

Improving CPU Fault Injection Simulations

Insights from RTL to Instruction-Level Models

*Jasper van Woudenberg, Rajesh Velegalati,
Cees-Bart Breunese, Dennis Vermoen*

FDTC 2024

Problem statement

Programmers have a hard time developing FI resistant code

Can we help programmers create more resistant code through simulation?

Countermeasure refs:

https://riscureprodstorage.blob.core.windows.net/production/2017/08/Riscure_Whitepaper_Side_Channel_Patterns.pdf

<http://hardwarehack.ing>

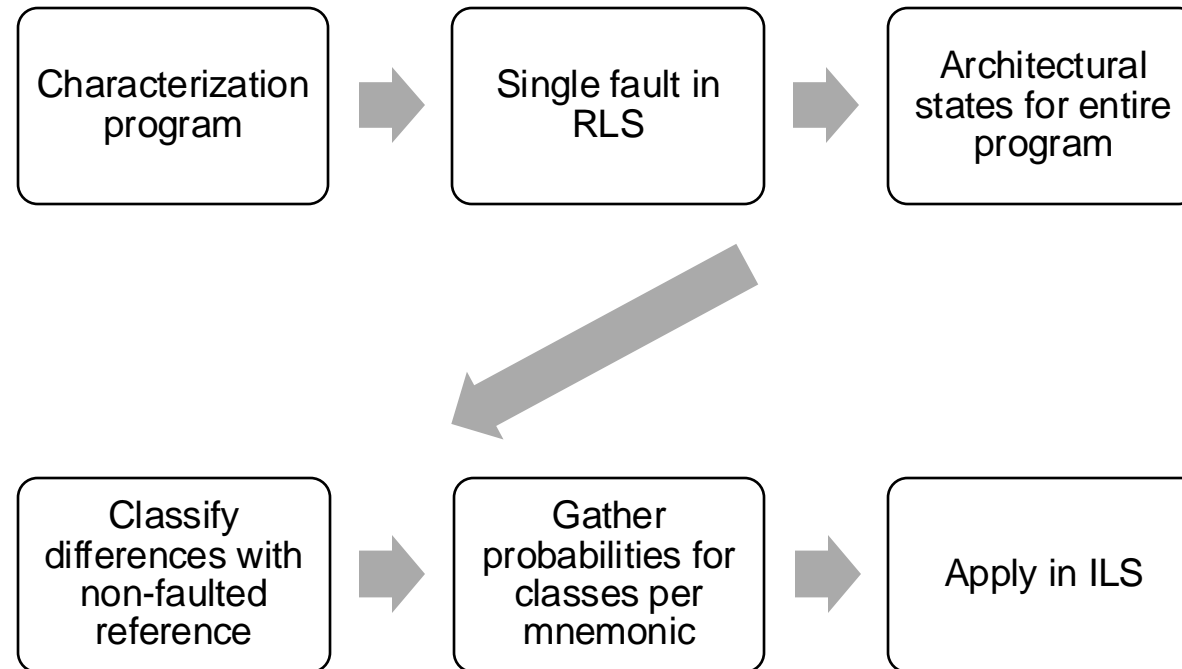
FI approach comparison

Rankings; higher is better

Can we make this “2”?

	RTL level sim (RLS)	Instruct. level sim (ILS)	Post-si fault injection
Accuracy	1	0	2
Runtime	0 (but scales with compute)	1 (also scales)	2 (scales poorly)
Ease of triage/fixing for sw dev	0	2	1
ROM fixing possibilities	1	2	0

