The Riscy Processor

Karan Bavishi, Mark Mansi, Suhas Pai, Preyas Shah December 22, 2016

1 Introduction

We implement the RV64I set of instructions in the RISCV ISA[1] with a superscalar, out-of-order core. Our implementation is optimized for correctness. Our second goal was performance. Our third goal was synthesize-ability.

In Section 2, we provide a brief overview of our design. In Section 3, we talk about the design of each individual module in detail. Section 4 describes some of the policies we designed and implemented to boost the performance of our processor. In Section 5, we discuss some interesting design decisions we took and the rationale behind it. In Section 6, we talk about a few challenges that we ran into, which forced us to rethink our design. Finally, in Section 7, we discuss the performance achieved by our processor for certain benchmarks.

2 Overview

- 4-wide pipeline
- Out-of-order execution; In-order commit
- Fetch
 - 3 cycle latency for cache hit (2 cycle I-cache latency)
 - Branch prediction (TODO): 2-level predictor, BTB, global history register.
 - I-Cache: Blocking, 4 KB, 2-way associative, 1 read/write port.
 - Prefetch Buffer: 4 entries, 32 bytes each
- Decode and Rename
 - 64-entry ROB
 - 2 cycle latency
- Issue and FUs
 - 4 ALUs

- 4 issue queues, 16 entries each, 1 ALU per queue
- Load-balancing arbiter places new instructions in issue queues.
- ALUs also used to compute ld/st addresses.

• Address Queue

- 32 entries
- Tracks load/store dependencies.
- Non-speculative memory disambiguation.
- Loads access cache out-of-order between stores.
- Stores access cache on commit.
- D-Cache (TODO)

• Writeback

- Processor supports back-to-back execution of dependent instructions.
- Writeback structure (a.k.a ROB WB or more affectionately, FooPP) is designed to avoid the massive tangle of wires created by broadcast-based writeback among 4 ALUs and a LSQ.
- Stall: we implement a distinct top-level module to handle stall coordination among all stages. The two stall producers in the pipeline are the issue queue arbiter and the ROB. A stall is generated when there is not enough room in the issue queues or the ROB. Each stage is a consumer of stalls produced by stages before it.

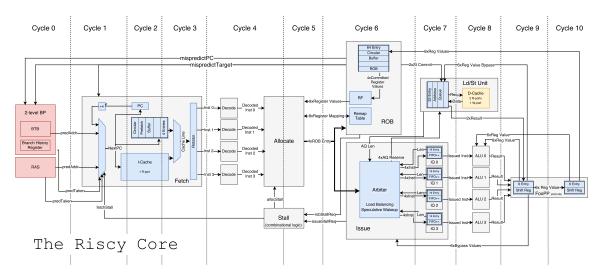


Figure 1: Top-level diagram of our processor

A top-level diagram of our design is shown in Figure 1. A bigger version of the diagram is also available for view[2].

3 Modules

Our processor has a pipeline depth of 8. The pipeline is divided into several independently designed and implemented blocks (implemented as modules in Chisel). This section presents an overview of the whole processor. Then, we present the design of each block.

3.1 Fetch

Our fetch module encapsulates the logic of selecting the next PC, issuing requests to the instruction memory hierarchy (including the I-cache), and presenting instructions to the subsequent stage (decode).

We implement a pipelined, blocking, 2-way associative instruction cache with 2-cycle latency. And additional cycle of latency comes from computing the next PC.

The instruction cache returns a hit if it finds the requested entry. If not, it issues a refill request to memory for the missing address. Once the memory responds, it refills its data and tag arrays with the response. The cache line replacement policy is random replacement.

We also added a prefetch buffer, containing 4 entries. It does cache block prefetching ie. it stores the additional cache lines sent by memory as a part of its response to a refill request. This is primarily done to get around the single-port limitation of the I-cache. This prefetch buffer can be easily changed to do next-line prefetching as a part of future work.

The output from the instruction cache is rotated and presented to the rest of the pipeline.

3.2 Decode and Allocate

The next two cycles are consumed in decoding and renaming the instructions from fetch.

Decode takes less than one cycle because it consists entirely of separating wires. Thus, we combine decode with the first cycle of renaming.

Renaming proceeds in two phases. The first cycle decodes the instruction and renames all destinations. The second cycle renames operands and produces and ROB entry which can just be latched to the ROB next cycle.

