



HAL
open science

From low-level fault modeling (of a pipeline attack) to a proven hardening scheme

Sébastien Michelland, Christophe Deleuze, Laure Gonnord

► To cite this version:

Sébastien Michelland, Christophe Deleuze, Laure Gonnord. From low-level fault modeling (of a pipeline attack) to a proven hardening scheme. *Compiler Construction (CC'24)*, Mar 2024, Edinburgh (Scotland), United Kingdom. 10.1145/3640537.3641570 . hal-04438994v2

HAL Id: hal-04438994

<https://hal.science/hal-04438994v2>

Submitted on 28 Jun 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - ShareAlike 4.0 International License

From Low-Level Fault Modeling (of a Pipeline Attack) to a Proven Hardening Scheme*

Sébastien Michelland

sebastien.michelland@lcis.grenoble-
inp.fr
UGA, Grenoble INP, LCIS
Valence, France

Christophe Deleuze

christophe.deleuze@grenoble-inp.fr
UGA, Grenoble INP, LCIS
Valence, France

Laure Gonnord

laure.gonnord@grenoble-inp.fr
UGA, Grenoble INP, LCIS
Valence, France

Abstract

Fault attacks present unique safety and security challenges that require dedicated countermeasures, even for bug-free programs. Models of these complex attacks are made workable by approximating their effects to a suitable level of abstraction. The common practice of targeting the Instruction Set Architecture (ISA) level isn't ideal because it discards important micro-architectural information, leading to weaker security guarantees. Conversely, including micro-architectural details makes countermeasures harder to model and reason about, creating a new challenge in validating and trusting protections.

We show that a semantic approach to modeling faults makes micro-architectural models workable, and enables precise cooperation between software and hardware in the design of countermeasures. We demonstrate the approach by designing and implementing a compiler/hardware countermeasure, which protects against a state-of-the-art pipeline fetch attack that generalizes multi-fault instruction skips. Crucially, we provide a formal security proof that guarantees faults are detected by the end of every basic block. This result shows that carefully embracing the complexity of low-level systems enables finer, more secure countermeasures.

CCS Concepts: • **Software and its engineering** → *Semantics*; **Software notations and tools**; • **Security and privacy** → **Formal methods and theory of security**.

Keywords: Compilation, semantics of fault models, security proofs, co-design

ACM Reference Format:

Sébastien Michelland, Christophe Deleuze, and Laure Gonnord. 2024. From Low-Level Fault Modeling (of a Pipeline Attack) to a

*This work has been partially funded by the ARSENE project of the “France 2030” government investment plan managed by the French National Research Agency, under the reference ANR-22-PECY-0004.

CC '24, March 2–3, 2024, Edinburgh, United Kingdom

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of the 33rd ACM SIGPLAN International Conference on Compiler Construction (CC '24)*, March 2–3, 2024, Edinburgh, United Kingdom, <https://doi.org/10.1145/3640537.3641570>.

Proven Hardening Scheme. In *Proceedings of the 33rd ACM SIGPLAN International Conference on Compiler Construction (CC '24)*, March 2–3, 2024, Edinburgh, United Kingdom. ACM, New York, NY, USA, 23 pages. <https://doi.org/10.1145/3640537.3641570>

1 Introduction

An attacker with access to a physical device can perform *fault injection attacks*. Physical interference such as a clock glitch, a power supply voltage glitch, or an electromagnetic pulse, can cause hardware to behave erroneously [4], sometimes just enough to bypass an application's security. The development of fault injection attacks [30] makes them a tangible threat to modern safety- and security-critical systems. Countering them is uniquely challenging due to the unpredictable effects of low-level interference on high-level security properties — a leap that traditional development tools meticulously avoid by building upon a clean abstraction stack from hardware to programming languages.

In order to conquer the complexity of these attacks, security engineers construct *fault models* by approximating faults' effects to a desired level of abstraction. These span from bit flips in RTL (Register Transfer Level) latches [32] to failures in pipeline forwarding [21] to corrupted ISA registers [7] and branch inversion directly in source code [28]. Countermeasures are then based on these models, so in a sense secure programs resist *fault models* rather than *faults*. The clear trade-off is one of accuracy versus simplicity; low-level descriptions are more true to practical attacks, but high-level approximations make it practical (in many cases *possible*) to reason about and protect against them.

In practice, most existing works study faults at the ISA level, based on mis-executions of assembler programs (instruction skips, wrong jumps, corrupted registers, etc. [18]), with countermeasures as *transformations* of assembler programs. This is a natural choice as assembler is the lowest software abstraction, and dealing with software has benefits such as ease of deployment, board-independence, compiler automation, and the ability to protect only critical sections of programs (compared to fixed costs in e.g. die surface). Hardware protections [19] are less common, but better equipped to deal with local and remote side-channel attacks [31], which share many aspects with fault attacks (see f.i. [35]).

The key issue with ISA-level fault models is that the approximation is quite crude; [22] shows that faulted behaviors often depend on micro-architectural features and cannot be described accurately without including hardware details. Pipeline analysis in [37] further shows that targeted fault attacks can and do defeat many ISA-level countermeasures by exploiting unmodeled low-level effects.

Naturally, using low-level models widens the *abstraction gap* between the attack and the countermeasure (often applied during compilation at an IR or back-end level). This creates a risk that protections could be altered or defeated by the compiler’s late stages. These cross-layer concerns (commonly avoided by disabling optimizations or basing security claims on exhaustive injection campaigns) resurface when attempting to formally prove a countermeasure’s security.

The issue of proving security for countermeasures at the ISA level or lower has received little attention compared to traditional testing. Works that reach proven levels of security usually focus on specific boards (e.g. RTL model checking in [32]) or high-level languages (e.g. symbolic analysis in [28]). Proving traditional countermeasures would provide an appropriate response to [37] by guaranteeing the absence of attack opportunities, across multiple abstraction layers starting at the micro-architectural level.

This approach contrasts with recent progress in working out sophisticated hardening schemes that combine multiple types of protections or resist multiple consecutive or independent attacks (*multi-fault attacks*) [14, 34]. The applicability of these solutions still hinges on ISA-level fault models adequately capturing faults’ effects, which is not guaranteed. To ensure security at a lower abstraction level, we choose to refocus on single but accurate models, which we address with a software/hardware co-design approach.

Contribution summary. We present the first formally-proven countermeasure addressing a low-level fault model. We formalize the fault in an extended assembler language in Section 2, then propose in Section 3 a countermeasure based on a hardware extension and a compile/link-time program transformation. From the formal semantics for the fault, we derive a security theorem in Section 4. We implement the scheme in LLVM and GNU ld, then evaluate its security with emulated fault campaigns and its performance with processor simulations in Section 5. Finally, we compare this work with recent literature in Section 6, and conclude in Section 7. The long version of this article contains a detailed operational semantics of the extended assembler and the complete proof of the security theorem.

2 Formal Description of Fetch Skips

We protect against a slightly generalized version of the fault model described by [3], adapted for RISC-V¹. The target system is a 32-bit little-endian RISC-V processor with the “C”

¹RISC-V is an open standard ISA; details can be found at <https://riscv.org>.

```
int g(int x) { return f(x) + 1; }
g:
# ▼ Push register ra to stack
24: 41 11      c.addi  sp, sp, -16
26: 06 c6      c.sw   ra, 12(sp)
# ▼ Call f (linker later inserts address of f)
# ▼ ra is both target address and return address
28: 97 00 00 00  auipc  ra, 0
2c: e7 80 00 00  jalr   ra, 0(ra)

# ▼ Add 1 to return value a0 of f
30: 05 05      c.addi  a0, a0, 1
# ▼ Pop ra from stack and return a0
32: b2 40      c.lw   ra, 12(sp)
34: 41 01      c.addi  sp, sp, 16
36: 82 80      c.ret

24: | c.addi | c.sw |
28: | auipc (1/2) | auipc (2/2) |
2c: | jalr (1/2) | jalr (2/2) |

30: | c.addi | c.lw |
34: | c.addi | c.ret |
```

Figure 1. C code, object code, and memory layout of a simple function $g(x) = f(x) + 1$.

extension for compressed instructions (such as RV32IC [1]), which means it has a mix of 16-bit and 32-bit instructions.

The instruction set includes the usual arithmetic instructions (add, xor...), memory instructions (lw, sw...), conditional (beq...) and unconditional (j, jal...) branches. All instructions have a 32-bit encoding unless prefixed with “c.”, which denotes a 16-bit encoding from the “C” extension.

For brevity, we use “aligned” to mean a multiple of 4 (bytes) and “unaligned” for a multiple of 4 plus 2. We also use the term “line” to refer to 4 bytes of data at an aligned address.

2.1 RISC-V Programs and Their Execution

This section introduces the notions needed to describe the fault and state the security theorem in Section 4.

We model RISC-V programs as collections of *blocks* made of non- or conditionally-branching instructions, terminated by an unconditional branch.² Figure 1 shows a function $g(x) = f(x) + 1$ consisting of two blocks, one that calls f (implicitly forwarding the argument) and one that increments the return value $a0$ then returns.

We model the execution as a sequence of *execution steps* in which the CPU obtains the next instruction (with a combination of a buffer read and/or a memory fetch), decodes it,

²These aren’t classical basic blocks; not counting conditional branches as terminators works better for this fault. This is of minor importance except in the security proof.

and executes it. For instance, function g in Figure 1 executes in 8 steps (ignoring the call to f), with each step running one instruction and consuming either 2 or 4 bytes of code.

The central behavior of interest in this paper is that *there is not always one fetch for one step*: certain instructions are fetched in advance, and others are read by the CPU over two consecutive fetches. This is because most CPUs *always fetch 4 aligned bytes in memory*, that is, a *line* in the memory layout table. For example, the first step in running g fetches the 4 bytes at address 24 and puts them in a decoding buffer, but only consumes the first 2 to execute $c.addi$. Then, the second step executes $c.sw$ from the buffer without needing another memory fetch.

Figure 2 shows all the ways a step may use buffered data and/or fetch from memory when executing one instruction. When PC is aligned (left), the next instruction is on a new line, so a fetch is performed. When PC is unaligned (right), 2 bytes of instruction data are left in the decoding buffer. If they form a 16-bit instruction (case c), the CPU runs it immediately without a fetch. If they form the first half of a 32-bit instruction (case d), a fetch is performed at $PC + 2$ to obtain the second half and piece together the opcode. Each step then increments PC by 2 or 4, setting up the next cycle.

2.2 Fault Model

We are interested in protecting programs against “*fetch skip*” attacks described by [3]. Fetch skips refine the traditional *instruction skip* fault model by analyzing the effect of physical attacks (here, clock glitches) on the micro-architecture rather than on the program only. The proposed fault model captures 80% of the observed behaviors under clock glitches attack ([2]). Because the model is more accurate with respect to physical effects, a countermeasure against fetch skips will be more secure in practice than a countermeasure against plain instruction skip.

The different types of fetch skips are illustrated in Figure 3. In this model, attempts by the CPU to fetch at an address L may result in:

- **Skip 32 bits**, k times ($S_{32(k)}$): Instead of returning the memory contents at address L , the memory contents at address $L + 4k$ are returned, and PC is incremented by $4k$.³
- **Skip and repeat 32 bits** (S&R32): The contents at address L are requested from memory, but due to the short cycle the decoding stage triggers before the response is available. A copy of the previous line (usually at $L - 4$) is returned to the decoder for the current step, and the requested value arrives in the decoding buffer later.

³[3] observes that a the CPU may run a nop during a $S_{32(k)}$ attack; we omit this detail here but account for it in the countermeasure.

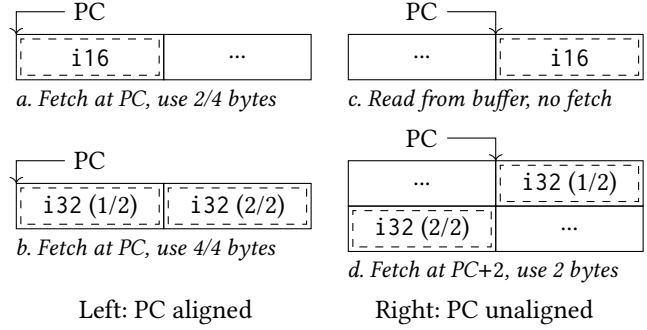


Figure 2. Execution step patterns. i_{16}/i_{32} represent arbitrary 16/32-bit instructions. The dashed region is the next instruction to be run.

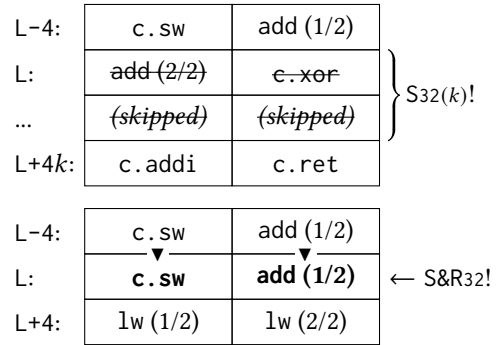


Figure 3. Effect of fetch skip attacks on fetch requests.

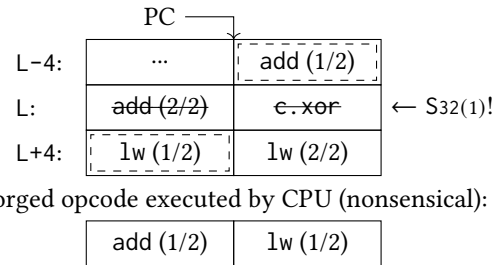


Figure 4. Forging a 32-bit instruction with a $S_{32(1)}$ attack.

We consider a “level N attacker” that can independently attack every fetch with either one $S_{32(k)}$ attack ($1 \leq k \leq N$) or one S&R32 attack.⁴

The connection between fetch skips and the memory layout of instructions creates a new effect, which we call **instruction forging** and illustrate in Figure 4, that cannot be captured by the usual “instruction skip” model. Here, $PC = L - 2$ and the decoding buffer contains the first half of a 32-bit unaligned add instruction. The CPU fetches at L to obtain the second half (Figure 2, case d). Attacking this fetch with $S_{32(1)}$ results in an unrelated 16-bit value (the start of an lw) being used to complete the add , causing the execution of an opcode not originally in the program. This opcode can be “anything” (including illegal). [3] demonstrate

⁴We use these patterns because the behavior of more complex combinations cannot be inferred from [3] without more hardware details.

