

A Red Team/Blue Team Assessment of Functional Analysis Methods for Malicious Circuit Identification

Adam Waksman

Jeyavijayan Rajendran

Matthew Suozzo

Simha Sethumadhavan

Computer Architecture and Security Technologies Lab
Department of Computer Science
Columbia University
New York, NY, USA
{waksman,simha}@cs.columbia.edu
{ms4249@columbia.edu}

Department of Electrical and Computer Engineering
New York University
New York, NY, USA
{v.ece}@nyu.edu

ABSTRACT

Recent advances in hardware security have led to the development of FANCI (Functional Analysis for Nearly-Unused Circuit Identification), an analysis tool that identifies stealthy, malicious circuits within hardware designs that can perform malicious backdoor behavior. Evaluations of such tools against benchmarks and academic attacks are not always equivalent to the dynamic attack scenarios that can arise in the real world. For this reason, we apply a red team/blue team approach to stress-test FANCI's abilities to efficiently detect malicious backdoor circuits within hardware designs.

In the Embedded Systems Challenge (ESC) 2013, teams from research groups from multiple continents created designs with malicious backdoors hidden in them as part of a red team effort to circumvent FANCI. Notably, these backdoors were not placed into *a priori* known designs. The red team was allowed to create arbitrary, unspecified designs. Two interesting results came out of this effort. The first was that FANCI was surprisingly resilient to this wide variety of attacks and was not circumvented by any of the stealthy backdoors created by the red teams. The second result is that frequent-action backdoors, which are backdoors that are not made stealthy, were often successful. These results emphasize the importance of combining FANCI with a reasonable degree of validation testing. The blue team efforts also exposed some aspects of the FANCI prototype that make analysis time-consuming in some cases, which motivates further development of the prototype in the future.

Categories and Subject Descriptors

B.6.2 [Hardware]: Logic Design—*Security and Trust*

Keywords

hardware; security; backdoors; functional analysis; intellectual property

1. INTRODUCTION

Hardware security and trust is a subject of rapidly increasing global concern [1, 2, 3, 4]. The economic drive for newer and better computing technology demands global cooperation and the sharing of intellectual property. However, as third-party IP is increasingly used by technology companies, major trust issues exist [5, 6, 7]. A single malicious circuit, often called a backdoor or trojan, hidden in hardware design IP, can have catastrophic effects [8, 9].

Recently, static analysis was proposed as a method for combating third-party IP backdoors. A tool called FANCI (Functional Analysis for Nearly-Unused Circuit Identification) was developed that specifically targets a large class of backdoors in digital designs, called *stealthy* backdoors [10].

FANCI has performed extremely well on academic designs and benchmarks, but the current set of benchmarks is limited. As further evaluation, we present in this paper a red team/blue team approach to stress-testing FANCI. In this approach, a variety of red teams design different backdoors, to stress FANCI's abilities, both for the types of attacks it was designed to stop and for alternative types of attacks. Our blue team applied the FANCI tool to the red teams' designs and applied minimal manual analysis (about one hour per design) to attempt to track down the backdoor in each design.

The goal of this type of red team/blue team approach is to both stress-test the tool to see if it achieves everything the implemented algorithm should achieve and to violate different axioms of the system to see how the tool responds to new types of attacks. This was all performed as part of the 2014 Embedded Systems Challenge (ESC). We discuss the results and our observations, as well as comments on future directions for functional analysis-based security.

2. OVERVIEW OF BLUE TEAM METHODS

Our blue team used the recently proposed FANCI [10]. This tool is a prototype, semi-complete version of a new algorithm for detecting any and all stealthy logic with a digital hardware design or gatelist. The idea behind FANCI is the following. An organization wants to acquire a third-party hardware design, either as source code or a gatelist. Some degree of validation and/or verification will be applied to the produced design to make sure that it matches or is at least similar to that specification. A malicious provider may include backdoors in this third-party design, but by necessity the design is likely to be similar to the true specification. Thus, the added backdoor circuitry is likely to be *stealthy*, *i.e.*, largely unused. It may not be the case that all stealthy circuits are malicious, but it is usually the case that all malicious circuits are stealthy (this

