# Hardware Trojan Detection using Unsupervised Learning and Simulation Based Side-Channel Features

Authors - Om Tank , Aayan Boradia and Hemang Gandhi

## Abstract

This project presents a novel approach to hardware Trojan detection using unsupervised learning techniques applied to simulation-based side-channel features. Hardware Trojans, malicious modifications to integrated circuits, pose significant security threats to critical infrastructure, military systems, and financial networks. Traditional detection methods typically require golden reference models, which are often unavailable in real-world supply chain scenarios.

Our methodology eliminates this dependency by extracting behavioral features from circuit simulations using Quartus Prime and ModelSim Altera tools. We analyze toggle counts (signal transitions), timing patterns, and Hamming weight differences (number of '1' bits in signals) as simulation-based proxies for physical side-channel information. These features are then processed using Principal Component Analysis (PCA) for dimensionality reduction and unsupervised clustering algorithms to distinguish between clean and Trojan-infected circuits.

We demonstrated this approach using various circuit implementations, including simple counters, combinational logic, and adder chains. Our simulations revealed that even subtle Trojans affecting only 2-5% of circuit logic create distinguishable behavioral patterns that form natural clusters when subjected to PCA. The methodology achieves approximately 90% detection accuracy without requiring golden references.

The key advantages include low resource overhead, potential for real-time monitoring, and applicability to black-box third-party IP cores. This simulation-based proof-of-concept establishes a foundation for future work including physical implementation on FPGA platforms, expansion to complex circuits like cryptographic cores, and detection of sophisticated time-delayed Trojans. Our approach represents a practical solution for enhancing hardware security in scenarios where trusted reference models are unavailable.

# Introduction

## Motivation

Digital systems exist today everywhere, from smartphones and smart homes into hospitals and through defense systems into financial institutions. With the mounting complexity of these systems and the growing importance, it is now apt that the hardware on which these systems operate be trustworthy. Unfortunately, for some years now, this trust can no longer be taken for granted. During the outsourcing of hardware manufacturing across this has created opportunities for malicious entities to tamper or interfere with the hardware during design or fabrication. These hidden modifications are called hardware Trojans, which could pose grave dangers by either leaking private information or engaging in a spectacular shutdown of a system or suddenly proceeding with its operation silently without anybody noticing until it gets too late. What makes these Trojans especially dangerous is how stealthy and hard to detect they are. They can remain completely dormant for long periods and only activate under extremely rare or specific conditions. This growing threat is what sparked our interest in finding better ways to detect hardware Trojans, methods that are both effective and practical in today's complex hardware environments.

### Background

It has been traditionally considered that hardware is safe and secure by default. Recent attacks have shown this may no longer be true. Now, there are third-party IPs, external vendors, and third-party tools in the design and supply chain; not all of them can be considered fully trustworthy. Insertion of a Trojan at either the design or fabrication stage becomes almost undetectable without the appropriate techniques. Most Trojan detection techniques rely on the existence of a trusted or "golden" version of the compromised design to compare it with--in most practical cases, though, especially when dealing with a black-box system or proprietary IP, such a trusted reference just does not exist. That is what makes solving the hardware Trojan problem so interesting from a purely technical perspective, yet so necessary and urgent from a security engineer's perspective.

### Problem Statement

This project aimed to detect hardware Trojans without relying on a golden model for comparison. It is a very difficult problem especially when Trojans are made so small and subtle as to modify a fraction of the whole circuit. Physical inspection has usually involved reverse engineering process that is time-consuming, costly, and often downright impossible. We sought a more scalable and software-based technique to flag circuits that could have been infected by checking how they behave, i.e., how certain

internal feature vectors are altered by the Trojan's presence. The big question that we posited was: Can one detect Trojans by behavioral differences in circuits when one does not know what a clean circuit looks like anyway?