Research overview



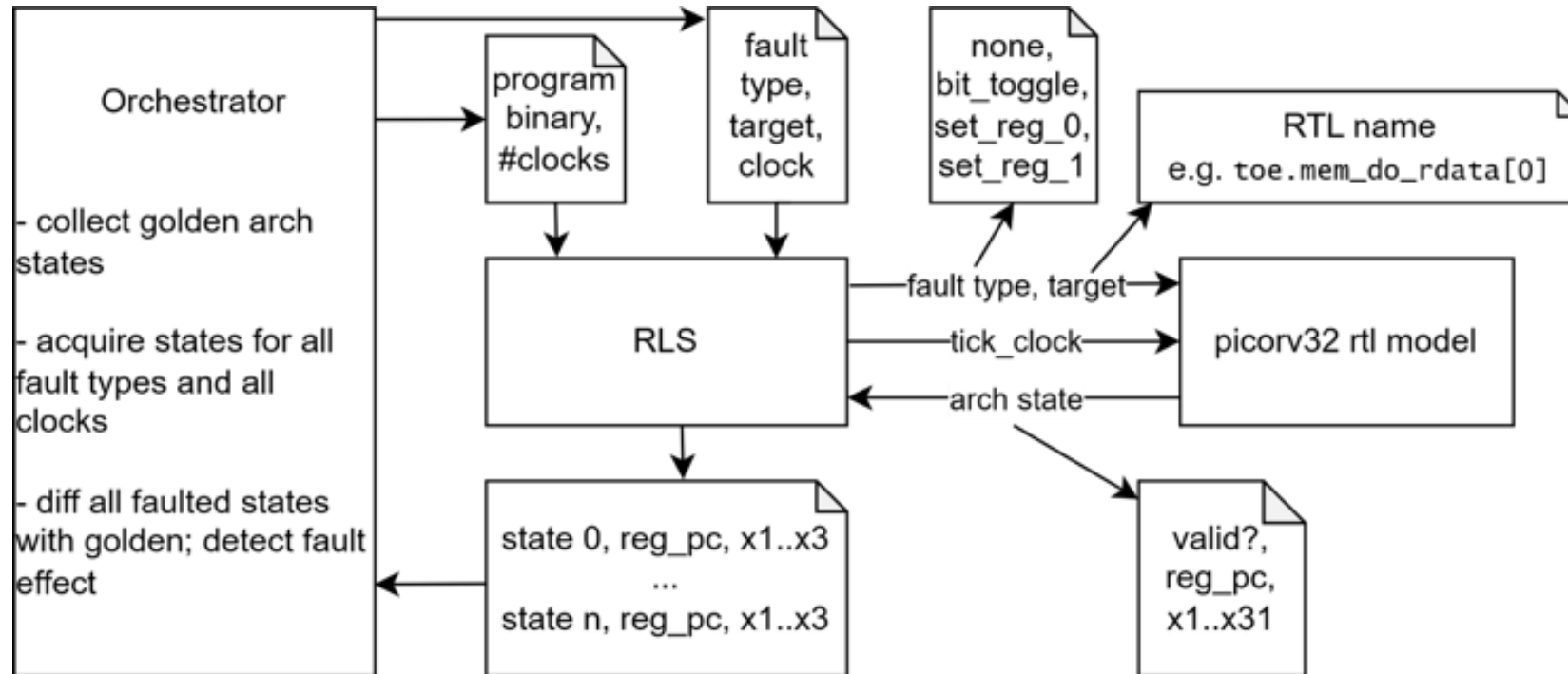
Characterization program

- We want fault statistics per mnemonic, presumably different statistics
- Put in (almost) all rv32imc instruction (mnemonics)
- Ensure each instruction changes the architectural state
 - Arch state we look at is registers + PC only

```
...
    sltu x13, x1, x2
        lui x1, 0x12345
        sra x14, x1, x2
        srl x15, x1, x2
        sub x16, x1, x2
        xor x17, x1, x2
        xori x18, x1, 0xff
        beq x1, x1, label_beq
            .word 0x00000013 // NOP in 32 b
label_beq:
    bge x1, x2, label_bge
        nop
label_bge:
    bgeu x1, x2, label_bgeu
        nop
label_bgeu:
    blt x17, x1, label_blt
        nop
label_blt:
    bltu x1, x2, label_bltu
        nop
label_bltu:
    bne x1, x2, label_bne
...

```

RTL level simulator

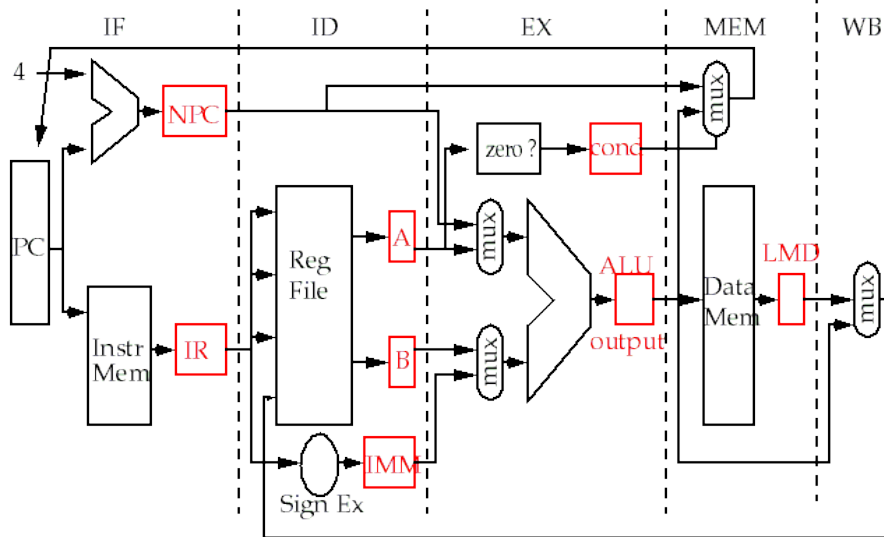


- Verilated Picorv32 core (configured for rv32imc), enhanced with FI capabilities
- Not faulting memories, only logic.
- Generates #clock cycles x #targets architectural state traces

Mapping faults to architectural states

- There is no 1:1 mapping between clock cycles and program counter (even in picorv32!)
- But, we can detect which state gets affected, and map back to that PC

http://ece-research.unm.edu/jimp/611/slides/chap3_1.html



Golden run

State	X1	PC
0	ab	00
1	cd	02
2	38	06

Faulted run

State	X1	PC
0	ab	00
1	cd	02
2	39	06

Fault differences

- For each RLS fault, we now know the instruction it was caused on and the arch state diff

```
},
"213": {
  "state": 37,
  "instr": "lw",
  "diff": [
    "reg_pc: 0000009a -> 0000009e",
    "reg_r5: 00000000 -> 00010001"
  ],
  "fi_diff": {
    {'reg_pc': '0x98'}: [
      "b'Single bit-flip in true_code_interface.toe.reg_pc[1]'"
    ],
    {'reg_pc': '0x9e'}: [
      "b'Single bit-flip in true_code_interface.toe.reg_pc[2]'"
    ],
    {'reg_pc': '0x92'}: [
```


Classify the fault differences

- Characterize the effect into 6 classes
- Aggregate counts per mnemonic
- The resulting model we call Architectural Fault Effect Model (AFEM)
 - We ignore “unknown”

```
"nop": {  
  "unknown": 29,  
  "nop": 0,  
  "reg_flip": 42069,  
  "reg_zero": 829,  
  "reg_one": 1279,  
  "reg_swap": 0  
},  
"addi": {  
  "unknown": 0,  
  "nop": 14,  
  "reg_flip": 23707,  
  "reg_zero": 671,  
  "reg_one": 719
```