Algorithm 1 Compute Control Value

```
1:  $count \leftarrow 0$ 
2:  $c \leftarrow \text{Column}(w_1)$ 
3:  $T \leftarrow \text{TruthTable}(w_2)$ 
4: for all Rows  $r$  in  $T$  do
5:    $x_0 \leftarrow \text{Value of } w_2 \text{ for } c = 0$ 
6:    $x_1 \leftarrow \text{Value of } w_2 \text{ for } c = 1$ 
7:   if  $x_0 \neq x_1$  then
8:      $count++$ 
9:   end if
10: end for
11:  $result \leftarrow \frac{count}{size(T)}$ 
```

observation is supported by analysis of hardware benchmarks) [10]. Thus, FANCI has been designed as a tool that identifies all stealthy circuits within a design. Once this has been done, we have a set of stealthy circuits which should be a superset of the set of malicious circuits. There can be false positives within that set, but if the tool works correctly there should not be false negatives. This is the principle idea behind FANCI.

FANCI works by detecting stealthy circuits by finding circuits that do not often impact other wires they are connected to. We call these wires *weakly-affecting*, because they drive values into other wires but do not often change the digital results or outputs of those wires. We quantitatively measure the degree of impact one wire has on other using a novel metric called *control value*. The control value of an input w_1 on an output w_2 quantifies how much the truth table representing the computation of w_2 is influenced by the input column corresponding to w_1 . Specifically, control value is a number between zero and one quantifying what fraction of the rows in the truth table for a circuit are directly influenced by w_1 . Importantly, this is independent of particular tests inputs that might be supplied during validation. The algorithm to compute the control value of w_1 on w_2 is presented as Algorithm 1. We note that in step 3, we do not actually construct the exponentially large truth table. We instead construct the corresponding boolean function. Since the sizes of truth tables grow exponentially, we approximate control values in practice by only evaluating a constant-sized subset of the rows in the truth table.

Once we have computed all of the control values for a given wire (a wire here means an output of any internal gate), we have a vector of floating point values that we can combine to make a judgement about stealth. We have found that using simple aggregating heuristics are effective for identifying stealthy wires. For ESC, we focus on two metrics. The first is simply the arithmetic mean, which is helpful usually for detecting wires that are part of the trigger. The other metric is called triviality and tends to be helpful in detecting payload wires. Triviality is described in full in [10] but can be thought of simply as the measure of how often the output is equal to zero (or by symmetry one). For example, a circuit that always outputs zero is completely trivial, while a circuit that outputs zero in half of cases is more benign. Other aggregation metrics may be able to achieve even better results in the future. Our overall algorithm is summarized in Algorithm 2.

For an example of a backdoor circuit and how FANCI detects it, consider a classic ‘ticking timebomb’ trigger. A backdoor is designed to go off after 2^{40} cycles of operation. This is implemented with a 40-bit counter and a comparator against the hard-coded value $2^{40} - 1$. The triviality of the output of the comparator would be $\frac{1}{2^{40}}$. When looking at the payload circuit, it would have all 40 input wires from the counter as inputs, meaning its vector of control

Algorithm 2 How FANCI Flag Suspicious Wires in a Design

```
1: for all modules  $m$  do
2:   for all gates  $g$  in  $m$  do
3:     for all output wires  $w$  of  $g$  do
4:        $T \leftarrow \text{TruthTable}(\text{FanInTree}(w))$ 
5:        $V \leftarrow \text{Empty vector of control values}$ 
6:       for all columns  $c$  in  $T$  do
7:         Compute control of  $c$ 
8:         Add control( $c$ ) to vector  $V$ 
9:       end for
10:      Compute heuristics for  $V$ 
11:      Denote  $w$  as suspicious or not suspicious
12:    end for
13:  end for
14: end for
```

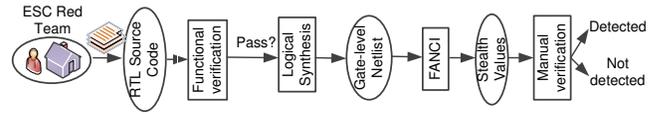


Figure 1: The steps involved in the contest, including both red team and blue team actions.

values would contain 40 low values, making the median value $\frac{1}{2^{40}}$. In this case, both the trigger and payload circuits are caught due to the stealthy nature of the triggering logic.

3. RED TEAM METHODS AND COMPETITION RULES

We outline the rules of the contest. Since the goal of the contest was to attack FANCI at its axioms, in the end we allowed most of the original rules to be broken. The main rule that we did not allow to be broken was that all attacks had to be digitally defined in the design. This means that we did not process any analog logic or allow for backdoors that require environmental or physical factors to turn on. We also only processed Verilog code, although adapting the FANCI tool for VHDL or Bluespec would be possible. Each submitted design contained exactly one backdoor. The blue team was given roughly three days to process and analyze the submissions. The flow of the contest is depicted in Figure 1.