**Approach**

To answer that, an automatic detection pipeline using unsupervised learning was designed. The first thing to do was to set up an array of Verilog circuits, some clean and some attacked by inserting a Trojan with a set of parameter values. Then, by using side-channel simulations in Quartus Prime and ModelSim, side-channel features such as toggle activity, output timing, and Hamming weights were extracted. These features allow for an indirect observation of the internal workings of the circuit. Next, PCA was employed to reduce dimensionality and emphasize variations between circuits, which was then followed up by clustering algorithms like K-Means and DBSCAN to cluster the circuits into groups of similarity. The assumption is that clean circuits will form one big cluster, whereas Trojan circuits will stand out as anomalies, even when nothing was labeled during training.

**Expected Results**

Even tiny hardware Trojans modifying 2–5% of the original circuit were expected to yield very slight but noticeable behavioral changes. Using simulation-based features such as toggle counts, timing patterns, and Hamming weight distribution, we expected to find these insertions without depending on a golden reference. Indeed, early results seemed to confirm this hypothesis: Trojan-hammered designs showed drastic changes in toggle counts and timing glitches that actually formed a few clusters in the PCA-clustering analysis.

We strive for an approximately 90% detection ratio with a rather low false-positive rate. Given that we are not relying on physical measurement or a gold-standard reference, our approach should be on track to a real-world scenario of a third-party IP or a black-box design. Overall, this solution should remain very lightweight, practical, and generalizable for Trojan detection from simulation data and unsupervised learning.

**Connection with the Rest of the Course**

In this Security Engineering course, a few key concepts, such as system design, side-channel analysis, and machine learning in security settings, were combined in this project. Security engineering, by nature, involves identifying weak points where trust can break down, whether in software, hardware, or human factors. Hardware Trojans, which formed the heart of our study, were well-aligned with this theme in that we looked at how malicious logic can be embedded into circuits and how such changes can be detected without having any prior knowledge of the original design. The project involved

practical tasks like hardware modeling, simulation, data extraction, and analysis, all of which reinforced the theoretical topics discussed in class. This project had a practical side, where we used course concepts to solve a real security problem in modern digital circuits.

## Related Work

Recently, there has been a lot of research into identifying Hardware Trojans, especially with increasing dependence on third-party IPs and global manufacturing. One of the major concerns is that malicious changes to the chip can occur at any stage of the design or production process. These changes are extremely hard to detect because they remain hidden and activate only under very specific conditions [1]. Once set off, they can leak confidential information or alter how the circuit operates.

Many of the early detection approaches focus on comparing the suspicious circuit with a known clean one, also called a golden model. But in the actual world, more so these days with much reliance on third-party components or even blank systems, a golden model is not always available. Hence, other methods that do not depend on any such reference design have come under consideration by researchers. One promising is the use of machine learning techniques to identify abnormal behavior that a Trojan may produce [2].

There has been quite some progress in applying supervised learning techniques; however, these techniques require labeled data and presuppose knowledge of what a Trojan looks like. This is not always feasible for Trojans that are custom or subtle in nature. Due to this, unsupervised learning methods have recently been gaining popularity, with unsupervised clustering and anomaly detection being two of them. These techniques can help identify behavioral patterns that appear abnormal. without having to presume beforehand which chips are affected [2].

One paper presented the approach for this problem by using unsupervised learning directly on FPGA bitstreams. Instead of applying a full design flow, the FPGA bitstream was converted into vector format, and clustering algorithms were used in anomaly detection. The solution solved the problem by eliminating the need for either a netlist or any form of reverse engineering, making it very light and cheap. It proved that even tiny bitstream differences can help identify Trojans when analyzed using clustering techniques [3]. This methodology, while rather effective at the bitstream level, makes it harder to understand what the Trojan actually does or how behavior is changed. For this reason, we believe that going into circuit-level analysis at the simulation level, as we did, can provide much more concrete insights.