3.3 ROB

The ROB is implemented as a 64-entry circular buffer. All committing logic is implemented in the ROB module. 4 instructions can commit in a single cycle. In the same cycle, 4 new ROB entries can also be added to the tail of the ROB.

We choose to use an ROB as opposed to a merged register file because we prioritized performance and ease of implementation over energy.

3.4 Issue and FUs

In the same cycle, an ROB entry both latches in the ROB and goes to the issue queue arbiter, which puts it in a queue. All instructions go through the arbiter regardless of their type.

Our issue stage consists of an arbiter feeding 4 issue queues. The arbiter does load balancing to attempt to keep any queue from becoming filled while other ALUs are available.

Each issue queue is independent and feeds a single ALU. The issue queues can each independently wake up a single instruction to be sent an ALU. Our implementation is capable of speculatively waking up ALU operations when their operands are known to be nearly ready. This allows our processor to execute instruction with back-to-back dependencies. Our writeback mechanism helps with this also...

3.5 Writeback

Each cycle, our processor may produce up to 6 values that need to be written back (4 from ALUs and 2 from memory). To avoid creating $O(n^2)$ wires between all of these value producers/consumers, we designed a single writeback structure, which we call FooPP. All values are written back to the FooPP in the cycle they produced. Any consumers of those values, including the ROB and issue queues, take the values directly from the FooPP.

To avoid adding a 1 cycle bubble between instructions with back-to-back dependencies, the FooPP keeps 1 cycle of history, which the FUs can then consume from easily.

The writeback structure also filters out any values coming from squashed instructions. It does this by comparing the *era* of the input values with the *current era* according to the ROB. This allows us to get rid of the squashed instructions easily by letting them percolate through the pipeline.

3.6 Address Queue

When a load or store instruction reaches the arbiter, the arbiter asks the address queue to reserve a spot for it. If the address queue is full, the pipeline stalls. The arbiter then puts the load or store in a normal issue queue to wait until an ALU can compute its address.

When the address is written to the Writeback structure, the address queue will consume it.

Loads can access the D-cache as soon as their address is known and it is also known that there is no preceding store to the same address.

Stores can access the D-cache as soon as the ROB signals the address queue that they have committed.

To limit the number of ports need on the D-cache, we instead assume only two ports. The ROB gives the guarantee that no more than 2 stores will commit in any cycle.

4 Performance Enhancement Policies

In this section, we talk about some of the design decisions we took which were aimed at improving the performance of our processor.

- 1. Fetch We added a prefetch buffer which stored cache blocks when a refill response was received from memory. This helped us get around the limitations of having a single-ported I-cache, which allowed us to only write one cache line when the refill response was received. Having a prefetch buffer allowed us to store multiple cache lines.
 - The prefetch buffer also helped improve performance by having enough instructions to keep the backend busy. Without the prefetch buffer, we were only able to serve instructions every 2 cycles before incurring another hit.
- 2. Decode & Allocate Both Decode and Allocate stages were done in the same pipeline stage to avoid wasting a cycle
- 3. Issue & ROB Renamed instructions are sent to both the ROB and Issue modules in parallel to save a pipeline stage. We also added support for speculative wakeup of instructions whose operands would be ready in the subsequent cycle via the ALU output. This allowed us to perform back-to-back execution for dependent instructions and boost performance.
- 4. Execute Apart from sending the ALU results to the ROB, the Writeback module was also responsible for storing the results for 2 cycles for bypass. This bypass support allowed us to execute dependent instructions in backto-back cycles.
- 5. ROB mis-prediction eras Incase of a mis-prediction, there are two ways to squash the speculated instructions. The first involves adding $O(n^2)$ valid bit wires everywhere, and resetting them when a mis-prediction occurs. The second way is to let the squashed instructions percolate through the pipeline.

The problem with the latter approach is that the ROB can not let new instructions pass through unless all the squashed instructions have percolated. This is because we can end up having new and squashed instructions with the same ROB tag. In order to get around this issue, we add misprediction eras. Whenever ROB realizes a misprediction, it increments the era counter. This era can help us distinguish between new instructions and to-be squashed instructions.

5 Design Decisions

In this section, we talk about the interesting design choices we made for our processor.