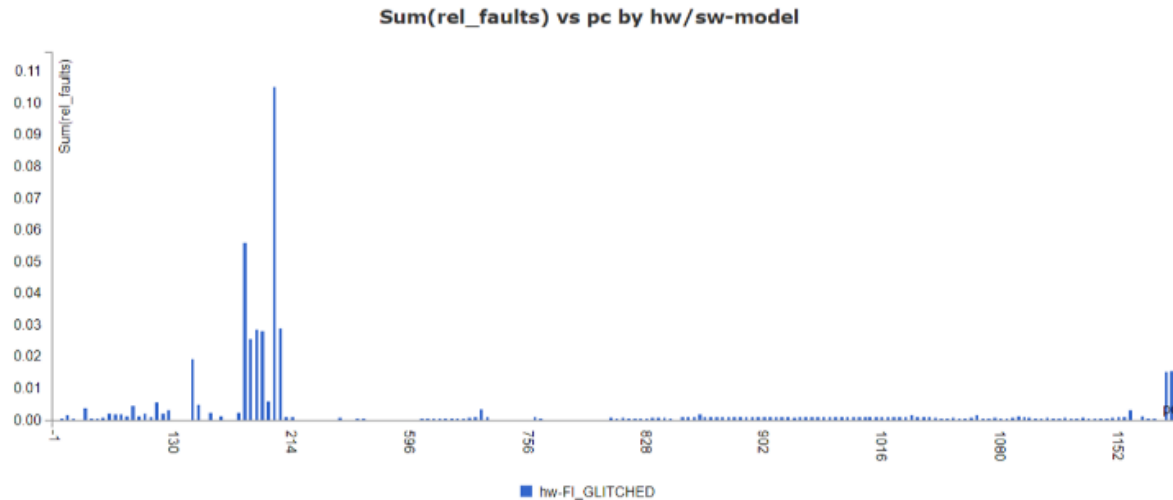
Instruction Level Simulator using AFEM

- Obtain an architectural state to fault
- Flip weighted coin to choose a class
- Further flip coin where applicable to select a register and/or bits therein

```
"nop": {
  "unknown": 29,
  "nop": 0,
  "reg_flip": 42069,
  "reg_zero": 829,
  "reg_one": 1279,
  "reg_swap": 0
},
"addi": {
  "unknown": 0,
  "nop": 14,
  "reg_flip": 23707,
  "reg_zero": 671,
  "reg_one": 719
```

Evaluation for accuracy and runtime

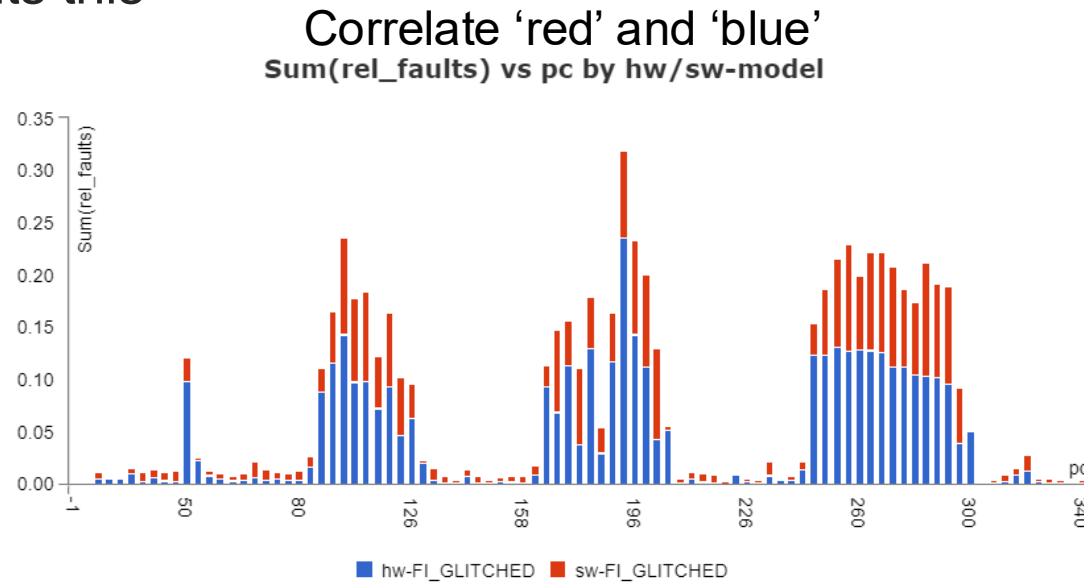
- Compare ILS AFEM \$x samples to RLS, NOP, NOP2, instruction bitflip
 - \$x = how often to do a program run with fault per instruction
- Test program from FIRM*: branch test, memory test, register test
- Correlation between vectors PC -> successful fault probability; RLS baseline



* FI Resistance Metric developed at Riscure by Carlo Maragno, Praveen Vadnala, Pierre-Yves Peneau, Chris Berg, Nisrine Jafri, Marc Witteman