The original rules were that a submission should include a single module with 10,000 or less gates, including both source code and gatelist, using only standard logical gates and flip flops, as well as documentation of the module’s specification. We ended up making a variety of exceptions, including 1) designs larger than 10,000 gates, 2) designs composed of many modules, 3) designs lacking source code, 4) designs lacking specification, 5) backdoors that rely on non-standard clock/reset usage, and 6) designs that use non-standard flip flops.

One aspect of this contest that was unique was that red teams were allowed to choose or create their own baseline designs. They were not forced to insert backdoors into previously existing designs. This added an interesting dimension to the blue team analysis, because some of the designs being analyzed had not been seen before, and some of those were also without source code. The red teams did not intentionally obfuscate their gatelists, which is a direction that might be interesting in the future. Code obfuscation should not impact FANCI, since the tool ignores wire names, but it would add difficulty to the manual analysis of the results that ultimately must happen.

Red teams were given points for backdoors that went undetected. They received more points for more stealthy backdoors, ranging from one point (for always or frequently on backdoors) to five points (for backdoors stealthy enough to go unnoticed even by high quality validation testing). The winners of the contest were determined by a panel of judges based on the score and how interesting the design choices were.

4. RESULTS AND EFFICACY OF BLUE TEAM ANALYSIS

The results of the blue team analysis were overall better than expected. While many of FANCI’s axioms were targeted by attacks, FANCI still held up quite well. All highly stealthy attacks were detected with relative ease. We found overall that the best way to go after FANCI in this contest was to either use a frequently-on (or always-on) backdoor or to use a ‘pathological pipeline.’¹ These are the two cases where a backdoor is likely to bypass FANCI. In the first case, frequently-on backdoors are not stealthy. Hence, FANCI will not flag them as such. In the second case, extreme pipelining makes it harder to understand the data and stress-tests those portions of the tool prototype that are least developed. We include further details for the three red teams that were deemed the best by the judges.

- Red Team A:** Red Team A created three backdoors, which highlight different aspects of the FANCI tool and algorithm. They created two frequent-action backdoors (which were not detected) and one stealthy backdoor (which was detected). The stealthy backdoor uses a 66-bit cheat code that comes from a pseudo-random input source. This backdoor highlights the importance of using intermediate wires, which was one of the main contributions of FANCI. If FANCI only computed the impact of module inputs on module outputs, this type of backdoor might not be caught, but since we analyze all intermediate wires, noticing a 66-input AND gate is easy. The other two backdoors were frequent-action backdoors that are nearly-always on. One is turned on by the *reset* signal, and one is turned on only during the negative clock edge. Both of these are frequently enough on that they have no low control values and are not stealthy (and thus undetectable). These are types of backdoors that highlight the need to still use minimal validation testing in addition to FANCI. They also illuminate the fact that signals like reset and the clock, which might be given less attention during validation testing, are still relevant, especially when dealing with an adversary (as opposed to simply looking for unintentional bugs).

- Red Team B:** Red Team B submitted three backdoors, one that we evaluated for the contest and two that we looked at post-competition due to logistical constraints. This team produced a variety of backdoors on independent baseline designs. The first backdoor was designed into a DES accelerator and was built off of an aggregation of finite state machines (FSMs), somewhat similar to UCONN’s approach². They used the state machines to control round scheduling and thus weaken the strength of encryption, which would expose the ciphertext to software-based attacks upon release. This was a case where we were able to detect the trigger circuit because of the combinational logic that controls the state machines, though we note that this backdoor took more manual analysis to track down

¹Pathological pipeline is the term used in the original FANCI paper [10] to describe backdoor triggers that rely on heavy pipelining to prevent significant stretches of combinational logic from existing. These backdoors may also be referred to as FSM-heavy.

²UCONN did not win the contest but did construct an FSM-based, heavily pipelined backdoor that took a long time to find. The various teams worked independently and simultaneously.