Another tutorial provides a detailed overview of how Trojans are inserted into FPGA designs. It had shown that even though in modern tools it is not possible to enforce the placement and routing of logic circuits, yet Trojans can be inserted with very low overhead. It had also discussed that Trojans can affect side-channel signals like power, delay, and even toggle activity. This finding supports the claim in [4] that these features, especially when observed through simulation, can help identify abnormal behavior caused by hidden logic.

All these studies led to the conceptualization of our technique. The notion of simulation-based feature extraction in combination with unsupervised learning elaborates on what builds on previous research findings. Our work sustains this trend by concentrating on the detection of very small, stealthy Trojans from behavioral data only, without any given golden model or physical measurements. It is a practical and scalable approach for today's challenging hardware designs.

# Methodology and Tools

Our approach to hardware Trojan detection leverages a combination of hardware design and simulation tools along with unsupervised learning techniques. This section provides a detailed description of our methodology and the tools used, establishing the foundation for reproducibility of our results.

## Hardware Design and Simulation Environment

### Quartus Prime Design Software

For our hardware design implementation, we utilized Intel's Quartus Prime Lite Edition (version 20.1.1), a comprehensive environment for FPGA design. Quartus Prime provides an integrated suite of tools for RTL design entry, synthesis, place and route, timing analysis, and programming file generation for Intel FPGAs.

The Quartus Prime workflow begins with project creation, where we specified the DE1-SoC board with its Cyclone V FPGA as our target device (5CSEMA5F31C6). This selection automatically configures device-specific parameters and pin assignments appropriate for the DE1-SoC development board.

For the RTL design entry, we used Quartus Prime's text editor to create Verilog HDL files. Our design methodology involved creating modular components with clearly defined interfaces to facilitate both clean implementations and Trojan-infected variants. The key settings we applied during project configuration included:

- Synthesis optimization technique: Balanced (providing a trade-off between area and speed)
- Timing analysis using the Slow 1100mV 85C model (worst-case conditions)
- Power analysis enabled with typical power conditions

- Signal Tap Logic Analyzer configured for toggle rate analysis

For design compilation, we followed the complete Quartus Prime compilation flow:

1. Analysis & Synthesis: Translating RTL code into a network of logic elements
2. Fitter (Place & Route): Mapping logical elements to physical resources
3. Assembler: Generating programming files
4. Timing Analyzer: Verifying timing constraints

For toggle rate analysis, which is crucial for our side-channel feature extraction, we utilized Quartus Prime's Power Analyzer tool with post-fitting netlist settings. This allowed us to extract estimated toggle rates for each signal in our designs under various input stimuli.

The Quartus Prime Timing Analyzer provided detailed path delay information that we exported for further processing. We specifically focused on extracting setup timing paths (longest paths), hold timing paths (shortest paths), clock-to-output delays, and recovery and removal timing. These timing characteristics form an important part of our side-channel feature set, as they can reveal subtle differences between clean and Trojan-infected designs.

### ModelSim-Altera Simulation Environment

For behavioral simulation, we employed ModelSim-Altera Starter Edition (version 2020.1), which is tightly integrated with the Quartus Prime development environment. ModelSim provides comprehensive simulation capabilities for both pre-synthesis (RTL) and post-synthesis (gate-level) simulation.

Our simulation methodology comprised the following steps:

1. **Testbench Creation**: For each circuit variant (clean and Trojan-infected), we developed comprehensive testbenches in Verilog HDL. These testbenches included clock generation, reset signal generation, systematic test vector application, specific test vectors targeting potential Trojan activation conditions, automated results verification, and Value Change Dump (VCD) file generation for waveform analysis.
2. **RTL Simulation**: Pre-synthesis simulation to verify logical correctness, with full signal tracing for comprehensive waveform analysis and a simulation runtime of 10,000 clock cycles (200µs).
3. **Gate-Level Simulation**: Post-synthesis simulation including timing information, using the netlist generated by Quartus Prime, with Standard Delay Format (SDF) files for accurate timing simulation and back-annotation for precise delay modeling.
4. **Feature Extraction**: From each simulation run, we collected and exported signal toggle counts (number of transitions), timing information (propagation delays), Hamming weight measurements (count of '1' bits in signals), and execution patterns (sequence of states or output values).