Name	Rule statement
NoFAULT	$(PC, \rho) a \Rightarrow [a] (PC, [a])$
S32(k)	$(PC, \rho) a \Rightarrow [a + 4k] (PC + 4k, [a + 4k])$
S&R32	$(PC, \rho) a \Rightarrow \rho (PC, [a])$

Figure 5. Short description of fetch rules.

$(PC, \rho) a \Rightarrow d (PC', \rho')$ indicates that the fetch of address a with current PC and ρ returns value d to the decoder, along with updated state PC', ρ' . $[a]$ denotes the memory contents at address a .

how this enables new vulnerabilities (e.g. by forging control flow instructions).

Once PC is out of sync, forging can continue without repeated fault injection. Continuing with Figure 4, after running the forged instruction we get $PC = L + 6$. Now the second half of `lw` is interpreted as its own 16-bit or (first half of) 32-bit instruction. Thus, the attack carries over to the next line; in the worst case, execution might not resynchronize with the intended sequence of instructions until a jump.

2.3 Program Model

To design and prove a countermeasure against such a low-level attack, it is helpful to reflect the faults' effects into language features, and study the updated language with semantic tools. In this section, we build a language model that captures the memory layout of instructions while ignoring most concrete effects (apart from jumps and select registers).

Syntax. A *program* is a collection of *blocks*. A block `bb` is a sequence of *instructions* $bb = [i_1, \dots, i_{|bb|}]$ where each instruction i_j is a 16-bit or 32-bit value, representing an assembler instruction encoded as described by the RISC-V ISA (or, as will be described in Section 3.2, a 32-bit “checksum literal”). We write $\|i\|$ for the size in bytes of instruction i (which is either 2 or 4).

Program state. Due to the stateful nature of S&R32, which duplicates the “previous” fetched line, a full description of the program state must include the decoding buffer. We model this state with a 4-tuple $\langle PC, \rho, \sigma, \alpha \rangle$, containing:

- The current PC value (always a multiple of 2);
- ρ , the last 4-byte value fetched by the CPU (decoding buffer). ρ is used to decode at unaligned PCs (Figure 2, right) and replaces the targeted line during a S&R32 attack;
- σ , the values of the Xccs registers associated with our countermeasure, whose role is described in Section 3.2;
- α , the rest of the architectural state (other registers, memory...), which we abstract away as an opaque value.

The attacks shown in Figure 3 are formalized in terms of this state in Figure 5. To illustrate alignment dynamics, attacking the fetch at $0x30$ during the execution of `g` (just upon returning from `f`) produces the following execution:

```

g: # PC >= 0x40000 is a protected region
400e4: 41 11      c.addi  sp, sp, -1
400e6: 06 c6      c.sw   ra, 12(sp)
# ▼ auipc was relocated
400e8: 97 00 00 00  auipc  ra, 0
# ▼ Checksum was set and LSB flipped
400ec: 0b 64 00 00  ccscallb 8
400f0: e2 75 06 c6  .word   0xc60675e2
# ▼ jalr was relocated
400f4: e7 80 00 fb  jalr   ra, -80(ra)
400f8: 02 90      c.break # 8 times

40108: 05 05      c.addi  a0, a0, 1
4010a: b2 40      c.lw   ra, 12(sp)
4010c: 41 01      c.addi  sp, sp, 16
4010e: 01 00      c.nop
# ▼ Checksum was set, no flip needed
40110: 0b 10 00 00  ccs
40114: 51 16 b3 40  .word   0x40b31651
40118: 82 80      c.ret
4011a: 02 90      c.break # 8 times

```

Figure 6. Linked object code for `g` after hardening ($N = 2$). Numbers on the right refer to lines of Algorithm 1 that added the instructions.

```

⟨0x30, ρ0, σ0, α0⟩
├ S32(1): (0x30, ρ0) 0x30 ⇒ [0x34] (0x34, [0x34])
├ Run instruction: c.addi sp, sp, 16
├ ⟨0x36, ρ1 = [0x34], σ1, α1⟩
├ No fetch. Run instruction: c.ret
├ ⟨α1.ra = 0x30, ρ2 = ρ1, σ2, α2⟩

```

`c.ret` jumps to the address stored in register `ra`, which is still $0x30$ because we skipped the load from the stack, later leading to a typical stack corruption crash.

3 A Co-designed Countermeasure

In this section, we propose a software/hardware countermeasure, based on code instrumentation with hardware support.

3.1 Overview

The key ideas of the countermeasure are as follows. Machine code is augmented (“hardened”) during compilation with checksum protections that react to a fault attack by forcing execution to trap before the end of the current block. This limits exploits by ensuring a sufficiently tight window between attack and detection (like most countermeasures it doesn’t prevent side-effects immediately resulting from the fault, which is inherently difficult due to timing). Figure 6 shows the hardened code for the `g` function from Figure 1.

Hardware is modified to automatically maintain a running checksum (in fact, a simple *sum*) of every line of instruction data fetched from memory during the execution of a block, independent of instruction alignment. Blocks are compiled so

that every exit is guarded by a `ccs` instruction (from our ISA extension), which traps if the running checksum is not equal to a reference value computed at compile-time. Blocks thus act as autonomous “jails”, in that faulty executions causing the checksum to deviate from its expected value cannot leave their current block.

The co-designed nature of the countermeasure is a result of the cross-checking of information between hardware, whose monitoring produces a trace (checksum) sensitive to fault attacks, and software, which provides reference checksum values to interpret that trace.

Section 3.2 describes the hardware extension that we rely on for the countermeasure. 3.3 details the hardening algorithm. 3.4 highlights the subtleties of implementing the hardening algorithm in LLVM. Finally, 3.5 discusses the design choices in a more general context.

As always, one can attempt to attack the countermeasure itself. The security theorem in Section 4 proves that no attempt at skipping, repeating or forging `ccs` instructions or jumps can succeed, leading to a strong security guarantee.

3.2 ISA and Hardware Extensions

We support the countermeasure with a custom ISA extension named `Xccs` for *Code CheckSum*⁵. Throughout this section, all examples refer to Figure 6, and u_N denotes the type of N -bit unsigned integers.

`Xccs` introduces four new Control and Status Registers (CSR), which together form the σ field of the program state:

- `CCS` : u_{32} is the running checksum for the current block. For example, the region of `g` from address `40108` through `40114` adds up in little-endian to $0x40b20505 + 0x00010141 + 0x0000100b = 0x40b31651$ so the value of `CCS` at `40114` will be `0x40b31651` if no attack is performed.
- `CCSPROT` : u_{32} indicates the PC value at which a protected instruction (jump) is expected to execute. It is zero most of the time and non-zero for a single step after a checksum is validated. In Figure 6, it is set at `400f4` and `40118`.
- `CCSD` : $\{E : u_1, JO : u_5\}$ holds control information; the *enable* bit (`E`) indicates whether checksum protection is active for the current block (which currently is whenever $PC \geq 0x40000$), and the *jump offset* field (`JO`) is set by `ccscall` as described below.
- `CCSDS` : u_{32} (“delay slot”) is used to implement the decoding of `Xccs` instructions’ 32-bit checksum arguments.

The meat of `Xccs` is the addition of four guard instructions whose encoding is shown in Figure 7, all of which are followed by a 32-bit `<checksum>` argument:

- `ccs` `<checksum>` compares the `CCS` register with the provided argument and traps if they differ. Otherwise, it sets $CCSPROT = PC + 8$ so a jump or function call can execute at the next step. For instance, at `40114` in Figure 6,

31	25	24	20	19	15	14	12	11	7	6	0	
funct7	rs2	rs1	funct3	rd	opcode							
0	0	0	001	0	XCCS							<code>ccs</code>
0	0	0	010	N	XCCS							<code>ccscall N</code>
0	0	0	101	0	XCCS							<code>ccsb</code>
0	0	0	110	N	XCCS							<code>ccscallb N</code>

Figure 7. `Xccs` instructions (32-bit). `XCCS` is `0001011`.

the dynamic value of the `CCS` register is compared to `0x40b31651`. When no faults are injected, these are equal, so execution proceeds to the `c.ret` instruction.

- `ccscall N` `<checksum>` is similar; it is used before function calls. It sets $CCSD.JO = N$, which causes the next function call’s return address to increase by $2N$ bytes. For instance, the `ccscallb 8` at `400ec` changes the return address after the `jalr` (call) instruction from `400f8` to `40108`, skipping over the `c.ebreak` block.
- `ccsb` and `ccscallb` are variations that flip the least significant bit (LSB) of the checksum before comparing it. This is leveraged to ensure that no `<checksum>` argument decodes as a jump or `Xccs` instruction, as detailed in Section 3.3. Again with Figure 6, the sum of the first block up to `400f0` is `0xc60675e3`, but that decodes as a `bltu`, creating a vulnerability. The LSB is flipped to avoid this.

All `Xccs` instructions further trap when run at unaligned PC, which prevents most attempts at forging them. Finally, existing CPU behavior is modified as follows:

1. All branch instructions trap if $PC \neq CCSPROT$. If a branch occurs, `CCS` and `CCSPROT` are reset to 0.
2. Every other instruction traps if $CCSPROT \neq 0$.
3. Every value retrieved from a fetch is added to `CCS`.
4. Call instructions add $2 \times CCSD.JO$ to the return address and clear `CCSD.JO` before jumping.

Hardware support in this countermeasure serves an important dual purpose: it enforces the jail, and it updates the checksum through a system that is not vulnerable to attacks, so that checks of the trusted `CCS` value can be guaranteed even though the reference `<checksum>` is itself vulnerable.

3.3 Hardening Algorithm

Algorithm 1 shows the hardening process for a single block in pseudocode, with our canonical example in Figure 6. (As all blocks are independent, this algorithm is executed by the compiler for every block in the program.)

The main `for` loop (line 15) iterates over the instructions of the original block. All instructions are copied to the hardened block (line 26). Jump instructions are preceded by a guard, which is `ccscall/ccscallb` (line 22) for function calls and `ccs/ccsb` for other jumps (line 24). `nops` are used to ensure that guards are aligned and jumps are always separated by at least one non-jump instruction, both of which prevent subtle attacks against the countermeasure.

⁵The “X” is standard notation for unofficial RISC-V extensions.

Calling convention: $a0$ serves as both argument and return value for f and g . ra is used for function calls and is caller-saved (push/pop omitted).

<pre> 1 g: (...) # push ra to stack 2 PseudoCALL @f 3 # ^a0 forwarded implicitly 4 5 6 \$a0 = nsw ADDI \$a0, 1 7 (...) # pop ra from stack 8 9 PseudoRET \$a0 10 </pre>	<pre> 1 g: (...) # push ra 2 PseudoCALL @f, .LBB1_0 3 PseudoCCSTRAP 2 4 5 .LCCS_Region_Start0: 6 \$a0 = nsw ADDI \$a0, 1 7 (...) # pop ra 8 CCS .LCCS_Region_Start0 9 PseudoRET \$a0 10 PseudoCCSTRAP 2 </pre>	<pre> g: e4: 41 11 06 c6 (...) # push ra e8: 97 00 00 00 auipc ra, 0 ec: 0b 24 00 00 ccscall 8 # ▼ R_RISCV_CHECKSUM: g f0: 00 00 00 00 .word 0x00000000 f4: e7 80 00 00 jalr ra, 0(ra) f8: 02 90 c.ebreak # 8 times (2N+4) .LCCS_Region_Start0: 108: 05 05 c.addi a0, a0, 1 10a: b2 40 41 01 (...) # pop ra 10e: 01 00 c.nop # ccs alignment 110: 0b 10 00 00 ccs # ▼ R_RISCV_CHECKSUM: .LCCS_Region_Start0 114: 00 00 00 00 .word 0x00000000 118: 82 80 c.ret 11a: 02 90 c.ebreak # 8 times </pre>
(a) Machine IR before back-end pass.	(b) Machine IR after back-end pass.	(c) Object code before linking.

Figure 8. Stages of hardening the program from Figures 1 and 6 ($N = 2$) in our LLVM implementation.

The second **for** loop (line 27) adds a barrier of `c.ebreak` instructions, which raise a distinctive exception when executed. Their role is to prevent control from leaving the block by skipping over the terminator. Up to $2N + 4$ are needed to address the worst case where control reaches the middle of a checksum whose second half is a 32-bit opcode, in which an attack of a $S\&R_{32}$ followed by a $S_{32(N)}$ would reach $4N + 6$ bytes past the checksum.

Procedure `addToBlock` appends instructions to the hardened block while computing the reference checksum value sum by summing instruction’s opcodes. This accounts for instructions’ alignment with the function

$$\text{realign}(\text{offset}, i) = \begin{cases} i & \text{if } \text{offset} \equiv 0 \pmod{4} \\ 2^{16}\text{LSH}(i) + \text{MSH}(i) & \text{otherwise} \end{cases}$$

where $\text{LSH}, \text{MSH} : u_{32} \rightarrow u_{16}$ are the Most and Least Significant Halves respectively. This process is equivalent to summing the lines of the final layout table. Finally, `addChecksum` selects whether to use `ccs/ccscall` or their `-b` variants. The `-b` variants are selected when the checksum value is an “invalid checksum literal”, i.e. it decodes as a jump, `Xccs` instruction or `c.ebreak`. Such values could be misused as instructions if an attacker were to skip the guard that precedes it. Flipping the LSB ensures that these sensitive values do not appear in code. Our running example in Figure 6 is obtained by executing this algorithm on both blocks of g ’s original code from Figure 1 and linking it.

3.4 LLVM Implementation

Algorithm 1 cannot be implemented as-is in a single pass in a standard compiler, because reference checksum values depend on the exact bit-level encoding of each instruction, which is not decided until the linker relocates references to globals and functions. See for instance how the call in

Figure 8c (before linking) has placeholder zero-offsets but the one in Figure 6 (after linking) has a proper target offset.

We see this as a benign occurrence of a fundamental issue: the progressive lowering that standard toolchains go through is designed around *functional* invariants, not *security* invariants. Encodings are decided late because they don’t matter to the compiler, which is only interested in the functional specification of instructions. The addition of a security countermeasure to the compiler breaks this assumption and creates a new concern in *threading the security transform through the abstraction gap of a functional lowering*.

In this case, we are able to implement the algorithm in two steps: a late Machine IR⁶ pass followed by an extension to the linker relocation process.