- 1. Pipelined Decode & Rename In our Decode and Allocate design, we decode and rename the destination in the first cycle. Renaming of operands is done in the second cycle. This allows us to allocate new entries in the ROB in the second cycle, in parallel with the renaming of operands.
- 2. Use Writeback to flush values from squashed instructions As mentioned in Section 4, we decide to let the instructions to be squashed percolate through the pipeline. An interesting question was to how to filter out the results from these squashed instructions and not make them be mistakenly used for bypass. We ultimately decided to add some filter logic to the Writeback module inputs. This made sure that the squashed values were invalidated and thus never available for bypass.
- 3. Arbiter load-balancing and bypass Our original idea was to have the arbiter build dependency graphs and put dependent instructions in the same queue. The rationale was that putting dependent instructions in the same queue would reduce the overhead of bypass across queues. However we soon realized that this plan would result in a multi-cycle arbiter design. We ultimately decided to have a simple arbiter which would issue instructions based on load balancing, i.e. instructions would be issued based on the emptiness of the queues. This allowed us to design a single-cycle arbiter. However, this brought back the original problem of complex bypass. For this, we designed the Writeback structure which would allow bypass of computed values to the 4 ALUs.
- 4. LSQ: Use Issue queue wakeup logic for load/store instructions Our original plan was to have 4 issue queues and 1 load-store queue, using 4 ALUs and 1 AGU. But we decided to re-use the existing ALUs for address computation for loads and stores. This would allow us to re-use the complex instruction wakeup logic already present in the Issue Queues, and not re-implement from scratch.
- 5. LSQ: Serialize loads and stores so that D-cache can have 1 read port and 2 write ports
- 6. Unified Stall logic We decided to come up with a unified stall module which would take stall requests as input from ROB and Arbiter modules and generate stall outputs to the consumers such as the Fetch and the I-cache modules. This unified logic simplified the interface between all these modules, and avoided complex wiring.

6 Design Challenges

In this section, we talk about certain challenges that we encountered which forced us to rethink our design.

- 1. Cancelling memory requests due to mispredictions After adding misspeculation recovery support, we found that our processor had a 5-cycle bubble before instructions were fetched from the new target. On investigating more, we realized that it was because the I-cache was busy requesting an address from memory that would no longer be needed. To get around this, we added support for cancellation of requests to memory incase of mispredictions.
- 2. ROB deadlock issues During testing, we saw the ROB running into weird deadlock issues. This forced us to rethink the question what does it mean to stall for the ROB? We realized that even if the ROB is stalled, it must continue to commit instructions to avoid deadlocks.
 - This helped solve our original problem where we saw our processor being deadlocked when the ROB stalled everything as soon as a stall request was issued by the LSQ via the Arbiter when it got full.
- 3. Timing issues because of multiple fan-ins to Issue Queues We had to carefully design the Issue Queue logic to ensure robustness to timing issues which were caused because of multiple fan-ins. If not done carefully, this could result in weird situations where the instructions would simply bounce off and not get stored in the Issue Queues.
- 4. Zombie instructions eating up space While adding support for misspeculation recovery through percolation, we realized that the squashed instructions (or zombie instructions) could end up staying much longer in the Issue Queue because of operands not being ready. Consider the case where a zombie load instruction is kept waiting for 100 cycles waiting for the memory to respond to the request. This can happen because in our percolation policy, we execute the zombie instructions but decide to filter the results.

In order to get around this, we decided to add a special priority for zombie instructions in the instruction wakeup logic. This meant that zombie instructions lying in the Issue Queue were given the highest priority for wakeup, and thus would be percolated as soon as possible.

7 Performance

7.1 Methodology

We ran a series of workloads to show the upper and lower bounds of IPC that our processor could achieve. We ran the workloads in two main settings to test the impact of lack of a data-cache and data-prefetching:

- 7 cycle memory latency
- 4 cycle memory latency (smallest possible latency supported by our design)

Additionally, we also ran the same workloads with instruction prefetching turned on. In Section 7.3, we talk about the results achieved with the two abovementioned settings, and along with instruction-prefetching wherever correctness of the benchmarks was verified.