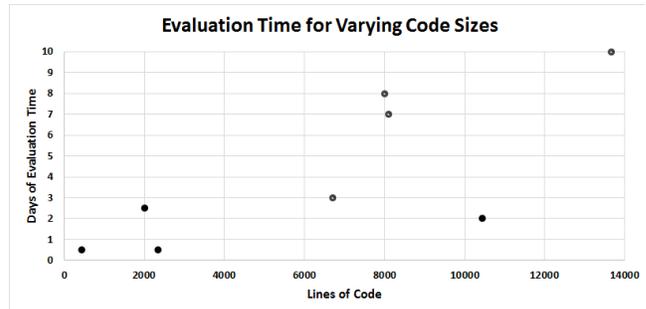


Figure 2: Trade-off between lines of code analyzed and runtime. The black points denote designs that finished completely in the contest time frame. The lighter points represent designs that were analyzed partially for the contest. For those points, the full runtime is estimated.

than any of the other ones we detected.

The second backdoor was a similar round-scheduling attack, this time on XTEA. This one acquires its stealth by counting the number of one-bits in the input, which is a good way to generate high stealth. However, the combinational circuit for doing the counting was flagged. The third backdoor also attacked crypto hardware, this time a cellular automata PRNG. This one was a sequence cheat code where the user supplies one bit of the trigger at a time. This attack would be extra hard for validation to catch, because it hides on the negative edge of the clock (similar to one of the Red Team A attacks). However, since FANCI ignores clock timing, it can detect the trigger logic. Making this attack more frequently on might be more effective against both FANCI and validation tests if the validation team is known to not check negative clock edges.

- Red Team C:** Red Team C produced a backdoor in a stream cipher module. This attack is essentially one of the ones suggested in the original FANCI paper [10] and highlights the need for basic validation and oversight when applying FANCI. The attack is a combination of a frequent-action backdoor and a pathological pipeline backdoor. The trigger fires frequently but not always, meaning the stealth scores are not particularly low. The design is also heavily pipelined, with roughly one flip flop between each pair of complex logic gates for critical portions of the design. Looking at the gatelist, it is immediately obvious that the design has been compromised, but identifying what the exact backdoor payload is can be quite hard, and FANCI does not detect this type of attack. As mentioned in the original paper, going after this type of attack requires either validation tests or oversight from an integration engineer.

5. OBSERVATIONS ON RESULTS

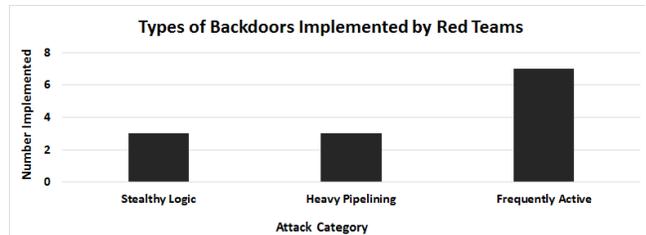


Figure 3: The types of backdoors implemented by the red teams.

We include a few observations and takeaways based on the results of the contest and our experiences.

- **Runtime and Scalability** Figure 2 shows the runtime of the tool as a function of the number of lines of code in the various designs, using one primary design from each red team. Naturally, FANCI runs slower on larger designs, but the slowdown is more or less linear, which makes the analysis feasible. These tests were done on a single core of a commodity machine.

- **Attack Categorization:** A positive result of the contest was the discovery that many of the red teams designed attacks very similar to the types we anticipated when first designing FANCI. In [10], we mentioned three general attack avenues against FANCI: frequently active (non-stealthy) backdoors, heavily pipelined backdoors, and false positive flooding. While the third option was not employed by the red teams, the first two were used by multiple teams. The breakdown is shown in Figure 3. This evidence supports that a more formal taxonomy could be derived from this survey. Additionally, it supports our belief that FANCI and validation testing should be used together synergistically. Ideally, validation testing should be designed with the assumption that FANCI will detect anything stealthy. This would allow validation teams to focus their efforts on other avenues, such as some of the attacks we saw that target reset or the negative clock edge. Figure 4 shows the overall breakdown by red team and also divides the frequently active backdoors from the stealthy backdoors. While FANCI caught all the stealthy backdoors, it caught only a few of the frequently active ones.

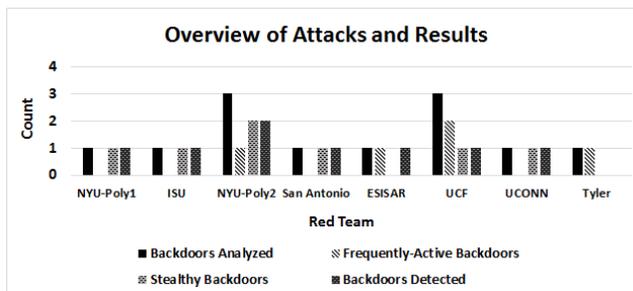


Figure 4: Overview of red team attacks and blue team results.