- **Machine IR pass:** the program’s Machine IR representation is first transformed late in the back-end (from Figure 8a to Figure 8b). This pass handles all tasks that *add* code into the program, including:
 - Aligning functions and blocks to 4-byte boundaries;
 - Adding aligned `Xccs` instructions before all Machine IR instructions that expand into RISC-V jumps, such as `PseudoRET`. A label indicates the start of the region that the checksum must cover;
 - Adding the trap barrier with the `PseudoCCSTRAP` N pseudo-instruction, which later expands into a series of $2N + 4$ `c.ebreak` instructions.

At this stage some jumps are still hidden in pseudo-instructions, like the function call in `PseudoCALL`. This is because far jumps in RISC-V are implemented with a pair of instructions, `auipc` and a jump, to overcome the limited distance that can be encoded into single jump instructions.

⁶The back-end intermediate representation used by LLVM, which is aware of the target architecture, not the usual LLVM IR.

Algorithm 1 Algorithm: HARDEN

Input: A source block $[i_1, \dots, i_n]$
Input: Upper bound N for $S_{32}(k)$ rule ($k \leq N$)
Output: A hardened block hb .

```

1:  $hb \leftarrow []$ 
2:  $sum : u_{32} \leftarrow 0$ 
3:  $offset \leftarrow 0$  ▷ Blocks are 4-aligned.
4: procedure addToBlock( $i$ )
5:    $hb.append(i)$ 
6:    $sum \leftarrow sum + \text{realign}(offset, u_{32}(i))$ 
7:    $offset \leftarrow offset + \|i\|$ 
8: procedure addChecksum( $i, ib$ )
9:   if  $sum + i$  is not a valid checksum literal then
10:     addToBlock( $ib$ )
11:     addToBlock( $sum \oplus 1$ )
12:   else
13:     addToBlock( $i$ )
14:     addToBlock( $sum$ )
15: for  $i$  in  $[i_1, \dots, i_n]$  do
16:   if  $i$  is a jump instruction then
17:     if  $offset = 0$  then ▷ No empty sections.
18:       addToBlock( $\text{encode}(nop)$ )
19:     else if  $offset \equiv 2 \pmod{4}$  then ▷ Force alignment.
20:       addToBlock( $\text{encode}(c.nop)$ )
21:     if  $i$  is a function call then
22:       addChecksum( $\text{encode}(ccscall(2N + 4)),$   

23:          $\text{encode}(ccscallb(2N + 4))$ )
24:     else ▷ Jumps/branches.
25:       addChecksum( $\text{encode}(ccs), \text{encode}(ccsb)$ )
26:      $offset \leftarrow 0$ 
27:   addToBlock( $i$ )
28: for  $j = 1$  to  $2N + 4$  do ▷ Add trap barrier.
29:   addToBlock( $\text{encode}(c.ebreak)$ )
30: return  $hb$ 

```

In LLVM, this is expanded later in the code emitter due to limitations in the back-end structure; we count this towards the Machine IR pass for simplicity of exposition. The Machine IR pass is followed by static branch relaxation (the unfolding of far jumps into multi-instruction sequences and converse compaction of near jumps into shorter instructions), which is the main reason why the addition of new code cannot be delayed more; it would break short jumps.⁷

During object file generation, the 8-byte CCS Machine IR instruction is replaced with a 4-byte Xccs opcode and a placeholder zero-checksum. A custom relocation of type `R_RISCV_CHECKSUM` (marked by a comment in Figure 8c) is added to mark the checksum region for the linker.⁸

⁷We disable linker relaxation for simplicity to maintain jumps' alignment, but it could be enabled after adding appropriate alignment relocations.

⁸The `R_RISCV_CALL` relocation for `auipc/jump` pairs is also replaced with a custom type to inform the linker of the newly-added `ccs` in the pair.

- **Linking:** the linker follows relocation entries to compute checksums and insert them in the provided spaces. The linker script is also updated so that hardened objects are linked to a different virtual address (`0x40000`) than non-hardened libraries and runtime files (`0x10000`), which we use as the basis to enable Xccs protection.

3.5 Discussion

We discuss some of our design choices and how the scheme interacts with other micro-architectural features.

Fault's effect on CCS updates. Our countermeasure relies on CCS updates not being vulnerable to attacks. This is a reasonable inference based on the fault model: recall that fetch skips are induced by clock glitches, which are known to cause problems locally along critical signal propagation paths. [3] identify that such paths are mostly in the fetch stage of the pipeline, but CCS can be updated in the execution stage using the `realign sum` technique, and is thus unlikely to be affected.

Interrupts. Common interrupts and signal handlers (that are invisible from the main thread) would not interfere with Xccs protections (with the only OS support needed being to save Xccs registers, which can be viewed as an extension of PC, to the CPU context structure). However, a non-returning interrupt (such as a signal exiting) would leave the current block without a check. We assume such a no-returning action implies abandoning the critical section where the interrupt occurred; otherwise, there might be a vulnerability.

Effect of faults on complex architectures. The fault model from [3] does not describe hardware responses to fault attacks during speculative or out-of-order execution. The study and design of fault models at the micro-architectural level is already state-of-the-art, and applying it to these complex features is a completely open problem. While Xccs is amenable to speculative execution (mispredictions would not lead to false checksum exceptions because the checksum resets at the beginning of every block) and out-of-order execution (the checksum update is associative-commutative, allowing for reordering within each block) it remains unlikely that clock glitches would affect such complex designs in the same way as the simple processors from which fetch skips are derived.

Possibility of a hardware-only solution. Hardware-only countermeasures against fault attacks present their own challenges [12]. Detecting the clock glitch at the (hardware) source using a *fault detector* [15] creates a performance trade-off between sensitivity and the rate of false alarms, which limits the approach to security-critical systems that can afford the performance loss. Detecting corruption in the instruction stream requires extra hardware logic that risks being itself faulted (increased attack surface). Software/hardware propositions like our Xccs countermeasure have contrasting benefits. We minimize exposure to the fault because the detection only relies on CCS updates (which occur in

the execution stage, away from disrupted fetch logic) and a checksum check made after the fault’s transient effect has subsided. (We carefully discussed the safety of these operations with authors of [3].) We also get the benefits of using more software and less hardware, such as incurring costs only in critical sections and the possibility of implementing hardware support late or during minor hardware revisions.

Specificity of the attack model. Many (mostly early) works in fault literature attempt to protect against *all* program misbehaviors, described as “soft errors”. By contrast, we target a single vulnerability, which might appear overly specific. However, the fetch skips model coined by [3] results from extensive physical injection campaigns, where it described the impact of 80–90% of clock glitches on Cortex-M boards [2], making this countermeasure useful against common attack vectors on real boards. In addition, we argue that the lack of a precise definition for soft errors leads to tricky vulnerabilities⁹ preventing any proof-based security standard from being met. This is why we focus on fetch skip attacks, for which we can formally prove security.

4 Security Theorem

We now state the security theorem for our countermeasure and the main steps of the proof. We omit some invariants about σ/ρ whose significance is only clear with the full formalization and proof (in the research paper version of this article).

We define a *step* as a statement

$$s = \langle PC, \rho, \sigma, \alpha \rangle \rightarrow \langle PC', \rho', \sigma', \alpha' \rangle$$

materializing the change in execution state associated with the execution of a single instruction, and an *execution* as a sequence of steps $[s_0, \dots, s_n]$ such that the initial state of any s_i ($i > 0$) is the final state of s_{i-1} .

Theorem (simplified).

Let $P = [\text{HARDEN}(\text{bb}_1), \dots, \text{HARDEN}(\text{bb}_{|P|})]$ a program hardened by Algorithm 1, and $e = [s_0, \dots, s_{|e|}]$ an execution such that

- s_0 starts at the top of a block of P ;
- $s_{|e|}$ successfully executes a program exit instruction.

Then e can be partitioned into sub-sequences $(s_{t_i} \dots s_{b_i})_{1 \leq i \leq m}$ each executing one hardened block $\text{HARDEN}(\text{bb}_i)$, such that

1. Each s_{t_i} starts at the top of $\text{HARDEN}(\text{bb}_i)$;
2. Each s_{b_i} ($i \neq m$) runs a legitimate (non-forged) jump instruction of bb_i , and $\sigma_{b_i}.\text{CCS}$ is the correct associated checksum;

⁹For instance, triplication countermeasures such as SWIFT-R [10] and NEMESIS [13] tend to assume that a single “soft error” only affects one of the three execution streams, but this is not true of e.g. the kind of decoding errors mentioned in [13].

3. If each segment $s_{t_i} \dots s_{b_i}$ contains at most one faulted fetch and the last segment $s_{t_m} \dots s_{b_m}$ has none¹⁰, then the entire execution e contains no faults.

In essence, this theorem says that the countermeasure guarantees the execution of the original sequence of instructions (as witnessed by the checksum) at every control flow edge. The key property is the *local security* at the block level; the security property of passing checksums is checked at every block, and deviations from that property are detected before the block ends. This facilitates both formal analysis (by removing difficulties associated with control flow) and testing (by providing an easy way to detect failures in the countermeasure).

The progression of the argument is as follows:

1. Control cannot reach the end of a protected block due to the `c.ebreak`-based trap barrier;
2. Control can only leave a protected block through a jump if it passes the associated checksum;
3. Using exactly one faulted fetch in the execution of a block always invalidates the checksum.

The checksum method by itself does not always guarantee that the execution cannot be compromised, as a vulnerable block could contain a *checksum collision*: an unintended execution path, reachable with multiple faults, whose checksum is the same as the intended execution. Such collisions can be detected statically after linking, but they are hard to solve because the linker is unable to add or remove instructions to prevent the collision once it is discovered, as the first could break short jumps and both would change jump offsets thus more checksums. In this case, the toolchain’s abstraction stack that we worked around fairly easily in Section 3.4 poses a greater challenge.

Note that the countermeasure still guarantees that checksums must pass even in multi-attack scenarios. Collisions are rare (as showcased in Section 5) despite the apparent weakness of the summing algorithm, mainly because the attacker cannot feed arbitrary inputs to the sum; they can only skip or repeat predetermined inputs. Thus, the scheme retains value as a practical protection mechanism against multi-fault attacks.

5 Implementation Setting and Evaluation

We now describe the implementation and the experimental settings we used to evaluate the countermeasure. We examine functional correctness, security guarantees, and cost.

¹⁰We do not protect the exit of the last block because in practice the last instruction is a `syscall` invocation in the non-hardened `libc` function `exit()`.

5.1 Implementation and Experimental Setting

We implemented the hardening scheme in LLVM [20] 12.0,¹¹ linking programs with an off-the-shelf GNU C library using a modified GNU ld 2.40 linker. The hardware extension Xccs was prototyped in the system emulator QEMU 8.0 [8] and the processor simulator Gem5 22.1 [24]. The implementation is about 1150 lines in LLVM/ld for the countermeasure, and 850 lines in QEMU/Gem5 for evaluation.¹²

We evaluate our countermeasure on programs from the MiBench benchmark suite [16], which is designed to be representative of embedded applications. We do not specifically use a cryptography- or security-oriented benchmark as we protect code at a low enough level that the nature of the program is not significant. Programs were hardened with $N = 2$, i.e. assuming a strong attacker that can inject double-skips every fetch. We perform three types of experiments, on a standard x86_64 GNU/Linux machine:

- Exhaustive single-fault injection campaigns in QEMU;
- Random multi-fault injection campaigns in QEMU;
- No-fault performance simulations in Gem5.

We instrument QEMU to support fault injections by skipping or replacing data during the transpilation step. We take advantage of the locally-correct design of the countermeasure and raise an explicit “countermeasure bypassed” exception at every successful exit of a block where a fault was triggered. This ensures that successful fault attacks are reported and cannot be masked by program logic. We similarly extend Gem5 with the decoding and execution of Xccs instructions, but leave out exceptional conditions since we use it only to evaluate performance without fault injections.

5.2 Functional Correctness

We first assert that hardening preserves the functional behavior of programs when no fault is injected. This is checked by running hardened MiBench programs and comparing their output to a non-cross-compiled x86 build and a non-hardened RISC-V build.¹³ All programs pass this test.

5.3 Security Guarantee

We then subject each program to two injection settings, both performed by emulation in QEMU:

1. *Single-fault exhaustive injection*: we attack every PC in the protected text segment of the program with single-fault `S32(1)`, `S32(2)`, and `S&R32`.
2. *Multi-fault random injection*: we attack 2000 code intervals of 64 bytes with randomly-selected sequences of 2 to 6 faults in close succession.

¹¹All the modified tools used in this paper are publicly available at gricad-gitlab.univ-grenoble-alpes.fr/michelse/fetch-skips-hardening.

¹²Gem5 currently only emulates 64-bit RISC-V instructions; we modified it to support 32-bit instructions for this evaluation.

¹³For some programs, comparison with the native build is skipped due to differences in the precision of the math library.

This showcases the guarantees proven in Section 4. Outcomes of injections for which the targeted PC was reached are classified in four categories in Figure 9:

- **Fault reported**: An invalid checksum, illegal jump, or other Xccs-imposed constraint was violated; or `c.ebreak` was executed.
- **Segfault**: The program segfaulted from an incorrect memory access as a result of the injected fault.
- **Other crash**: Illegal instructions being decoded by the CPU (SIGILL) or other rare crashes such as SIGBUS.
- **Countermeasure bypassed**: The program successfully exited a block after attacking it. This did not happen in any of the tests.

Figure 9 shows that no attack was able to bypass the protection from the countermeasure, which in the multi-fault case means that no checksum collision occurred.

A significant majority of faults that are reached result in Xccs violations, meaning that failing checksums and invalid attempts to jump out of blocks are adequately reported. Because all exits of attacked blocks were guarded with “Countermeasure bypassed” exceptions, which were not reached, every crash we encountered also occurred within the attacked block. This shows that the countermeasure fulfills its goal of containing faults within blocks.

Interestingly, crashes alone do not provide sufficient security: unprotected programs (not depicted) experience ~90% crashes but only ~30% in the block where the fault is injected. So like with instruction skip, a system facing random failures could leave the code unprotected, but the threat of targeted attacks requires a countermeasure like Xccs.