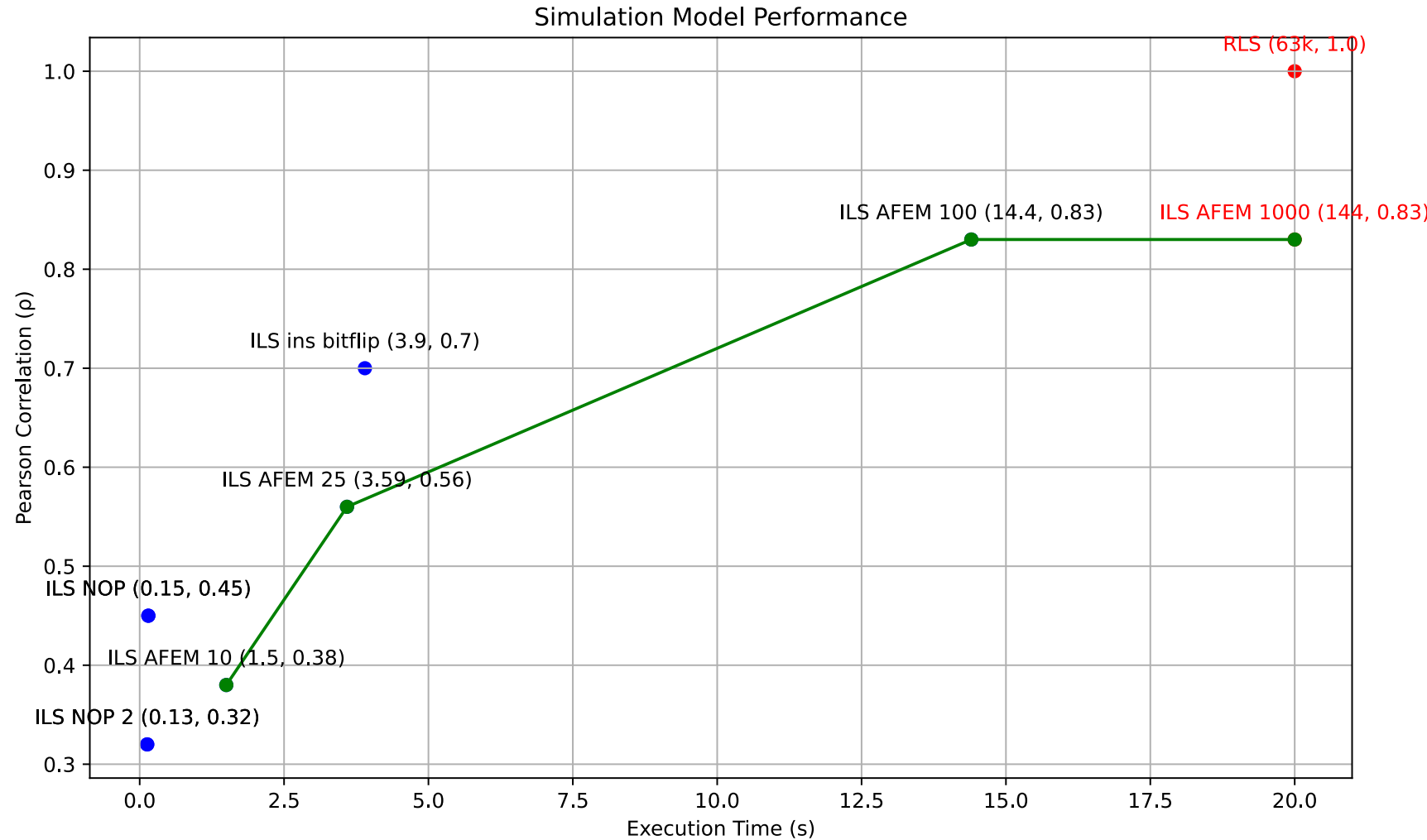
Correlation rationale

- We don't care about absolute height of peaks (hw prob << sw prob)
- We do care about relative heights (sensitive areas in program)
- Per PC, prediction on hw prob given sw prob
- Pearson correlation fits this



Accuracy vs time tradeoffs

Simulation model	ρ	sim time (s)
RLS	1	63k
ILS AFEM 1000	0.83	144
ILS AFEM 100	0.83	14.4
ILS ins bitflip	0.70	3.9
ILS AFEM 25	0.56	3.59
ILS NOP	0.45	0.15
ILS AFEM 10	0.38	1.5
ILS NOP 2	0.32	0.13