- **Algorithm vs. Implementation:** While the contest did not expose any deficiencies in the FANCI algorithm, the tool itself was stressed in some cases. Two issues stand out. First, runtime became an issue for large designs. Some modules would have taken more than the given three days to analyze, and so incomplete analysis runs were done. The tool is configurable for this, allowing for hasty passes. However, in the future, parallelization could do a much better job of alleviating this problem. The second issue is the way the tool handles pipelining. The core of the tool works on combinational logic, so flip-flops have to be dealt with. We believe the best way to handle flip-flops is to treat them as identity gates, so that simply inserting dummy flip-flops does not hide stealthy logic. On the other hand, this creates loops in the logic, which have to be dealt with. For most cases, our tool currently treats flip-flops as a barrier and does not analyze past them. This did not prevent us from catching any stealthy backdoors in this contest, but it made manual analysis more difficult. Improving the tool for this case would be beneficial.

- **Primary Takeaways:** The primary takeaways from the contest appear to be that FANCI handles a wide variety of backdoors and that effort should be spent on improving the usability of the proto-

type and on making static analysis and validation testing a tandem for future designs.

6. CONCLUSIONS

The ability to identify and understand hardware backdoors during the design phase using static analysis is critical for allowing continued use of third-party hardware intellectual property. Using a red team/blue team approach, we stress-tested FANCI, a state-of-the-art tool for identifying stealthy backdoors. Using this approach, we saw examples of stealthy attacks designed to target FANCI specifically, as well as examples of non-stealthy backdoors. Overall, the results of the contest were promising, as they demonstrated the effectiveness and flexibility of the FANCI approach. However, the results continue to highlight the importance of security awareness when integrating hardware designs, including FANCI, validation testing, and reasonable oversight.

References

- [1] Sally Adee. The Hunt for the Kill Switch. *IEEE Spectrum Magazine*, 45(5):34–39, 2008.
- [2] Marianne Swanson, Nadya Bartol, and Rama Moorthy. Piloting Supply Chain Risk Management Practices for Federal Information Systems. In *National Institute of Standards and Technology*, page 1, 2010.
- [3] Adam Waksman and Simha Sethumadhavan. Silencing Hardware Backdoors. In *Proceedings of the 2011 IEEE Symposium on Security and Privacy*, pages 49–63, Oakland, California, 2011.
- [4] Hassan Salmani, Mohammad Tehranipoor, and Jim Plusquellic. New Design Strategy for Improving Hardware Trojan Detection and Reducing Trojan Activation Time. In *Hardware-Oriented Security and Trust, 2009. HOST '09. IEEE International Workshop on*, pages 66–73, 2009.
- [5] M Tehranipoor, H. Salmani, Xuehui Zhang, Xiaoxiao Wang, R. Karri, J. Rajendran, and K Rosenfeld. Trustworthy Hardware: Trojan Detection and Design-for-Trust Challenges. *Computer*, 44(7):66–74, 2011.
- [6] Adam Waksman, Julianna Eum, and Simha Sethumadhavan. Practical, Lightweight Secure Inclusion of Third-Party Intellectual Property. In *Design and Test, IEEE*, pages 8–16, 2013.
- [7] E. Love, Yier Jin, and Y Makris. Proof-Carrying Hardware Intellectual Property: A Pathway to Trusted Module Acquisition. *Information Forensics and Security, IEEE Transactions on*, 7(1):25–40, 2012.
- [8] Samuel T. King, Joseph Tucek, Anthony Cozzie, Chris Grier, Weihang Jiang, and Yuanyuan Zhou. Designing and Implementing Malicious Hardware. In *Proceedings of the 1st Usenix Workshop on Large-Scale Exploits and Emergent Threats*, pages 5:1–5:8, Berkeley, CA, USA, 2008. USENIX Association.
- [9] R. Karri, J. Rajendran, K Rosenfeld, and M Tehranipoor. Trustworthy Hardware: Identifying and Classifying Hardware Trojans. *Computer*, 43(10):39–46, 2010.
- [10] Adam Waksman, Matthew Suozzo, and Simha Sethumadhavan. FANCI: Identification of Stealthy Malicious Logic Using Boolean Functional Analysis. In *ACM Conference on Computer and Communications Security*, pages 697–708, 2013.