5.4 Performance

On average, hardening with $N = 2$ increased the size of protected functions (i.e., excluding libraries and runtime files) by a factor of 2.46 (which is large but usual for security applications especially against skips, due to duplication). Individual differences are given in Figure 10; unsurprisingly, programs with longer straight-line sections see less of an increase, while short, loop-intensive programs like `dijkstra` and `bitcount` get the largest overhead.

We run Gem5 simulations for each program¹⁴ to estimate the overhead in execution time, also reported in Figure 10. The simulated system mainly consists of a 1-GHz RISC-V core with a 1600 MHz DDR3 controller simulated by Gem5’s *timing* model.

The main cost is the execution of Xccs instructions themselves, which incur extra fetch-decode-execute cycles. As expected, programs with denser control flow in hot sections such as `dijkstra` and `bitcount` see by far the largest difference. We expected the code size increase to reduce cache

¹⁴Except `patricia`, due to a 32- vs. 64-bit compatibility issue related to reused opcodes in the RISC-V ISA.

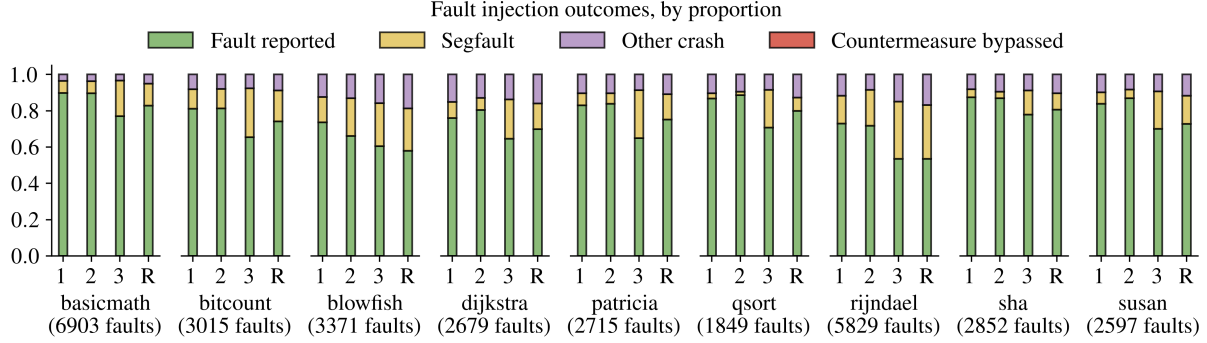


Figure 9. Outcomes of exhaustive S32(1), S32(2) (bars 1 and 2) and S&R32 (bars 3) single-fault injections, and random multi-fault injections on MiBench programs (bars R). The total number of injected faults is given for each program.



Figure 10. Evaluation of the cost of our counter-measure: static and dynamic metrics.

efficiency, and ran simulations with an 8-kB 4-way instruction cache (the median size of hardened sections in our test, which is a lot smaller than the total text segment). This speeds up executions significantly (not pictured) but the ratio of hardened to non-hardened speed remains consistent. The overall performance hit is 10.2%, consistent with the knowledge that MiBench programs average out 6–7 instructions per basic block (keeping in mind that our blocks are longer than traditional basic blocks). Unsurprisingly, cryptographic schemes have the best performance results due to being computation-heavy with long blocks.

Existing countermeasures without hardware support have a comparatively much larger overhead. [14] compare multiple combinations of countermeasures against single- to quadruple-instruction skips on a secure boot program. Only 3 combinations (nZDC + RACFED, NEMESIS + RACFED, and CompaSec) close more than 75% of vulnerabilities, all with space overhead above $\times 4.85$ and time overhead above $\times 5.01$.

6 Related Work

The contribution of this paper follows the now numerous developments in the area of compilation for security. Recent countermeasure designs have consistently involved compilation tools [5, 29, 35], and issues stemming from cross-layer

abstractions have been raised previously [6]. Our work differs in its deliberate involvement of hardware, whereas existing designs that account for non-ISA details (such as timings in [35]) still end up performing pure assembly transformations.

Countermeasure designs that lay more heavily towards hardware also exist. [25] show how hardware support can improve a classic countermeasure against data corruption which consists in executing instructions multiple times and comparing results. This is traditionally very costly due not only to repeated execution, but also program size, register pressure, and the addition of many branches for comparisons. In this work, the ISA is extended with a new “replay” instruction `rp1 w n` which repeats the next w instructions n times, and compares in hardware with an exception in case of a mismatch. A benefit is that the compiler only needs to mark sensitive instructions, however this shifts more security burden to hardware and increases the effort for hardware qualification/certification compared to a more lightweight extension like Xccs.

Our countermeasure is based on classical redundancy ideas and draws inspiration from previous work on instruction skips [36] and control-flow integrity protections [27, 38]. Complications associated with multi-fault attacks incentivized us to aim for a design that guarantees security *locally*

(at every block) to keep formal reasoning simple. Many countermeasures in security literature (for instance [10, 13, 14]) achieve sophisticated and carefully-balanced compiler transforms, but lack a detailed enough programming and fault model to allow for a formal proof of the security property. We believe that our semantics approach that includes “just enough” details in the programming model, encodes the fault model in the semantics, and encourages formal proofs, could be generalized for other low-level countermeasures.

As for the correctness of securization schemes, apart from manual proofs (in which this work falls), automatic countermeasure verification techniques includes translation validation [9], static analysis [11], and symbolic execution [28]. These usually focus on proving properties of the form “either the attack is detected or the program crashed”, as we did here. Another remarkable piece of work in this area is a proof that the C compiler CompCert [23] preserves the constant-time countermeasure when applied to source code [17]. A possible extension of our work would be to further validate through one of these approaches.

7 Conclusion

We have presented a new countermeasure to a micro-architectural fault attack targeting the instruction fetch unit in RISC-V processors. Our approach combines a hardware extension of the RISC-V ISA, dubbed Xccs, and a machine code transformation performed at compile- and link-time. The design process was guided by a specially-tailored semantics that captures just the right set of architectural details to formalize the hardening process and enable a proof of security. We evaluated the scheme’s security by emulating both exhaustive single-fault and random multi-fault injection campaigns, and its performance with a processor simulation. To our knowledge, the present contribution is the first example of such a proven countermeasure for a micro-architectural fault model; which notably connects recent security experts’ results to compiler back-end techniques.

We believe that this approach is applicable more generally to the design of code hardening schemes that guarantee both security and safety. We implemented this particular countermeasure in the late back-end, which is pretty standard for low-level attacks, and keeps the abstraction gap between the attack and countermeasure fairly small. In future work we wish to explore countermeasures implemented at other (sometime multiple) levels in the compilation flow, widening that gap. This will escalate toolchain integration challenges already raised in this paper, as the compiler’s heavy descent in abstraction might break or hinder security properties. Partial solutions such as [33], which advocates for embedding security properties into functional properties, leads us to believe that cleaner, more methodical integration is possible. Similar improvements to the semantic model appear needed to mirror this property preservation in the security proof.

Data-Availability Statement

An artifact of this work is available for reproduction on Zenodo [26] and includes pre-built tools and reference outputs. Source code and documentation for the tools (LLVM, GNU binutils, QEMU, Gem5) is further available for reuse through the Git repository at gricad-gitlab.univ-grenoble-alpes.fr/michelse/fetch-skips-hardening.

References

- [1] 2019. The RISC-V Instruction Set Manual Vol. I.
- [2] Ihab Alshaer. 2023. *Cross-Layer Fault Analysis for Microprocessor Architectures (CLAM)*. Theses. Université Grenoble Alpes [2020-....]. <https://theses.hal.science/tel-04417620>
- [3] Ihab Alshaer, Brice Colombier, Christophe Deleuze, Vincent Beroulle, and Paolo Maistri. 2022. Variable-Length Instruction Set: Feature or Bug? *2022 25th Euromicro Conference on Digital System Design (DSD)*. <https://doi.org/10.1109/DSD57027.2022.00068>
- [4] H. Bar-El, Hamid Choukri, D. Naccache, Michael Tunstall, and C. Whelan. 2006. The Sorcerer’s Apprentice Guide to Fault Attacks. *Proc. IEEE* 94, 2 (2006), 370–382. <https://doi.org/10.1109/JPROC.2005.862424>
- [5] Thierno Barry, Damien Couroussé, and Bruno Robisson. 2016. Compilation of a Countermeasure Against Instruction-Skip Fault Attacks. In *Workshop on Cryptography and Security in Computing Systems (Proceedings of the Third Workshop on Cryptography and Security in Computing Systems)*. vienna, Austria. <https://doi.org/10.1145/2858930.2858931>
- [6] Thierno Barry, Damien Couroussé, Bruno Robisson, and Karine Heydemann. 2017. Automated Combination of Tolerance and Control Flow Integrity Countermeasures against Multiple Fault Attacks. In *European LLVM Developers Meeting*. Saarbrücken, Germany. <https://hal.sorbonne-universite.fr/hal-01660160>
- [7] Gilles Barthe, François Dupressoir, Pierre-Alain Fouque, Benjamin Grégoire, and Jean-Christophe Zapolowicz. 2014. Synthesis of Fault Attacks on Cryptographic Implementations. In *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security (Scottsdale, Arizona, USA) (CCS '14)*. Association for Computing Machinery, New York, NY, USA, 1016–1027. <https://doi.org/10.1145/2660267.2660304>
- [8] Fabrice Bellard. 2005. QEMU, a Fast and Portable Dynamic Translator. In *Proceedings of the Annual Conference on USENIX Annual Technical Conference (Anaheim, CA) (ATEC '05)*. USENIX Association, USA, 41.
- [9] Matteo Busi, Pierpaolo Degano, and Letterio Galletta. 2022. Towards Effective Preservation of Robust Safety Properties. In *Proceedings of the 37th ACM/SIGAPP Symposium on Applied Computing (Virtual Event) (SAC '22)*. Association for Computing Machinery, New York, NY, USA, 1674–1683. <https://doi.org/10.1145/3477314.3507084>
- [10] J. Chang, G.A. Reis, and D.I. August. 2006. Automatic Instruction-Level Software-Only Recovery. *International Conference on Dependable Systems and Networks (DSN'06)*, 83–92. <https://doi.org/10.1109/DSN.2006.15>
- [11] Maria Christofi, Boutheina Chetali, Louis Goubin, and David Vigilant. 2013. Formal Verification of a CRT-RSA Implementation Against Fault Attacks. *Journal of Cryptographic Engineering* 3, 3 (2013), 157–167. <https://doi.org/10.1007/s13389-013-0049-3>
- [12] Ruan de Clercq and Ingrid Verbauwhede. 2017. A survey of Hardware-based Control Flow Integrity (CFI). *CoRR* abs/1706.07257 (2017). arXiv:1706.07257 <http://arxiv.org/abs/1706.07257>
- [13] Moslem Didehban, Aviral Shrivastava, and Sai Ram Dheeraj Lokam. 2017. NEMESIS: a software approach for computing in presence of soft errors. In *Proceedings of the 36th International Conference on Computer-Aided Design (ICCAD '17)*. IEEE Press, Irvine, California, 297–304.
- [14] Johannes Geier, Lukas Auer, Daniel Mueller-Gritschneider, Uzair Sharif, and Ulf Schlichtmann. 2023. CompaSec: A Compiler-Assisted Security

