changed parameter	the_cpu.fetchWidth = 2	the_cpu.fetchWidth = 4	the_cpu.fetchWidth = 4	the_cpu.issueWidth = 12
Second changed parameter	the_cpu.dispatchWidth = 2	the_cpu.decodeWidth = 4	the_cpu.decodeWidth = 4	the_cpu.numIQEntries = 12
Third changed parameter		the_cpu.numIQEntries = 12	the_cpu.numIQEntries = 12	the_cpu.commitWidth = 12
Fourth changed parameter			the_cpu.renameWidth = 12	the_cpu.wbWidth = 12

Original CPI (no hardware optimization): 2.083105 CPI

	Configuration 1	Configuration 2	Configuration 4	Configuration 5
CPI	2.081982	1.010357	1.068256	0.998003

Speedup Original CPI (wrt)	1,000539	2,061751	1,950005	2,087273
---------------------------------------	----------	----------	----------	----------

Which is the best optimization in terms of CPI and speedup, why?

Your answer:

The best optimization in terms of CPI and speedup is the configuration number 5. In particular, by analyzing the *stats.txt* file, I have noticed that lots of cycles were being used for arithmetic integer operations and memory phase.

I tried to increase the number of arithmetic functional units and also the amount of L1 cache available, however this tuning didn't change the CPI.

Parameters explanation:

1. **cpu.issueWidth:** represents the number of instructions that the CPU can issue or dispatch simultaneously to the execution units in a single clock cycle. A higher issue width allows for more instructions to be fetched and executed concurrently, potentially increasing the CPU's performance by improving instruction-level parallelism.
2. **cpu.numIQEntries:** it refers to the number of entries in the Instruction Queue (IQ). The IQ holds instructions that have been fetched and are waiting for their operands or are awaiting execution. A larger IQ capacity can enable the CPU to handle more out-of-order execution and better manage instruction dependencies, reducing stalls and improving overall performance.
3. **cpu.commitWidth:** This parameter denotes the number of instructions that can be committed or finalized in the CPU in a single clock cycle. Committing instructions involves updating the architectural state of the CPU with the results of executed instructions. A higher commit width allows for more completed instructions to be integrated into the CPU's architectural state concurrently, potentially improving throughput.
4. **cpu.wbWidth:** The Write-Back Width specifies the number of instructions whose results can be written back to the register file or memory in a single clock cycle. It represents the capacity of the write-back stage of the CPU. A larger wbWidth can enhance the efficiency of the write-back phase by allowing more instructions' results to be updated in the CPU's architectural state simultaneously, potentially improving overall performance.