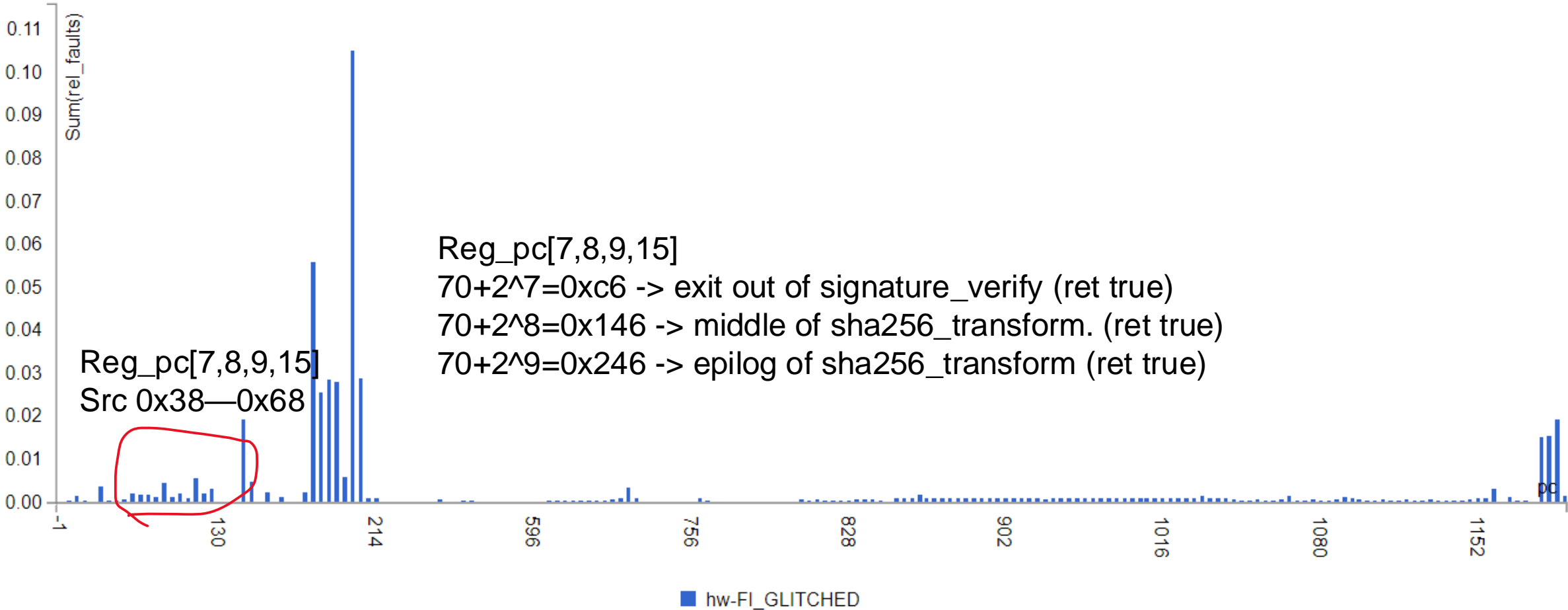


Application to boot code

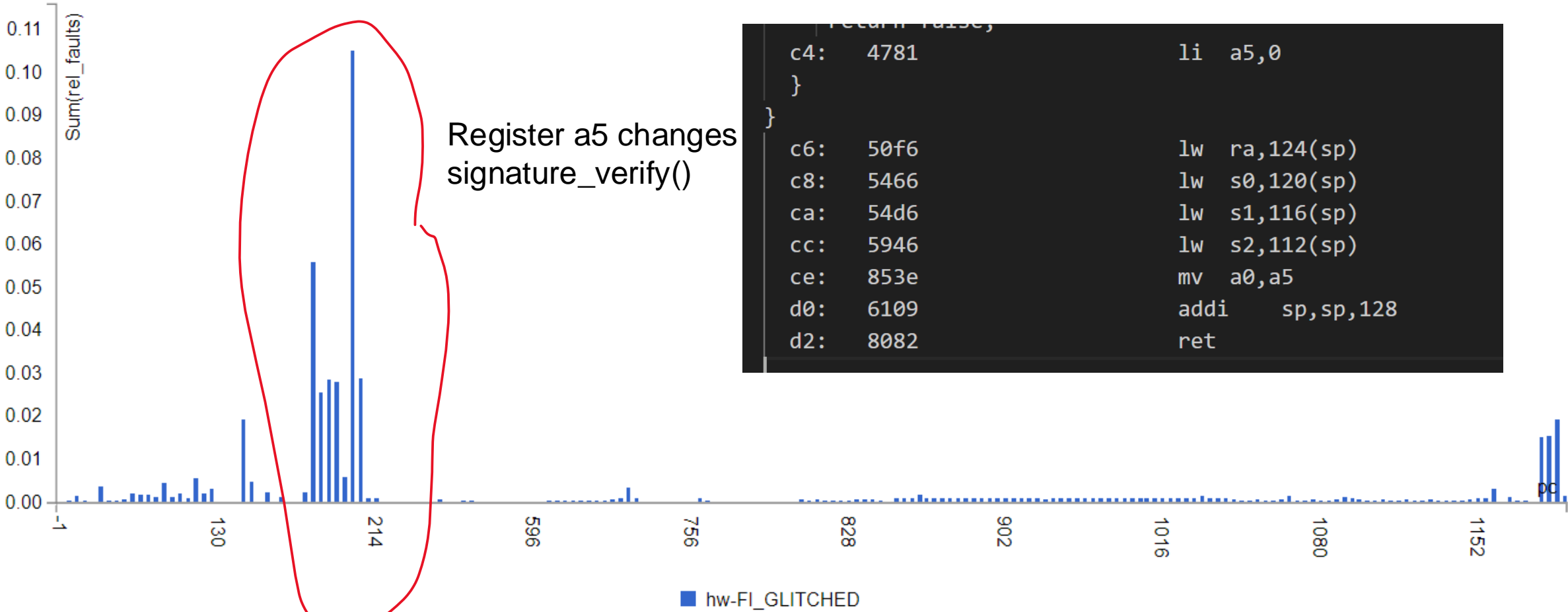
- Test code is test code and may not represent real applications
- Try boot code (... also test code, admittedly)