- Countermeasure to Address Instruction Skip Fault Attacks on RISC-V. In *Proceedings of the 28th Asia and South Pacific Design Automation Conference (Tokyo, Japan) (ASPDAC '23)*. Association for Computing Machinery, New York, NY, USA, 676–682. <https://doi.org/10.1145/3566097.3567925>
- [15] Kamil Gomina, Jean-Baptiste Rigaud, Philippe Gendrier, Philippe Candelier, and Assia Tria. 2014. Power supply glitch attacks: Design and evaluation of detection circuits. In *2014 IEEE International Symposium on Hardware-Oriented Security and Trust (HOST)*. 136–141. <https://doi.org/10.1109/HST.2014.6855584>
- [16] M.R. Guthaus, J.S. Ringenberg, D. Ernst, T.M. Austin, T. Mudge, and R.B. Brown. 2001. MiBench: A free, commercially representative embedded benchmark suite. *Proceedings of the Fourth Annual IEEE International Workshop on Workload Characterization. WWC-4 (Cat. No.01EX538)*, 3–14. <https://doi.org/10.1109/WWC.2001.990739>
- [17] Rémi Hutin. 2021. *Verified Secure Compilation against Timing Side-Channels*. Ph. D. Dissertation. <http://www.theses.fr/2021ENSR0029> Thèse de doctorat dirigée par Blazy, Sandrine et Pichardie, David Informatique Rennes, École normale supérieure 2021.
- [18] Andrea Höller, Armin Krieg, Tobias Rauter, Johannes Iber, and Christian Kreiner. 2015. QEMU-Based Fault Injection for a System-Level Analysis of Software Countermeasures Against Fault Attacks. In *2015 Euromicro Conference on Digital System Design*. 530–533. <https://doi.org/10.1109/DSD.2015.79>
- [19] Ronan Lashermes, Héléne Le Boudier, and Gaël Thomas. 2018. Hardware-Assisted Program Execution Integrity: HAPEL. In *NordSec 2018 : 23rd Nordic Conference on Secure IT Systems*. Oslo, Norway. https://doi.org/10.1007/978-3-030-03638-6_25
- [20] Chris Lattner and Vikram Adve. 2004. LLVM: A compilation framework for lifelong program analysis & transformation. In *International symposium on code generation and optimization, 2004. CGO 2004*. IEEE, 75–86.
- [21] Johan Laurent. 2020. *Modélisation de fautes utilisant la description RTL de microarchitectures pour l'analyse de vulnérabilité conjointe matérielle-logicielle*. Theses. Université Grenoble Alpes. <https://tel.archives-ouvertes.fr/tel-03167493>
- [22] Johan Laurent, V. Beroulle, C. Deleuze, Florian Pebay-Peyroula, and Athanasios Papadimitriou. 2018. On the Importance of Analysing Microarchitecture for Accurate Software Fault Models. (2018), 561–564. <https://doi.org/10.1109/DSD.2018.00097>
- [23] Xavier Leroy. 2009. Formal Verification of a Realistic Compiler. *Commun. ACM* 52, 7 (jul 2009), 107–115. <https://doi.org/10.1145/1538788.1538814>
- [24] Jason Lowe-Power, Abdul Mutaal Ahmad, Ayaz Akram, Mohammad Alian, Rico Amslinger, Matteo Andreozzi, Adrià Armejach, Nils Asmussen, Brad Beckmann, Srikant Bharadwaj, et al. 2020. The gem5 simulator: Version 20.0+. *arXiv preprint arXiv:2007.03152* (2020).
- [25] Noura Ait Manssour, Vianney Lapôte, Guy Gogniat, and Arnaud Tisserand. 2022. Processor Extensions for Hardware Instruction Replay against Fault Injection Attacks. In *2022 25th International Symposium on Design and Diagnostics of Electronic Circuits and Systems (DDECS)*. 26–31. <https://doi.org/10.1109/DDECS54261.2022.9770170>
- [26] Sébastien Michelland. 2024. *Replication package for article: From low-level fault modeling (of a pipeline attack) to a proven hardening scheme*. <https://doi.org/10.5281/zenodo.10440364>
- [27] Tanmaya Mishra, Thidapat Chantem, and Ryan Gerdes. 2022. Survey of Control-Flow Integrity Techniques for Real-Time Embedded Systems. *ACM Trans. Embed. Comput. Syst.* 21, 4, Article 41 (oct 2022), 32 pages. <https://doi.org/10.1145/3538275>
- [28] Marie-Laure Potet, Laurent Mounier, Maxime Puits, and Louis Dureuil. 2014. Lazart: A Symbolic Approach for Evaluation the Robustness of Secured Codes against Control Flow Injections. In *2014 IEEE Seventh International Conference on Software Testing, Verification and Validation*. 213–222. <https://doi.org/10.1109/ICST.2014.34>
- [29] Julien Proy, Karine Heydemann, Alexandre Berzati, and Albert Cohen. 2017. Compiler-Assisted Loop Hardening Against Fault Attacks. *ACM Trans. Archit. Code Optim.* 14, 4, Article 36 (dec 2017), 25 pages. <https://doi.org/10.1145/3141234>
- [30] Carlton Shepherd, Konstantinos Markantonakis, Nico van Heijningen, Driss Aboukassimi, Clément Gaine, Thibaut Heckmann, and David Naccache. 2021. Physical fault injection and side-channel attacks on mobile devices: A comprehensive analysis. *Computers & Security* 111 (dec 2021), 102471. <https://doi.org/10.1016/j.cose.2021.102471>
- [31] Stefan Tillich, Christoph Herbst, and Stefan Mangard. 2007. Protecting AES Software Implementations on 32-Bit Processors Against Power Analysis. In *Applied Cryptography and Network Security, 5th International Conference, ACNS 2007, Zhuhai, China, June 5-8, 2007, Proceedings (Lecture Notes in Computer Science, Vol. 4521)*, Jonathan Katz and Moti Yung (Eds.). Springer, 141–157. https://doi.org/10.1007/978-3-540-72738-5_10
- [32] Simon Tollec, Mihail Asavaoe, Damien Courroussé, Karine Heydemann, and Mathieu Jan. 2022. Exploration of Fault Effects on Formal RISC-V Microarchitecture Models. In *2022 Workshop on Fault Detection and Tolerance in Cryptography (FDTC)*. 73–83. <https://doi.org/10.1109/FDTC57191.2022.00017>
- [33] Son Tuan Vu. 2021. *Optimizing Property-Preserving Compilation*. Ph. D. Dissertation. <http://www.theses.fr/2021SORU5435> Thèse de doctorat dirigée par Heydemann, Karine et Cohen, Albert Henri Informatique Sorbonne université 2021.
- [34] Vincent Werner, Laurent Maingault, and Marie-Laure Potet. 2022. An end-to-end approach to identify and exploit multi-fault injection vulnerabilities on microcontrollers. *Journal of Cryptographic Engineering* (Oct. 2022), 1–17. <https://doi.org/10.1007/s13389-022-00292-z>
- [35] Hans Winderix, J. Mühlberg, and F. Piessens. 2021. Compiler-Assisted Hardening of Embedded Software Against Interrupt Latency Side-Channel Attacks. (2021), 667–682. <https://doi.org/10.1109/EuroSP51992.2021.00050>
- [36] Yuan Yao and Patrick Schaumont. 2018. A Low-Cost Function Call Protection Mechanism Against Instruction Skip Fault Attacks. In *Proceedings of the 2018 Workshop on Attacks and Solutions in Hardware Security (Toronto, Canada) (ASHES '18)*. Association for Computing Machinery, New York, NY, USA, 55–64. <https://doi.org/10.1145/3266444.3266453>
- [37] Bilgiday Yuce, Nahid Farhady Ghalaty, Harika Santapuri, Chinmay Deshpande, Conor Patrick, and Patrick Schaumont. 2016. Software Fault Resistance is Futile: Effective Single-Glitch Attacks. In *2016 Workshop on Fault Diagnosis and Tolerance in Cryptography (FDTC)*. 47–58. <https://doi.org/10.1109/FDTC.2016.21>
- [38] Anthony Zgheib, Olivier Potin, Jean-Baptiste Rigaud, and Jean-Max Dutertre. 2022. A CFI Verification System based on the RISC-V Instruction Trace Encoder. In *2022 25th Euromicro Conference on Digital System Design (DSD)*. 456–463. <https://doi.org/10.1109/DSD57027.2022.00067>

Supplementary material for the CC'24 article:

From low-level fault modeling (of a pipeline attack) to a proven hardening scheme

A Operational semantics

In this appendix, we formalize operational semantics for a minimal, representative subset of 32-bit RISC-V assembly which includes fetch skip attacks and our XCCs extension, and provide a detailed proof of the claims made in Section 4.

Notations. We write u_N the type of N-bit unsigned integers. We also use structure-like notations $\{\langle \text{field} \rangle : \langle \text{type} \rangle, \dots\}$ and $\langle \text{struct} \rangle . \langle \text{field} \rangle$ for bit fields and named collections of values.

A.1 Program and execution models

Definition 1 (Instruction, Block, Program).

An **instruction** i is a 16- or 32-bit integer (a u_{16} or a u_{32}) corresponding to a RISC-V opcode.¹⁵ We write $\|i\|$ the size of an instruction in bytes (2 or 4).

A **block** bb is a nonempty sequence of instructions, along with a 4-aligned address value $\text{blockAddr}(bb)$ indicating where it is loaded in memory.

A **program** P is a collection of non-intersecting blocks.

Since instructions in the same block are loaded contiguously, each block spans the interval

$$\text{blockSpan}(bb) = \left[\text{blockAddr}(bb), \text{blockAddr}(bb) + \sum_{i \in bb} \|i\| \right).$$

Programs are assumed to be well-formed, in that no two blocks' spans intersect and every jump points to the blockAddr of some block.

Definition 2 (Program state).

A program state in RISC-V assembler with XCCs and fetch skips is a quadruplet $\langle PC, \rho, \sigma, \alpha \rangle$, where

- $PC : u_{32}$ is the program counter;
- $\rho : u_{32}$ is the last row fetched from code memory;
- $\sigma = \{\text{CCS}, \text{CCSD}, \text{CCSDS}, \text{CCSPROT}\}$ holds XCCs registers:
 - $\text{CCS} : u_{32}$ the running checksum of the current block;
 - $\text{CCSD} : \{E : u_1, JO : u_5\}$ the XCCs status register;
 - $\text{CCSDS} : u_{32}$ the XCCs “delay slot” register;
 - $\text{CCSPROT} : u_{32}$ the protected PC register.
- $\alpha : u_{32}[32]$ is the architectural state (registers $x0 \dots x31$).

The program state tracks the progress of execution. For simplicity, we do not include memory, which would be handled like the registers in α , or side-effects, which could be recorded as a trace whenever a side-effecting syscall is invoked by the `ecall` instruction. We only consider the final state of registers.

¹⁵There is no risk of type confusion because 16- and 32-bit instructions differ on their low bits.

Note that ρ is a raw 32-bit value which is not necessarily a legal instruction because instructions are not always 32-bit and not always aligned in memory.

Definition 3 (Step, Execution).

A **step** is a statement representing the effect of executing one instruction, written $\langle PC, \rho, \sigma, \alpha \rangle \rightarrow r$ where r is either a program state or one of two termination reasons:

- \perp , denoting an exception or crash;
- $\text{end}(\alpha)$, denoting successful completion with final state α .

An **execution** is a sequence of program states ending with a termination reason, such that all pairs of consecutive elements are related by a step:

$$\langle PC, \rho, \sigma, \alpha \rangle \rightarrow \dots \rightarrow \langle PC_n, \rho_n, \sigma_n, \alpha_n \rangle \rightarrow \begin{cases} \perp \\ \text{end}(\alpha') \end{cases}$$

A.2 Fetch and step rules

We can now describe how execution steps are derived from the fetch-decode-execute cycle of the CPU, including fault injections during fetches.

Definition 4 (Fetch).

A **fetch** is a statement $(PC, \rho) a \Rightarrow d (PC', \rho')$ representing the CPU fetching 32 bits at 4-aligned address $a : u_{32}$ and getting value $d : u_{32}$ (“data”), with the PC and ρ members of the program state being updated in the process.

Each fault attack corresponds to one use of the $S_{32(k)}$ or $S\&R_{32}$ rules. Notice, for example, how rule $S_{32(k)}$ skips over a directly to $a + 4k$. In the following rules, $[a]$ denotes the 32 bits stored in memory at address a .

$$\begin{array}{c} \text{NoFAULT} \\ \hline (PC, \rho) a \Rightarrow [a] (PC, [a]) \\ \\ S_{32(k)} \quad 1 < k \leq N \\ \hline (PC, \rho) a \Rightarrow [a + 4k] (PC + 4k, [a + 4k]) \\ \\ S\&R_{32} \quad \rho \neq [a] \\ \hline (PC, \rho) a \Rightarrow \rho (PC, [a]) \end{array}$$

Figure 11. Fetch rules.

From there, we can formalize execution steps as an optional fetch, then the execution of a decoded instruction. This is shown in Figure 12. The $\llbracket \cdot \rrbracket$ and $\llbracket \cdot \rrbracket_{\text{ccs}}$ functions represent the semantics of individual instructions and are specified later. Both are oblivious to ρ so we re-insert ρ into their result after-the-fact using the \bullet function. The functions $\text{LSH}, \text{MSH} : u_{32} \rightarrow u_{16}$ extract the least and most significant halves of a 32-bit value.

These inference rules for steps formalize the patterns shown in Figure 2 (ALIGNED-* is 2.(a) and (b), UNALIGNED-16 is

$$\begin{array}{l}
\text{ALIGNED-16} \\
\frac{\text{PC aligned} \quad (\text{PC}, \rho) \text{ PC} \Rightarrow d(\text{PC}_F, \rho') \\
\sigma.\text{CCSDS} = 0 \quad \text{LSH}(d) \text{ is a 16-bit instruction}^{16}}{\langle \text{PC}, \rho, \sigma, \alpha \rangle \rightarrow \llbracket \text{LSH}(d) \rrbracket (\text{PC}_F, \sigma, \alpha) \bullet \rho'} \\
\\
\text{ALIGNED-32} \\
\frac{\text{PC aligned} \quad (\text{PC}, \rho) \text{ PC} \Rightarrow d(\text{PC}_F, \rho') \\
\sigma.\text{CCSDS} = 0 \quad \text{LSH}(d) \text{ is a 32-bit instruction leader}^{16}}{\langle \text{PC}, \rho, \sigma, \alpha \rangle \rightarrow \llbracket d \rrbracket (\text{PC}_F, \sigma, \alpha) \bullet \rho'} \\
\\
\text{UNALIGNED-16} \\
\frac{\text{PC unaligned} \quad \sigma.\text{CCSDS} = 0 \\
\text{MSH}(\rho) \text{ is a 16-bit instruction}^{16}}{\langle \text{PC}, \rho, \sigma, \alpha \rangle \rightarrow \llbracket \text{MSH}(\rho) \rrbracket (\text{PC}, \sigma, \alpha) \bullet \rho'} \\
\\
\text{UNALIGNED-32} \\
\frac{\text{PC unaligned} \quad \text{MSH}(\rho) \text{ is a 32-bit instruction leader}^{16} \\
\sigma.\text{CCSDS} = 0 \quad (\text{PC}, \rho) \text{ PC} + 2 \Rightarrow d(\text{PC}_F, \rho')}{\langle \text{PC}, \rho, \sigma, \alpha \rangle \rightarrow \llbracket 2^{16} \text{MSH}(\rho) + \text{LSH}(d) \rrbracket (\text{PC}_F, \sigma, \alpha) \bullet \rho'} \\
\\
\text{CHECKSUM-DELAY-SLOT} \\
\frac{\sigma.\text{CCSDS} \neq 0 \quad \text{PC aligned} \\
(\text{PC}, \rho) \text{ PC} \Rightarrow d(\text{PC}_F, \rho')}{\langle \text{PC}, \rho, \sigma, \alpha \rangle \rightarrow \llbracket \sigma.\text{CCSDS} \rrbracket_{\text{ccs}} (\text{PC}_F, \sigma, \alpha, d) \bullet \rho'} \\
\\
\text{Re-inserting } \rho \text{ into the } \langle \text{PC}, \sigma, \alpha \rangle \bullet \rho = \langle \text{PC}, \rho, \sigma, \alpha \rangle \\
\text{result of } \llbracket \cdot \rrbracket \text{ and } \llbracket \cdot \rrbracket_{\text{ccs}}: \quad \perp \bullet \rho = \perp \\
\text{end}(\alpha) \bullet \rho = \text{end}(\alpha)
\end{array}$$

Figure 12. Step rules.

2.(c) and UNALIGNED-32 is 2.(d)). For instance, in UNALIGNED-32, $\text{MSH}(\rho)$ identifies the two remaining bytes in the previously-fetched line as 32-bit instruction leader, so the rule fetches the next line d at $\text{PC} + 2$ and recombines ρ with $\text{LSH}(d)$ to form the complete 32-bit instruction and run it.

The extra rule CHECKSUM-DELAY-SLOT describes the special treatment of the 32-bit checksum value following an XCCS instruction, which is triggered by $\sigma.\text{CCSDS} \neq 0$ and is described later with the semantics of XCCS instructions.

