Who tests the TestRIG? Tooling for randomised tandem verification

Peter Rugg, Alexandre Joannou, Jonathan Woodruff, Franz A. Fuchs, and Simon W. Moore Department of Computer Science and Technology, University of Cambridge peter.rugg@cl.cam.ac.uk Project URL: cheri-cpu.org



TestRIG

TestRIG is an ecosystem for cross-verifying RISC-V implementations using a standard **RVFI-DII** interface. Verification Engines connect to the implementations over this interface. **QuickCheckVEngine** uses Haskell's QuickCheck library to **generate tests** and **automatically shrink** any divergences to a minimal reproducer. The RISC-V golden Sail model implements RVFI-DII, allowing implementations to be compared against this (hopefully) correct-by-definition executable simulator. Since initial publication, the TestRIG infrastructure has seen increasing community engagement, including users and contributors from **Microsoft Research**, **IowRISC**, and **SCI Semiconductor**. The repository now links to 10 RVFI-DII-extended implementations, and has several forks from other members of the community. TestRIG is in use to test **CHERI**, which adds unforgeable hardware capabilities for memory safety and compartmentalisation, in the **Toooba** and **CVA6** processors. The RISC-V community is in the process of **standardising** CHERI extensions.

How can we tell the

generated tests are

testing what we want

What if we have

encoding mistakes

in the generators?

them to?



Measuring Coverage

Attempt: code coverage

Measuring lines of code run is a good start, but has the problem that code may be run, but have **no effect** on the tested outputs, causing coverage to pass but **bugs to be missed**. For example, CHERI checks may get run on capabilities where the integrity tag is already clear, **hiding the error**.

Solution: mutation coverage

- Simulate real bugs by modifying the Sail code with incorrect behaviour
- **Definitively** answer: "would the test framework have caught that bug?"
- Automatically try classes of **common mistakes**

Implementation

- Python classes describe mutation types: pattern and transformation
- Coverpoints and run results are tracked in a SQLite Database
- Support for building and running many mutants in **parallel**
- Produces pretty HTML reports to highlight coverage issues

What if the processor never reaches an "interesting" state?

What if the results are thrown away before they are compared?



Minimal, single instruction annotated reproducer

Future Work

- Switch from Python text transformations to directly interacting with the **Sail compiler**
- Automatically inline Sail functions
- Extend supported mutations
- Compare effectiveness of different test suites
- Identify and close coverage **blind-spots**
- Extend to **microarchitectural** mutation coverage in a Verilog implementation
- Find lots of **bugs**?

Bonus: relaxing the "tandem"

Aligning implementation and model on every possible instruction input is hard! CSR legalisation, store conditional failures, misaligned accesses,...

TestRIG now supports a "single implementation" mode: run implementation on its own and assert more liberal properties over RVFI trace.

Example property: instruction count in = instruction count out, i.e. processor did not **lockup**. Template to inject **undirected** random instructions found and diagnosed:

- Processor mis-decoded and locked up on certain illegal instructions
- Subtle and rare compressed branch mispredict infinite loop
- Reachable fatal assert condition in a version of the Sail model

DII is required to allow such liberal testing without having to reason about control flow.



SRI International[®]