7.2 Benchmarks

We ran multiple workloads which were aimed at highlighting the lower and upper bounds of the IPC that could be achieved by our processor. In the following table, we list the benchmarks that were used:

Benchmark	ILP	Instructions	Description		
Counters	High ILP	8K insts	32 counters incremented inde-		
(optimal)			pendently.		
Fibonacci	Max ILP 1.5	150 insts	Compute the first 64 Fibonacci		
			numbers.		
Daxpy	High ILP	3K insts	Element-wise addition of arrays.		
	(50 % mem-				
	ory insts)				
Stall (worst	Low ILP	160 insts	Back-to-back dependent instruc-		
case)			tions force in-order execution		

We describe the benchmarks and the reasons for choosing them in more detail below:

- 1. Counters In this benchmark, we increment 32 counters independently. The reason for choosing this benchmark is to have as many independent instructions as possible and thus find out the upper limit of ILP that can be achieved by our design.
- 2. Fibonacci The Fibonacci benchmark is aimed at testing our design with a purely compute-based workload. Due to the dependencies in the Fibonacci recurrence relation, the maximum ILP that can be achieved is 1.5.
- 3. Daxpy We use a loop-unrolled version of the DAXPY benchmark. Since the original version involves multiply instructions which are not supported in our design, we multiply by 1 which is basically a no-op. This benchmark has a high ILP workload with a high memory-to-compute ratio of instructions, aimed at stressing the LSQ part of our design.
- 4. Stall In this benchmark, we execute 160 back-to-back dependent instructions. The maximum ILP that can be achieved is 1. The rationale behind choosing this to verify that our design can execute dependent instructions in back-to-back cycles with no bubbles.

7.3 Results

In the following table, we describe the ILPs achieved for all the benchmarks in the 2 chosen memory latency settings, with instruction prefetching enabled and disabled. We omit results for settings where we were unable to verify correctness.

Setting	Counters	Fibonacci	Daxpy	Stall
7 cycle, I-prefetch disabled	0.80	0.74	0.80	0.74
7-cycle, I-prefetch enabled	2.22	1.30	-	0.91
4 cycle, I-prefetch disabled	1.14	1.04	0.88	0.92
4 cycle, I-prefetch enabled	2.66	-	-	-

The following interesting observations can be made:

1. Instruction prefetching makes a huge difference - As can be seen from the ILP numbers achieved by our best-case Counters benchmark, it is hard to achieve a high ILP unless you have enough instructions to keep the backend busy. With instruction prefetching disabled, our processor incurs the memory latency penalty every 2 cycles.

Interestingly, even with prefetching enabled, our design is only able to reach an ILP of 2.66. The reason for this is that our prefetching design is fairly naive – it is not a next-line prefetcher. The prefetch unit does not fetch instructions well in advance and only stores additional cache lines when a penalty is incurred. If we increase the size of the prefetch buffer from 4 to 8, our projections show that we can achieve an ILP of about 3.2.

Having a more intelligent prefetcher and increasing the size of the prefetch buffer should help boost the ILP close to the maximum theoretical limit of 4.

- 2. Important to have a big instruction window Similar to our previous observation, it is extremely important to have a big instruction window to exploit the techniques implemented in an out-of-order design. Without enough instructions in the frontend, it is difficult to find enough parallelism and the processor remains hampered by dependencies and memory latencies.
- 3. Memory latency can significantly degrade the ILP As can be seen from the ILP numbers achieved by the Daxpy benchmark, it is hard to achieve a high ILP unless the memory latency can be hidden effectively.
 - The Daxpy benchmark has an inherently high-ILP workload, but since our design lacks a data-cache, it is unable to go beyond executing more than one instruction per cycle.
- 4. Close to optimal performance in Fibonacci We achieve an ILP of 1.30 for the 7-cycle, prefetching enabled setting, which is quite close to the theoretical limit of 1.5. We can actually reach an even higher figure of 1.42

- for the 4-cycle, prefetching enabled setting, but we decided not to show it here because we couldn't get the correct result from the benchmark.
- 5. Optimal performance in Stall Our Stall benchmark results are very close to the optimal ILP of 1.0. This highlights our processor's ability to execute dependent instructions in bac-to-back cycles with no bubbles.

References

- [1] RISCV ISA https://riscv.org/specifications/.
- [2] RISCY Core diagram https://github.com/mark-i-m/riscy/raw/master/riscy/doc/riscy_diagram.pdf