A.3 Instruction semantics

We now define instruction semantics through the functions

$$\begin{array}{l}
\llbracket \cdot \rrbracket : (\text{PC}, \sigma, \alpha) \mapsto \langle \text{PC}', \sigma', \alpha' \rangle \text{ or } \perp \text{ or } \text{end}(\alpha') \\
\llbracket \cdot \rrbracket_{\text{ccs}} : (\text{PC}, \sigma, \alpha, d) \mapsto \langle \text{PC}', \sigma', \alpha' \rangle \text{ or } \perp \text{ or } \text{end}(\alpha')
\end{array}$$

The two functions are similar in nature; $\llbracket \cdot \rrbracket_{\text{ccs}}$ is used when re-running XCCS instructions after their 32-bit checksum value (passed as 4th argument) has been found.

In the following, we assume a function `encode` which encodes assembler notation into 16/32-bit values (16-bit for mnemonics starting with “c.”, 32-bit otherwise), following

¹⁶As per the RISC-V ISA [1], a 32-bit instruction leader is a value $i : u_{16}$ such that $i \equiv 3 [4]$ and a 16-bit instruction is any other 16-bit value.

the RISC-V ISA [1]. Unspecified members in the output σ'/α' are implicitly kept from the input σ/α .

Instructions that are not jumps, traps, or XCCS instructions have their natural semantics, with two small changes: they cannot run when CCSPROT is set to a non-zero address, and they update CCS as a side-effect. We illustrate this category with the `add` and `c.add` instructions, and the `c.nop` instruction that we use as padding while hardening.

- $i = \text{encode}(\text{add } r_d, rs_1, rs_2)$,
 $i = \text{encode}(\text{c.add } r_d, rs_1, rs_2) :$

$$\llbracket i \rrbracket (\text{PC}, \sigma, \alpha) = \begin{cases} \perp & \text{if } \sigma.\text{CCSPROT} \neq 0 \\ \langle \text{PC} + \llbracket i \rrbracket, \sigma', \alpha' \rangle & \text{otherwise, with} \\ & \sigma'.\text{CCS} = \sigma.\text{CCS} + \text{realign}(\text{PC}, u_{32}(i)) \\ & \alpha'[r_d] = \alpha[rs_1] + \alpha[rs_2] \end{cases}$$

- $i = \text{encode}(\text{c.nop}) :$

$$\llbracket i \rrbracket (\text{PC}, \sigma, \alpha) = \begin{cases} \perp & \text{if } \sigma.\text{CCSPROT} \neq 0 \\ \langle \text{PC} + \llbracket i \rrbracket, \sigma', \alpha \rangle & \text{otherwise} \\ & \sigma'.\text{CCS} = \sigma.\text{CCS} + \text{realign}(\text{PC}, u_{32}(i)) \end{cases}$$

The CCS update adds the instruction’s opcode to the CCS register, adjusted for PC alignment:¹⁷

$$\text{realign}(\text{PC}, i) = \begin{cases} i & \text{if } \text{PC} \equiv 0 [4] \\ 2^{16} \text{LSH}(i) + \text{MSH}(i) & \text{otherwise} \end{cases}$$

Trapping instructions include `ebreak/c.ebreak` which we treat as non-recoverable, and `ecall` which invokes a syscall whose number is specified in register `a7 (x17)`. We treat an invocation of `syscall __NR_exit (1)` as program termination.¹⁸

- $i = \text{encode}(\text{ebreak})$, $i = \text{encode}(\text{c.ebreak}) :$

$$\llbracket i \rrbracket (\text{PC}, \sigma, \alpha) = \perp$$

- $i = \text{encode}(\text{ecall}) :$

$$\llbracket i \rrbracket (\text{PC}, \sigma, \alpha) = \begin{cases} \text{end}(\alpha) & \text{if } \alpha[17] = 1 \text{ and } \sigma.\text{CCSPROT} = 0 \\ \perp & \text{otherwise} \end{cases}$$

Among standard instructions, jumps are the most heavily affected. We show here the cases of `jal (jump-and-link)`, which in RISC-V is used for unconditional branches and function calls, and `beq (branch-if-equal)`, a representative conditional branch. The extension to `jalr` (which is used to return from functions) is immediate by replacing the target $\text{PC} + \text{imm}$ of `jal` with the value of a register in α .

- $i = \text{encode}(\text{jal } r_d, \text{imm}) :$

¹⁷Computing the sum of realigned instruction opcodes ends up being equivalent to summing the values returned by fetch rules in the program, see Lemma 3. We choose to update during decoding because of our intuition that a CCS update during the fetch cycle would be impacted by the fault.

¹⁸We provide this exit mechanism to have a reasonable program model, but in practice the `exit()` function is in the non-protected lib, so we don’t actually worry about protecting the last block.

$$\llbracket i \rrbracket(\text{PC}, \sigma, \alpha) = \begin{cases} \perp & \text{if } \sigma.\text{CCSPROT} \neq \text{PC} \\ \langle \text{PC} + \text{imm}, \sigma', \alpha' \rangle & \text{otherwise, with} \\ & \sigma'.\text{CCS} = 0 \\ & \sigma'.\text{CCSPROT} = 0 \\ & \sigma'.\text{CCSD.JO} = 0 \\ & \alpha'[r_d] = \text{PC} + \llbracket i \rrbracket + 2 \times \sigma.\text{CCSD.JO} \end{cases}$$

There are two changes compared to traditional jumps: (1) it is only allowed when CCSPROT is set to PC; (2) it accounts for the *jump offset* (JO) field of CCSD to jump over trap barriers after function calls. The jump also prepares the execution of the next block by resetting CCS and other registers to 0. We do not allow the compiler to compress `jal` into the 16-bit version `c.jal` to make sure the next block is also aligned.

Conditional branches are similar, except that if the branch is not taken the block continues and in particular CCS is updated rather than reset.

- $i = \text{encode}(\text{beq } rs_1, rs_2, \text{imm})$:

$$\llbracket i \rrbracket(\text{PC}, \sigma, \alpha) = \begin{cases} \perp & \text{if } \sigma.\text{CCSPROT} \neq \text{PC} \\ \langle \text{PC} + \text{imm}, \sigma', \alpha \rangle & \text{if } \alpha[rs_1] = \alpha[rs_2], \text{ with} \\ & \sigma'.\text{CCS} = 0 \\ & \sigma'.\text{CCSPROT} = 0 \\ & \sigma'.\text{CCSD.JO} = 0 \\ \langle \text{PC} + \llbracket i \rrbracket, \sigma', \alpha \rangle & \text{otherwise, with} \\ & \sigma'.\text{CCS} = \sigma.\text{CCS} + \text{realign}(\text{PC}, u_{32}(i)) \end{cases}$$

The last type of instruction is XCCS instructions. Because these have both a 32-bit opcode and a 32-bit argument, they cannot be executed in a single step. A “delay slot” mechanism with $\sigma.\text{CCSDS}$ and the CHECKSUM-DELAY-SLOT rule is used to solve this issue in two steps. At the first step, the XCCS opcode will be fetched and recorded in $\sigma.\text{CCSDS}$. At the second step, the checksum value will be fetched, and passed as argument to the appropriate semantics function $\llbracket \sigma.\text{CCSDS} \rrbracket_{\text{ccs}}$.

The first step which records the opcode in CCSDS proceeds as follows:

- $i = \text{encode}(\text{ccs}), i = \text{encode}(\text{ccsb}),$
 $i = \text{encode}(\text{ccscall } N), i = \text{encode}(\text{ccscallb } N) :$

$$\llbracket i \rrbracket(\text{PC}, \sigma, \alpha) = \begin{cases} \perp & \text{if PC is unaligned or } \sigma.\text{CCSPROT} \neq 0 \\ \langle \text{PC} + \llbracket i \rrbracket, \sigma', \alpha \rangle & \text{otherwise, with} \\ & \sigma'.\text{CCS} = \sigma.\text{CCS} + \text{realign}(\text{PC}, u_{32}(i)) \\ & \sigma'.\text{CCSDS} = i \end{cases}$$

The second step, which always uses the CHECKSUM-DELAY-SLOT rule, reads back from CCSDS and calls the $\llbracket \rrbracket_{\text{ccs}}$ function by passing the checksum d just fetched from memory as an extra 4th argument. All 4 XCCS instructions have similar semantics; the `-call` variants set $\text{CCSD.JO} = N$ and the `-b` variants flip the least significant bit of the checksum value before comparing. We factor them as follows.

- Checksum step (d is a checksum value):

$$\begin{aligned} \llbracket \text{encode}(\text{ccs}) \rrbracket_{\text{ccs}} &= \text{check}(0, 0) \\ \llbracket \text{encode}(\text{ccsb}) \rrbracket_{\text{ccs}} &= \text{check}(1, 0) \\ \llbracket \text{encode}(\text{ccscall } N) \rrbracket_{\text{ccs}} &= \text{check}(0, N) \\ \llbracket \text{encode}(\text{ccscallb } N) \rrbracket_{\text{ccs}} &= \text{check}(1, N) \end{aligned}$$

$$\text{check}(\text{mask} : u_{32}, \text{JO} : u_5)(\text{PC}, \sigma, \alpha, d) = \begin{cases} \perp & \text{if } \sigma.\text{CCS} \neq (d \oplus \text{mask}) \text{ or } \sigma.\text{CCSPROT} \neq 0 \\ \perp & \text{if } d \text{ is not a valid checksum literal} \\ \langle \text{PC} + 4, \sigma', \alpha \rangle & \text{otherwise, with} \\ & \sigma'.\text{CCSPROT} = \text{PC} + 4 \\ & \sigma'.\text{CCSDS} = 0 \\ & \sigma'.\text{CCSD.JO} = \text{JO} \end{cases}$$

The purpose of the “checksum literals” check is to close attack vectors in which an attacker feeds an incorrect checksum value by injecting a fault between the XCCS opcode and its checksum parameter. We avoid these by preventing a 32-bit value d from appearing verbatim as a checksum value if:

- it decodes as a jump instruction;
- it decodes as an XCCS instruction;
- it decodes as a pair of `c.ebreak`.

The countermeasure gets rid of such values by flipping their LSB, which is shown to be correct in Lemma 9.

B Hardening algorithm and proof of security

B.1 Structure of hardened programs

The security property that we want to establish is that in the execution of a hardened program, any single-fault injection leads to termination with \perp before the end of the current block. This relies on a particular structure for blocks, which we also use to formalize the hardening process.

Definition 5 (Source block).

A **source block** is a block whose sequence of instructions consists of:

- 0 or more **straight** (add, ebreak, ecall...) or **conditional branch** (beq) instructions, excluding c.ebreak and Xccs instructions; followed by
- One **unconditional branch instruction** (jal).

Unlike classical definitions, **we do not count conditional branches as block terminators**. This improves the correspondence between blocks and checksum regions; in particular, it ensures the invariant that CCS resets to 0 at the start of every block.

Definition 6 (Hardened block).

A **hardened block** is a block hb obtained by hardening a source block, i.e. $hb = \text{HARDEN}(sb, N)$ for some source block sb. HARDEN is detailed as Algorithm 2 (which is a copy of Algorithm 1).

The hardening algorithm processes instructions individually. All the original instructions are kept. Before each jump, a checksum check is added by procedure addChecksum(), in the form of an Xccs instruction (ccscall for function calls, ccs otherwise) followed by the reference checksum value. In order to close certain attack vectors, the checksum value must not decode as a jump or Xccs instruction; if that happens, the LSB of the checksum value is flipped and the Xccs instruction is replaced by its -b variant. Finally, a barrier of c.ebreak instructions is added at the end.

Now, conditional branches in the middle of a block allow for early exits, but don't reset the checksum. The idea is that upon exiting the block, the expected checksum value is always the sum of all lines from the beginning of the block to the exit instruction. It will nonetheless be useful to split the blocks at each exit point to help formalization.

Lemma 1 (Structure of source blocks).

A source block sb can be uniquely decomposed into $m \geq 0$ early sections followed by one final section, i.e.

$$sb = se_1 + \dots + se_m + sf \quad (\text{"+" is concatenation})$$

where

- each se_k (**source early**, $1 \leq k \leq m$) section consists of straight instructions followed by one conditional branch;

Algorithm 2 Algorithm: HARDEN

Input: A source block $[i_1, \dots, i_n]$

Input: Upper bound N for $S_{32}(k)$ faults each fetch

Output: A hardened block hb.

```

hb ← []
sum : u32 ← 0
offset ← 0 ▷ Blocks are 4-aligned.
procedure addToBlock(i)
  hb.append(i)
  sum ← sum + realign(offset, u32(i))
  offset ← offset + ||i||
procedure addChecksum(i, ib)
  if sum + i is not a valid checksum literal then
    addToBlock(ib)
    addToBlock(sum ⊕ 1)
  else
    addToBlock(i)
    addToBlock(sum)
for i in  $[i_1, \dots, i_n]$  do
  if i is a jump instruction then
    if offset = 0 then ▷ No empty sections.
      addToBlock(encode(nop))
    else if offset ≡ 2 [4] then ▷ Force alignment.
      addToBlock(encode(c.nop))
    if i is a function call then
      addChecksum(encode(ccscall (2N + 4)),
        encode(ccscallb (2N + 4)))
    else ▷ Conditional/unconditional jumps.
      addChecksum(encode(ccs), encode(ccsb))
      offset ← 0
      addToBlock(i)
  for j = 1 to 2N + 4 do ▷ Add trap barrier.
    addToBlock(encode(c.ebreak))
return hb

```

- the sf (**source final**) section consists of straight instructions followed by the block's unconditional jump.

Proof. Because there is exactly one conditional branch per se_k and none other, m must be the number of conditional branches in sb. The grouping is straightforward from here. \square

Lemma 2 (Structure of hardened blocks).

A hardened block $hb = \text{HARDEN}(se_1 + \dots + se_m + sf, N)$ can be uniquely decomposed into 4-aligned sections

$$hb = he_1 + \dots + he_m + hf$$

where

- each he_k (**hardened early**, $1 \leq k \leq m$) section consists of:
 1. the straight instructions of se_k ;
 2. an optional c.nop or nop, in a way that 1) and 2) combined are not empty;

3. a 4-aligned ccs or ccsb;
4. a 4-aligned 32-bit checksum value;
5. the conditional branch of se_k ;
- the hf (**hardened final**) section consists of:
 1. the straight instructions of hf;
 2. an optional c.nop or nop, in a way that 1) and 2) combined are not empty;
 3. a 4-aligned cscall or cscallb if the terminator is a function call, a ccs or ccsb otherwise;
 4. a 4-aligned 32-bit checksum value;
 5. the terminator of sf;
 6. $2N + 4$ c.ebreak instructions.