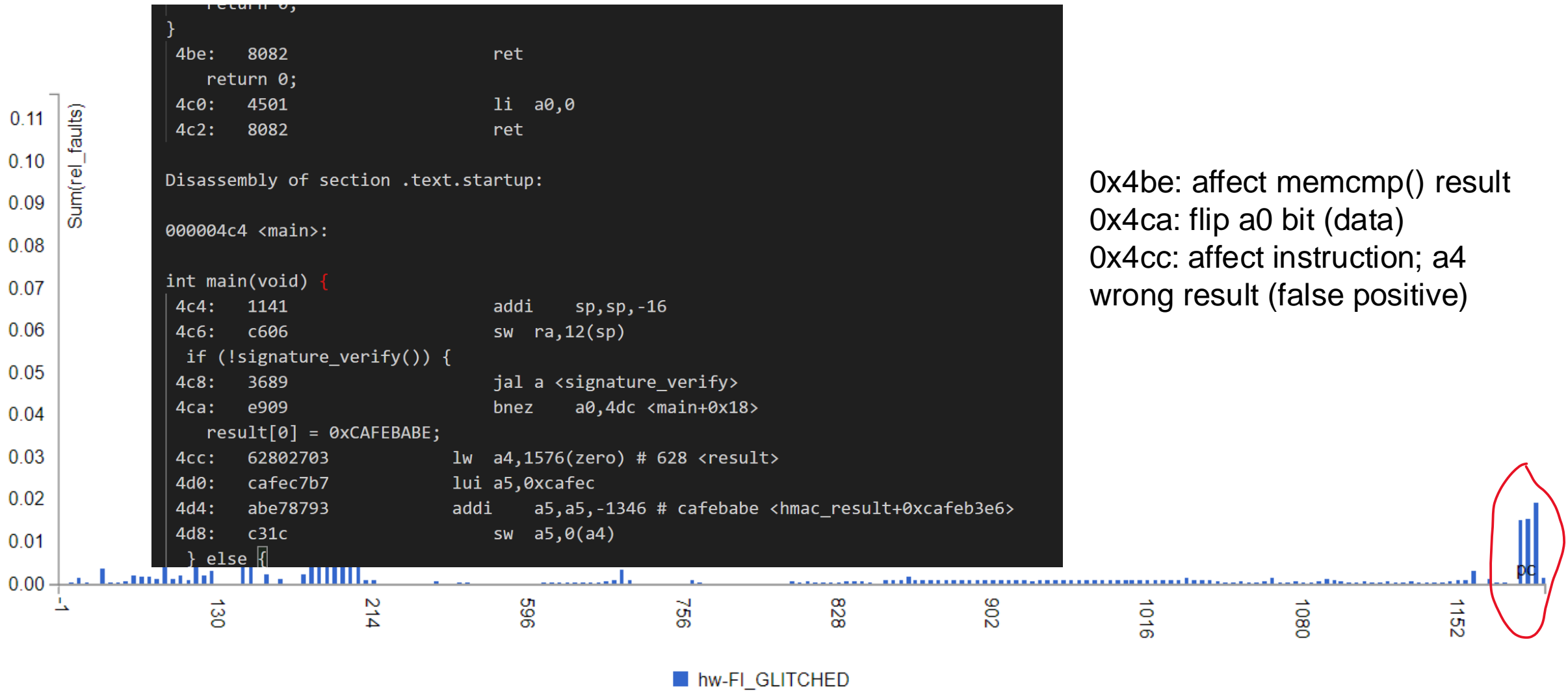
```
int main(void) {
    if (!signature_verify()) {
        result[0] = 0xCAFEBAFE;
    } else {
        application_entry_point();
    }
    // Do not return from this main function!
    while(1) {}
}
```