Proof. Decomposition: The main (**for i**) loop in HARDEN is a morphism for concatenation. Take for he_k/hf the hardened sequences corresponding to se_k/sf in the input; structurally, we get $hb = he_1 + \dots + he_m + hf$.

Section contents: Straight instructions are clearly preserved by HARDEN. For both types of sections, items 2 through 5 are generated by the handling of the section's jump by HARDEN. The trap barrier is always added at the end of the block; we can count it towards hf.

Alignment: First, all sections have a length that is a multiple of 4. This constraint is ensured by adding a c.nop before section's branch if the straight instructions finish on an unaligned address. The following Xccs instruction, checksum value, and branch instruction all occupy 4 bytes each. Then, because the block hb itself is 4-aligned and the sections follow each other in memory, each section individually must be 4-aligned. \square

Finally, we give an intuitive view of the checksum values by proving that the per-instruction summing process with realignment is equivalent to summing 4-aligned lines. We can do this by comparing both computations on the underlying sequences of u_{16} values.

Lemma 3 (realign sum computes the sum of fetched lines).

Let $h = (h_0, \dots, h_{2n-1}) : u_{16}^{2n}$ a sequence of u_{16} ("halfwords").

Assume $(w_0, \dots, w_{n-1}) : u_{32}^n$ is a sequence of u_{32} ("words") whose concatenation is h , i.e. such that $w_i = h_{2i} + 2^{16}h_{2i+1}$.

Let (i_0, \dots, i_{m-1}) a sequence of instructions, and write $offset_j = \frac{1}{2} \sum_{k=0}^{j-1} \|i_k\|$ the offset of instruction j in the sequence (in 16-bit units). Assume the concatenation of this sequence is also h , i.e. $offset_m = 2n$ and for all j ,

- $i_j = h_{offset_j}$ if i_j is a 16-bit instruction;
- $i_j = h_{offset_j} + 2^{16}h_{offset_j+1}$ otherwise.

Then,

$$\sum_{j=0}^{m-1} \text{realign}(2 \cdot \text{offset}_j, u_{32}(i_j)) = \sum_{i=0}^{n-1} w_i.$$

Proof. First, rewrite every instruction's realigned contribution in terms of h_i ; specifically, we have

$$\forall j, \text{realign}(2 \cdot \text{offset}_j, u_{32}(i_j)) = \sum_{i=\text{offset}_j}^{\text{offset}_{j+1}-1} 2^{16(i \bmod 2)} h_i.$$

This we can show by distinguishing four cases:

- i_j is 16-bit, $offset_j \bmod 2 = 0$:
 - LHS is i_j
 - RHS is $h_{offset_j} = i_j$
- i_j is 16-bit, $offset_j \bmod 2 = 1$:
 - LHS is $2^{16}i_j$
 - RHS is $2^{16}h_{offset_j} = 2^{16}i_j$
- i_j is 32-bit, $offset_j \bmod 2 = 0$:
 - LHS is i_j
 - RHS is $h_{offset_j} + 2^{16}h_{offset_j+1} = i_j$
- i_j is 32-bit, $offset_j \bmod 2 = 1$:
 - LHS is $MSH(i_j) + 2^{16}LSH(i_j)$
 - RHS is $2^{16}h_{offset_j} + h_{offset_j+1} = 2^{16}LSH(i_j) + MSH(i_j)$

Now, by summing from $j = 0$ to $m - 1$, we get

$$S := \sum_{j=0}^{m-1} \text{realign}(offset_j, u_{32}(i_j)) = \sum_{i=0}^{offset_m-1} 2^{16(i \bmod 2)} h_i.$$

Recalling that $offset_m = 2n$, this equivalent to unfolding w_i into a sum of u_{16} :

$$\sum_{i=0}^{n-1} w_i = \sum_{i=0}^{2n-1} 2^{16(i \bmod 2)} h_i = S. \quad \square$$

B.2 Program state upon leaving a hardened block

We will now show that the counter-measure ensures good properties in the execution of hardened blocks, namely that they can only be exited by a legitimate jump instruction and while validating the checksum associated to that jump.

In the following, we assume a program P and a successful execution $e = [s_1, \dots, s_{|e|}]$, with each step being written

$$s_i = \langle PC_i, \rho_i, \sigma_i, \alpha_i \rangle \rightarrow r_i,$$

which structurally implies that

$$\begin{cases} r_i = \langle PC_{i+1}, \rho_{i+1}, \sigma_{i+1}, \alpha_{i+1} \rangle & \text{if } i < |e|; \\ r_{|e|} = \text{end}(\alpha_{|e|}). \end{cases}$$

The first step in the proof is to show that the trap barrier prevents the end of a hardened block from being reached, meaning that any exit must happen through a jump instruction.

Lemma 4 (Hardened blocks must be exited by jumps).

Let hb be a hardened block of P , and assume the execution enters hb , meaning that there exists t ("top") such that $PC_t = \text{blockAddr}(hb)$.

Let n be the number of instructions executed without leaving hb , i.e. the largest number such that the following all hold:

1. $t + n \leq |e|$;
2. $\forall i \in [t, t + n)$, $PC_{i+1} \in \text{blockSpan}(\text{hb})$;
3. $\forall i \in [t, t + n)$, s_i is not a jump instruction.¹⁹

Then, either

- $t + n = |e|$ (hence s_{t+n} is a successful termination by an invocation of the `exit syscall`), or
- $t + n < |e|$ and s_{t+n} is a jump instruction.

Proof. By contradiction. An exit not covered by the lemma statement must be with $t + n < |e|$ and s_{t+n} not a jump instruction. As such, there must be no jump instruction in the entire sub-sequence $s_t \dots s_{t+n}$.

For a non-jumping step s_i , we always have $PC_{i+1} > PC_i$. This is because the only two changes to PC are the update by the instruction's semantics (which always adds $\|i\| > 0$) and the potential skip from the $S_{32(k)}$ fetch rule (which is $4k \geq 0$).

Therefore, $(PC_i)_{t \leq i \leq t+n}$ is a strictly increasing sequence, which means that every address in $\text{blockSpan}(\text{hb})$ is part of exactly one $[PC_i, PC_{i+1})$ interval (i.e. there is exactly one step that consumes it as part of its execution).

So then, we can look at hb 's ending, which has the following 4-aligned sequence (where each line covers 4 bytes of code):

```
ccs/ccsb/ccscall/ccscallb
<checksum value>
j32
```

L: `c.ebreak`; `c.ebreak # (repeat N+2 times)`

The `j32` marker represents the jump instruction, which is always a 32-bit instruction. The trap barrier starts at address L and contains $2N + 4$ `c.ebreak` instructions, therefore the block ends at $L + 4N + 8$.

From the previous argument, there is exactly one step s_i such that address L lies in the interval $[PC_i, PC_{i+1})$, meaning that the line at address L is loaded or skipped during the execution of step s_i . Step s_i starts at PC_i , which can be written in a unique way as either $PC_i = L - 4k$ or $PC_i = L - 4k - 2$ depending on its alignment.

Now by analyzing all cases for the value of k and alignment of PC, we can show that the execution always crashes at or shortly after s_i , completing the contradiction argument.

- $k > 0$. This case is only possible if step s_i uses the fetch rule $S_{32(k')}$ because $PC_{i+1} > L$, yet all other options for non-jump instructions result in $PC_{i+1} \leq PC_i + 4$. Since $k' \leq N$, the fetch will not reach beyond the trap barrier, so the fetch returns a pair of `c.ebreak` instructions. These then get consumed by s_i 's step rule, which we can identify from the alignment of PC.
 - PC aligned: s_i must then be using either `ALIGNED-16` or `CHECKSUM-DELAY-SLOT`. In the first case, the first `c.ebreak` fetched by s_i is executed, leading to a trap. In the second

case, the checksum check also traps because double `c.ebreak` is not a valid checksum literal.

- PC unaligned: since there was a fetch, `UNALIGNED-32` is being used. A 32-bit instruction is composed from `MSH(ρ_i)` and the first `c.ebreak`. Recall that s_i is not a jump, so this instruction either crashes or leads into s_{i+1} with PC_{i+1} unaligned. In this case, s_{i+1} will use `UNALIGNED-16` and run the second `c.ebreak`, trapping as claimed.
- $k = 0$. In this case, PC_i must be either L or $L - 2$.
 - $PC_i = L$ (aligned). Here, s_i can either use rule `ALIGNED-16`, `ALIGNED-32`, or `CHECKSUM-DELAY-SLOT`, all of which require a fetch that can only return $\rho_i = \text{encode}(j32)$ (if `S&R32` is used) or a pair of `c.ebreak`. Using `ALIGNED-16` and `ALIGNED-32` would either crash from running a `c.ebreak` or contradict the hypothesis that s_i is not a jump. Using `CHECKSUM-DELAY-SLOT` would also crash because neither fetch result is a valid checksum literal.
 - $PC_i = L - 2$ (unaligned). Because there is a fetch, s_i must be using rule `UNALIGNED-32`. Much like the similar case in $k > 0$, the instruction recomposed from `MSH(ρ_i)` and `LSH(d)` cannot be a jump, so if it doesn't crash execution continues to s_{i+1} with the new `LSH(ρ_{i+1})`.

If s_i uses fetch rule `NOFAULT` or $S_{32(k')}$, then ρ_{i+1} is a pair of `c.ebreak`, therefore s_{i+1} will use `UNALIGNED-16` and crash running the second `c.ebreak`.

Execution can only proceed past s_{i+1} if s_i uses fetch rule `S&R32`, in which case s_{i+1} is a repeat of s_i 4 bytes later ($PC_{i+1} = L + 2$). The analysis of s_i can be repeated, with two changes. First, $S_{32(k')}$ could now reach 4 bytes further, up to $L + 4 + 4N$, which hits the last 4 bytes of the barrier. This shows why $2N + 4$ `c.ebreak` are required. Second, using `S&R32` no longer succeeds because ρ_{i+1} is now also a pair of `c.ebreak`. Therefore, s_{i+2} crashes.

This guaranteed crash implies that a successful execution can only exit a protected block in one of the ways described by the theorem statement. \square

The next objective is to show that the `Xccs` instruction, checksum and jump instruction sequence is secure in that only *legitimate* jumps can be used to exit a hardened block. A key part of this is that only original instructions must be executed around the time of the jump, not forged instructions. This synchronization is guaranteed by the fact that an `Xccs` instruction is needed to jump, and `Xccs` instructions *cannot* be forged.

Lemma 5 (`Xccs` instructions can be repeated but not forged).

Let $\text{hb} = [i_1, \dots, i_{|\text{hb}|}]$ be a hardened block and

$$s = \langle PC, \rho, \sigma, \alpha \rangle \rightarrow \langle PC', \rho', \sigma', \alpha' \rangle$$

a step that successfully runs an `Xccs` opcode within hb , i.e. such that $[PC, PC'] \subseteq \text{blockSpan}(\text{hb})$. Further assume that we don't repeat upon landing in hb , i.e. either s doesn't use rule `S&R32` or $PC \geq \text{blockAddr}(\text{hb}) + 4$.

¹⁹For now we consider a forged jump that jumps back somewhere into the block as "leaving" the block; we will later show that this cannot happen.

Then, the instruction executed by s is an original instruction of hb , in the sense that there is an XCCS instruction $i_j \in hb$ (with address $L = \text{blockAddr}(hb) + \sum_{k=0}^{j-1} \|i_k\|$) such that one of the following is true:

1. s uses fetch rule NOFAULT and $PC = L$;
2. s uses fetch rule S32(k) and $PC = L - 4k$;
3. s uses fetch rule S&R32 and $PC = L + 4$.

Proof. XCCS instructions trap when executed with unaligned PC, so they can only be run by step rule ALIGNED-32. This rule always uses the value d returned by a fetch as an opcode, which is a 4-aligned value found in program memory at an address that depends on the fetch rule:

- With NOFAULT, d is the value at PC;
- With S32(k), d is the value at $PC + 4k$;
- With S&R32, d is the value at $PC - 4$.

Given the hypotheses and the length of the trap barrier, this value must be a 4-aligned value within hb , which can intersect the instruction sequence in three ways:

1. d matches a 32-bit instruction $i_j \in hb$: then the lemma is obviously true.
2. LSH(d) is the last two bytes of an instruction $i_j \in hb$ and MSH(d) is the first two bytes of i_{j+1} . This is impossible because d is an XCCS opcode so $MSH(d) = 0$. The 16-bit zero value is reserved in the RISC-V ISA as an invalid opcode, thus i_{j+1} could not be a valid instruction.
3. d matches a 32-bit checksum argument to an XCCS instruction: this is also impossible because XCCS opcodes are not valid checksum literals.

Since only case 1) is possible, s_t must indeed be executing a legitimate instruction $i_j \in hb$ (maybe after a fault). \square

Definition 7 (Legitimate entry).

A state $\langle PC, \rho, \sigma, \alpha \rangle$ is called a **legitimate entry** into a block bb if $PC = \text{blockAddr}(bb)$, $\sigma.CCSPROT = \sigma.CCSDS = 0$, and ρ does not decode as an XCCS instruction.

Definition 8 (Legitimate execution of a block).

A sequence of steps $s_t \dots s_b$ (“top”, “bottom”) is defined as a **legitimate execution of a hardened block hb** if the following conditions are met:

- s_t 's initial state is a legitimate entry into hb ;
- $\forall i \in [t, b)$, $PC_{i+1} \in \text{blockSpan}(hb)$;
- s_b is the first and only instruction in the sequence to be either a taken jump or an invocation of the `exit syscall`.

Lemma 6 (Jumps out of hardened blocks must be legitimate).

Let $hb = [i_1, \dots, i_{|hb|}]$ a hardened block of P , and $s_t \dots s_b$ a legitimate execution of hb where s_b is a jump. Then:

1. $b \geq t + 2$ and s_{b-2} is an XCCS instruction;
2. The jump is legitimate, i.e. there is an instruction $i_j \in hb$ such that $PC_b = \text{blockAddr}(hb) + \sum_{k=0}^{j-1} \|i_k\|$ and step s_b executes the opcode i_j .