Sum(rel_faults) vs pc by hw/sw-model

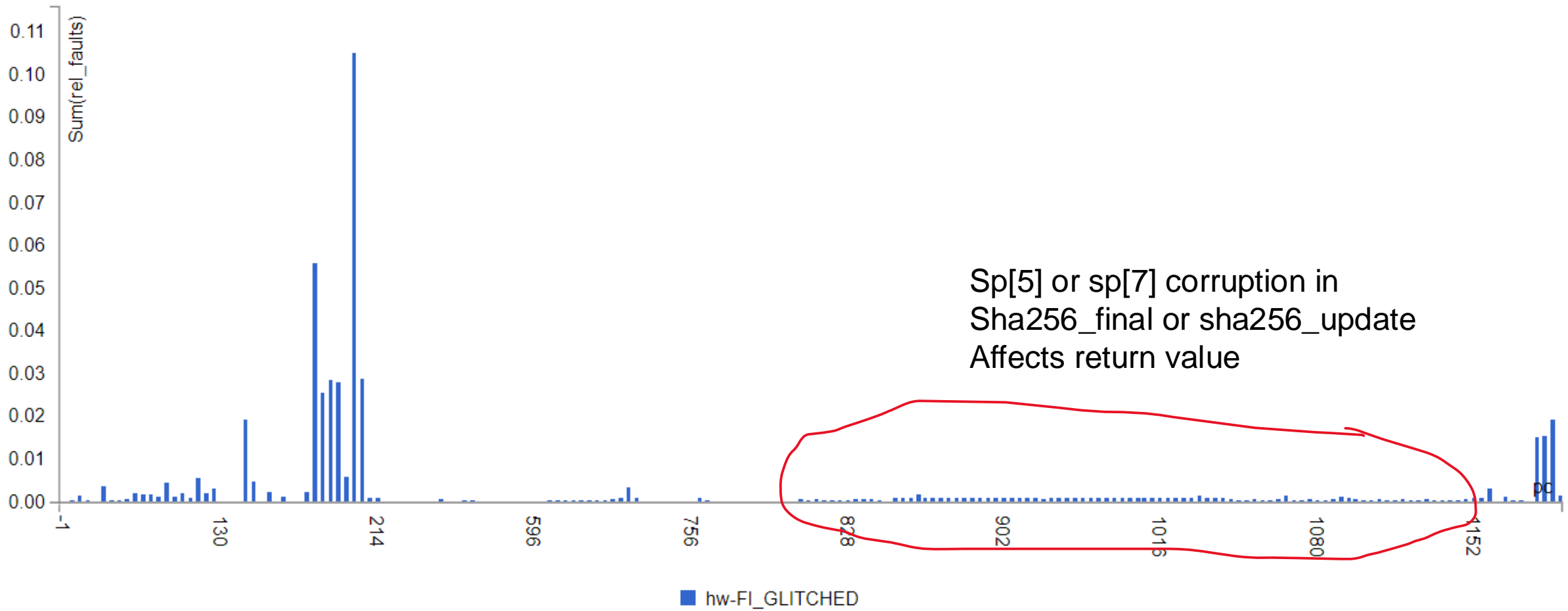


Sum(rel_faults) vs pc by hw/sw-model





Sum(rel_faults) vs pc by hw/sw-model

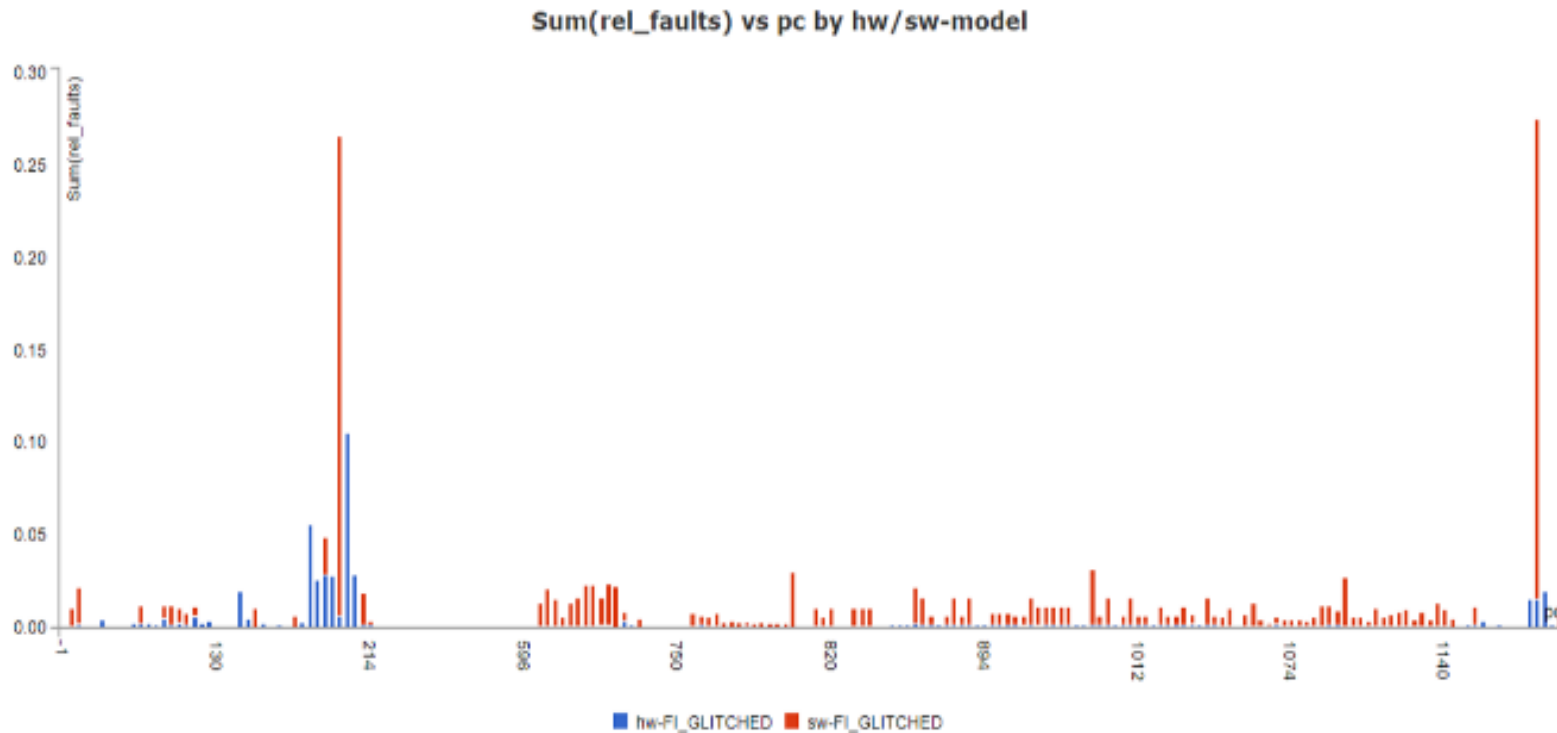


Improved AFEM: 'regonly'

- Insight: most faults directly affect the register state
 - Processor registers are a significant chunk of the RTL registers
 - Most relevant: src register, the PC, return address (RA), or Stack Pointer (SP).
- Try next model: randomly select one of these CPU registers and flip a random bit within the chosen register.
 - No longer conditional probability on the instruction mnemonic

Regonly results

- Identifies very similar regions to the RLS
- $\rho = 0.06$; 'regonly' predicts vulnerable locations, not accurately assigns probability



Discussion

- Generalization beyond picorv32 / rv32imc to larger designs
- Simple programs used; most bit flips cause relevant faults
- Did not test programs / CPUs with countermeasures
- RLS as baseline, instead of post-si testing
- Simple input fault models (no complex multi-bit faults)
- Larger designs / software require significant #CPUs

Conclusions

- Predicting relevant code paths for FI is non-trivial, simulation/testing can help
- We can learn FI models from arbitrary CPU implementations
 - Faster than RLS
 - More accurate than existing instruction skip/bitflip models
- For our CPU, we can improve by manually finetuning the model

Future work

Future:

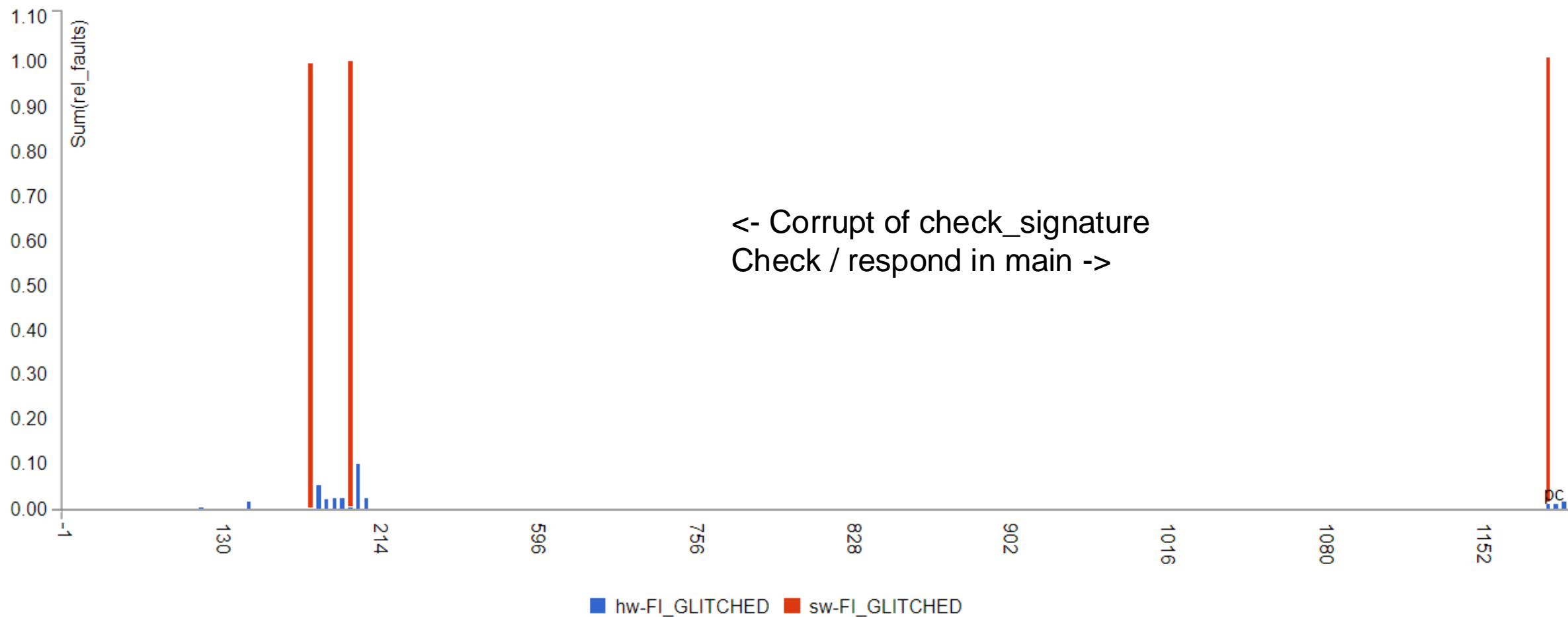
- Metrics that reflect indicating relevant code parts (as opposed to predicting the fault effects)
- Post-si evaluation of the models accuracy in sensitive code
- Validate how generic the models are

We would like to thank Google for their funding of this work

Thank you

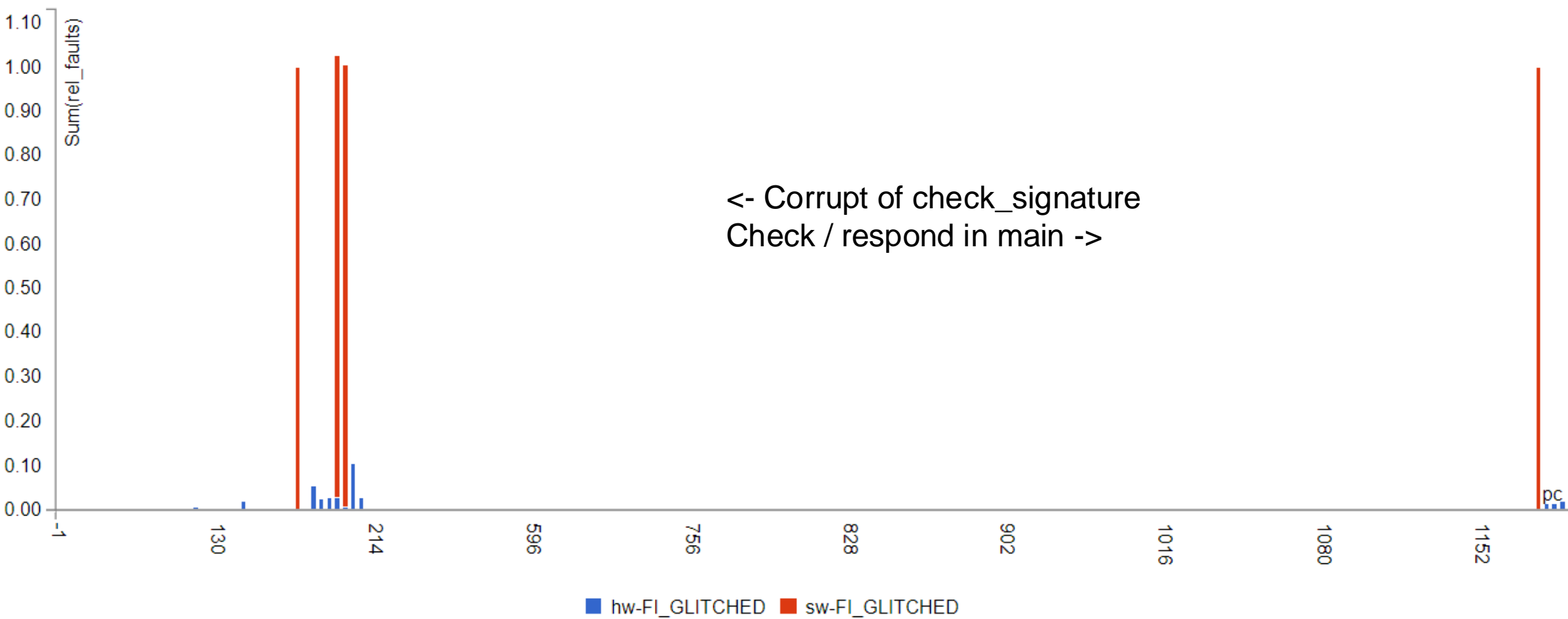
NOP

Sum(rel_faults) vs pc by hw/sw-model



nop2

Sum(rel_faults) vs pc by hw/sw-model



Metric

Property	Desire	Metric
Relevance	rank most sensitive parts of code	correlation hw and sw
Actionable	highlight lines of code Indicate fault model	no metric needed, straight output from sw simulator
Performance	run sim in reasonable time	instructions/sec sec/full test