3. The checksum is correct when leaving the block, i.e.

$$\sigma_b.CCS = \begin{cases} [PC_b - 4] & \text{if it's a valid checksum literal;} \\ [PC_b - 4] \oplus 1 & \text{otherwise.} \end{cases}$$

4. s_b 's final state is a legitimate entry into another block.

Proof. Walking back from the last few instructions, for s_b to be a jump and not trap, we must have $\sigma_b.CCSPROT \neq 0$. Therefore $b > t$ and s_{b-1} must use the CHECKSUM-DELAY-SLOT rule, as no other type of step can end with $CCSPROT \neq 0$. This implies that $\sigma_{b-1}.CCSDS \neq 0$, so once again $b - 1 > t$ and s_{b-2} must execute an XCCS instruction, because CCSDS can only be non-zero for one step and no other instruction sets it.

By Lemma 5, there must be an original XCCS instruction from hb being executed by s_{b-2} , and such instructions are only found at the end of early or final sections, which have the following structure:

```
L:      ccs/ccsb/ccscall/ccscallb
L+4:    <checksum value>
L+8:    j32
L+12:   # ... next section or c.ebreak ...
```

Note that while s_{b-2} runs the opcode found in memory at address L , PC_{b-2} is not guaranteed to be L since faults might be involved. The lemma focuses on showing that even then, in all execution scenarios starting from s_{b-2} the legitimate jump at $L + 8$ must be taken by s_b while also passing the checksum at $L + 4$.²⁰

First note that, as XCCS instructions set $CCSPROT = PC + 8$ and both s_{b-2} and s_b already increment PC by 4, neither s_{b-1} nor s_b can fetch with S32(k) as that would cause $PC > CCSPROT$ in s_b and trap. The only variable after s_{b-2} is whether the S&R32 fetch rule is used to change either fetched value.

- If s_{b-2} uses either of the fetch rules NOFAULT and S32(k), then $PC_{b-1} = L + 4$ and $PC_b = L + 8$.

Then, s_{b-1} uses the CHECKSUM-DELAY-SLOT rule after fetching a checksum literal. Attacking the fetch with S&R32 would return the XCCS opcode of s_{b-2} , and trap because it's not a valid checksum literal. Thus s_{b-1} must fetch with NOFAULT, passing the checksum.

Finally, s_b performs the jump. Fetching with S&R32 is again impossible as that would execute the checksum literal, which cannot be a jump. Therefore, s_b fetches with NOFAULT, and executes the intended jump at $L + 8$ while passing the intended checksum.

- If s_b uses the fetch rule S&R32, then $PC_b = L + 4$, which means that $PC_{b-1} = L + 8$ and $PC_b = L + 12$.

This time, s_{b-1} must fetch with S&R32 to obtain the intended checksum, as a NOFAULT fetch would return the jump opcode, which is not a valid checksum literal.

²⁰The analysis in this lemma remains correct even if rule S32(k) caused the execution of extra nop instructions, as that would reset CCSPROT to 0 and just prevent the jump altogether.

Then, s_b must also fetch with S&R32 to get the jump, since a NoFAULT fetch would the value at $L + 12$, which cannot be a jump.

Notice that in the second case the entire jump widget is “delayed” by 4 bytes by repeated attack with S&R32. This is still a legitimate exit for the purpose of this lemma. We will show later that in single-fault cases, the checksum would not pass in such a scenario.

Now, since the jump was executed properly, it is fairly easy to show that s_b is a legitimate entry into another block.

- The target of the jump is the `blockAddr` of another block by assumption about the validity of the entire program P .
- After the jump, we have $\sigma_b.CCSPROT = 0$ and $\sigma_b.CCSDS = 0$ as a result of the semantics of s_{b-1} and s_b .
- Finally, ρ_{b+1} is either a jump, a `c.ebreak`, or the first line of the next section (which is never an `XCCS` instruction). \square

B.3 Security guarantees

This key lemma implies that legitimate entries and exits from hardened blocks is an execution invariant. We can sum this up in the most general (multi-fault) case as every block guaranteeing its checksum must pass.

Theorem 1 (Security guarantee for multi-fault executions).

Let P a fully hardened program and $e = [s_0, \dots, s_{|e|}]$ an execution such that

- s_0 is a legitimate entry into a block of P ;
- $s_{|e|}$ ends successfully, returning some $\text{end}(\alpha)$.

Then there exists a sequence $[\text{hb}_1, \dots, \text{hb}_m]$ of blocks of P such that

1. The execution e can be partitioned into subsequences $(s_{t_i} \dots s_{b_i})_{1 \leq i \leq m}$ each a legitimate execution of hb_i ;
2. Each s_{b_i} ($i \neq m$) is a legitimate jump of hb_i and $\sigma_{b_i}.CCS$ is the correct checksum associated with that jump.

Proof. By induction, constructing (t_i) , (b_i) and (hb_i) along the way.

Assume s_{t_i} is a legitimate entry into hb_i (which is true for $i = 0$). By Lemma 4, there is b_i such that $s_{t_i} \dots s_{b_i}$ is an execution of hb_1 , with s_{b_i} leaving either by successful termination or by a jump. If s_{b_i} ends the program, we’re done. Otherwise, by Lemma 6, s_{b_i} is a legitimate jump out of hb_i with $\sigma_{b_i}.CCS$ passing the associated checksum, and its final state is a legitimate entry into another block hb_{i+1} . Define $t_{i+1} := b_i + 1$ and start over. \square

As we’ve discussed previously, this doesn’t completely rule out attacks, because for some blocks there exist multi-faulted paths whose checksum collides with the expected checksum. However, such paths do not exist when a single fault is injected during the execution of a block. To show

this, we can go back to the sum-of- u_{32} expression of the checksum and reason on the series of fetches that builds it.

Definition 9 (Fetch trace).

The *fetch trace* of a subsequence $s_t \dots s_b$ of e is the tuple

$$(N, S, R, \rho_{\text{entry}}, (a_i), (k_i), \text{next})$$

where fetches are partitioned by rule type:

- $N = \{i \mid s_i \text{ uses the NoFAULT rule}\}$
- $S = \{i \mid s_i \text{ uses the S32}(k) \text{ rule}\}$
- $R = \{i \mid s_i \text{ uses the S\&R32 rule and } i \neq t\}$
- $\rho_{\text{entry}} = \rho_t$ if s_t uses the S&R32 rule, 0 otherwise

and relevant information is recorded as follows:

- for $i \in N \cup S \cup R$, a_i is the address fetched by s_i ;
- for $i \in S$, k_i is the parameter to $S32(\cdot)$ (0 for other i);
- `next` gives the next step that includes a fetch, i.e.

$$\text{next}(i) = \min \{j \in N \cup S \cup R \mid j > i\}.$$

Lemma 7 (Relation between checksum and fetch trace).

Let hb a hardened block of P , $s_t \dots s_b$ a legitimate execution of hb and $(N, S, R, \rho_{\text{entry}}, (a_i), (k_i), \text{next})$ the fetch trace of $s_t \dots s_{b-2}$ (only up to the last `XCCS` instruction that counts towards the checksum).

Then fetched addresses are contiguous in the sense that

$$\forall i < b - 2, a_{\text{next}(i)} = a_i + 4k_i + 4$$

and the checksum upon leaving the block is

$$\sigma_{b-1}.CCS = \sum_{i \in N} [a_i] + \sum_{i \in S} [a_i + 4k_i] + \sum_{i \in R} [a_i - 4] + \rho_{\text{entry}}.$$

Proof. The checksum upon leaving the block is $\sigma_b = \sigma_{b-1}$ (since `CCS` is not updated during checks and jumps). By definition of instructions’ semantics, σ_{b-1} is the realigned sum of the instructions executed by $s_t \dots s_{b-2}$ (i.e. the values being passed to $\llbracket \cdot \rrbracket$ or $\llbracket \cdot \rrbracket_{\text{CCS}}$ in each step rule).

Because the sequence takes no branch, each instruction i ends with $\text{PC} \leftarrow \text{PC} + \llbracket i \rrbracket$ so each new fetch queries the data immediately following the previous fetch. Accounting for skips in $S32(k)$ means that $a_{\text{next}(i)} = a_i + 4k_i + 4$ (k_i being 0 for other rules).

In addition, the sequence starts on a 4-aligned boundary and ends with a 4-aligned `XCCS` instruction, so the concatenation of executed instructions’ encodings is exactly equal to the concatenation of values returned by fetches. Thus, by Lemma 3, $\sigma_{b-1}.CCS$ is equal to the sum of values returned by fetches. These can be determined for all four categories of fetches:

- For $i \in N$, NoFAULT returns $[a_i]$;
- For $i \in S$, $S32(k)$ returns $[a_i + 4k_i]$;
- For $i \in R$ ($i > t$), S&R32 returns ρ_i , which is always equal to $[a_i - 4]$ (because $\rho_i = [a_i]$ at the end of every step that includes a fetch);

- If s_t uses the S&R32 rule, then the fetch for s_t returns ρ_t (the last value fetched by the previous block), otherwise it's counted by N and S .

Adding these up yields the claimed formula for $\sigma_{b-1}.CCS$. \square

Lemma 8 (Single faults always invalidate the checksum).

Let hb a hardened block of P and $s_t \dots s_b$ a legitimate execution of hb ending in a jump. Assume the maximum number of skips allowed in a single $S32(k)$ rule is $N = 1$. Then there cannot be exactly one fault attack during the execution of the block, i.e.

$$\text{Card} \{i \in [t, b] \mid s_i \text{ uses } S32(k) \text{ or } S\&R32\} \neq 1.$$

Proof. By Lemma 6, a legitimate exit by a jump requires an XCCS instruction at s_{b-2} , and s_{b-1} and s_b must either both use NoFAULT or both use S&R32. The second option immediately implies the theorem; this leaves the first.

Let $(N, S, R, \rho_{\text{entry}}, (a_i), (k_i), \text{next})$ be the fetch trace of $s_t \dots s_{b-2}$. Because s_{b-1} and s_b both use NoFAULT, the fault attack must occur in $s_t \dots s_{b-2}$, which leaves 3 options.

- s_t uses S&R32: by Lemma 7, the checksum at the end of the block is

$$\sigma_{b-1}.CCS = \sum_{i=t+1}^{b-2} [a_i] + \rho_{\text{entry}},$$

with $a_i = \text{blockAddr}(hb) + 4i$. Remember that the expected checksum is

$$\sigma_{\text{expected}} = \sum_{i=t}^{b-2} [\text{blockAddr}(hb) + 4i].$$

The difference is $\rho_{\text{entry}} - [\text{blockAddr}(hb)] = \rho_t - [a_t]$, which is non-zero due to the condition on S&R32 (silent replacements do not count as faults). Therefore the checksum doesn't pass, contradicting the hypothesis that $s_t \dots s_b$ is a legitimate execution of hb .

- Some s_j uses S&R32 ($j \neq t$): by Lemma 7, the checksum is

$$\sigma_{b-1}.CCS = \sum_{i=t}^{b-2} [a_i] - [a_j] + [a_j - 4]$$

still with $a_i = \text{blockAddr}(hb) + 4i$. The difference with σ_{expected} is $[a_j - 4] - [a_j] = \rho_j - [a_j]$ which is again non-zero due to the condition on S&R32.

- Some s_j uses $S32(1)$: still by Lemma 7, the checksum is now

$$\sigma_{b-1}.CCS = \sum_{i \in N} [a_i] + [a_j + 4] = \sum_{i=t+1}^{b-2} [a_i] - [a_j],$$

since the execution is offset by the skip at s_j , leading to $a_i = \text{blockAddr}(hb) + 4i + 4 (i \geq j)$. The difference is $[a_j]$; for the checksum to pass we must have $[a_j] = 0$. This is impossible: for the same reasons as discussed in Lemma 5, a 4-aligned zero value cannot intersect any instructions, so it would have to be a checksum value... but then s_j would crash. This is because the execution was not faulted up

to this point, so s_j must be using rule CHECKSUM-DELAY-SLOT with $S32(1)$. As a result, s_j fetching $[a_j + 4]$, which is not a valid checksum literal (it's a 32-bit jump), leads to a crash. \square

We are now finally able to prove the security property that a successful execution with no more than one fault per block has no faults at all.²¹

Theorem 2 (Security guarantee for single-fault executions).

Let P a fully hardened program and $e = [s_0, \dots, s_{|e|}]$ an execution such that

- s_0 is a legitimate entry into a block of P ;
- $s_{|e|}$ ends successfully, returning some $\text{end}(\alpha)$.

Let $[hb_1, \dots, hb_m], (s_{t_i} \dots s_{b_i})_{1 \leq i \leq m}$ be the partition of e into block executions outlined by Theorem 1.

If each segment $s_{t_i} \dots s_{b_i}$ uses at most one faulted fetch rule and the last segment $s_{t_m} \dots s_{b_m}$ is faultless, then e is faultless.

Proof. By Theorem 1, segments $i = 1$ to $m - 1$ all validate their checksums, and by hypothesis, they use at most one faulted fetch rule. By Lemma 8, they must in fact use no faulted fetch. The last block is itself faultless by hypothesis. As a result, the entire execution is legitimate. \square

B.4 Feasibility of Algorithm HARDEN

Algorithm 2 relies on the ability to get rid of invalid checksum literals by flipping their Least Significant Bit (LSB). We still have to show that this indeed works.

Lemma 9 (Invalid checksum literals can be avoided).

If $d : u_{32}$ is an invalid checksum literal, then $(d \oplus 1) + 2^{14}$ is a valid checksum literal.

Proof. Invalid checksum literals are jumps, Xccs instructions, and pairs of $c.ebreak$. The low 7 bits of each invalid literal is listed in Table ??; note how flipping the LSB never yields a new pattern that's present in the table.

It is important that the validity of the checksum literal can be inferred from the low 7 bits only. This is because flipping the LSB of the checksum literal also comes with an update to the Xccs opcode (replacing ccs by $ccsb$ or $ccscall$ by $ccscallb$). This update amounts to flipping bit 14 of the opcode, thus affecting bits 14–31 of the checksum value. Since bits 0–6 are not affected, the final checksum literal is still valid even after accounting for this change. \square

Received 11-NOV-2023; accepted 2023-12-23

²¹We exclude the last block since we don't protect $ecall$; in practice it's in non-protected $libc$ and the entry into the $\text{exit}()$ function is a function call that is itself protected.