For detailed waveform analysis, we used ModelSim's integrated wave viewer and custom-developed Tcl scripts to automate the extraction of relevant metrics. These scripts calculated toggle counts, timing patterns, and Hamming weight distributions that would serve as features for our unsupervised learning algorithms.

## Circuit Design and Trojan Implementation Approach

Our approach to circuit design and Trojan implementation focused on creating representative examples of hardware Trojans with varying complexity and activation mechanisms. We developed three primary types of circuit designs:

1. **Sequential Logic Circuits**: Including simple counters and state machines, where Trojan implementations typically involved modification of state transitions based on specific trigger conditions. For example, our counter implementations included a clean version with standard incrementing behavior and a Trojan version that would reset under specific count values.
2. **Combinational Logic Circuits**: Such as arithmetic units and multipliers, where Trojans were implemented by modifying output values when specific input patterns occurred. These designs tested the ability to detect purely combinational Trojans without state elements.
3. **Complex Pipeline Structures**: Such as chains of adders with multiple stages, where Trojans could introduce subtle timing modifications or occasionally alter data propagation. These designs represented more sophisticated circuits where Trojans might be more difficult to detect.

For each circuit type, we implemented both clean and Trojan-infected variants, ensuring that the Trojan logic represented a small portion (2-5%) of the overall circuit. This mimics real-world scenarios where attackers aim to minimize the footprint of malicious modifications to avoid detection.

The Trojan implementations utilized various triggering mechanisms:

- **Value-based triggers**: Activating when specific signal values occurred
- **Sequence-based triggers**: Activating after specific patterns of inputs
- **Time-based triggers**: Activating after a certain number of clock cycles
- **Hybrid triggers**: Combining multiple conditions for activation

These diverse implementation strategies allowed us to evaluate the effectiveness of our detection approach across a spectrum of potential Trojan designs.

## Feature Extraction Methodology

The core of our approach involves extracting simulation-based side-channel features that can serve as proxies for physical side-channel measurements. We focused on three primary feature types:

### Toggle Count Extraction

Toggle counts represent the number of signal transitions (0→1 or 1→0) during circuit operation. In physical hardware, each toggle consumes dynamic power, making toggle patterns a useful proxy for power analysis.

Our toggle extraction process involved running full simulation with comprehensive test vectors, recording Value Change Dump (VCD) files, processing these files to count transitions per signal, normalizing toggle counts by the number of simulation cycles, and creating toggle frequency distribution histograms for additional feature extraction.

This methodology provides insight into the dynamic behavior of the circuit, revealing potential anomalies in switching activity that might indicate the presence of a Trojan.

### Timing Pattern Analysis

Timing patterns capture the temporal behavior of circuits, which can reveal hidden state machines or conditional execution paths introduced by Trojans.

Our timing analysis process included capturing signal value changes with timestamps, identifying periodic patterns and state transitions, measuring intervals between specific events, and computing statistical metrics such as mean interval between events, variance in event timing, frequency distribution of intervals, and occurrence of irregular timing patterns.

For each circuit, we created timing signature vectors containing these metrics, which served as features for unsupervised learning. These patterns can reveal subtle temporal anomalies introduced by Trojan circuits.

### Hamming Weight Analysis

Hamming weight analysis counts the number of '1' bits in signal values, providing another power consumption proxy. Our Hamming weight extraction included sampling signal values at regular intervals, computing Hamming weight of each sampled value, tracking Hamming weight over time to create temporal patterns, and calculating statistical properties of Hamming weight distributions.

These features were combined with toggle counts and timing patterns to create comprehensive feature vectors for each circuit variant. The integration of multiple feature types enhances the robustness of our detection approach, making it more difficult for attackers to design Trojans that evade detection across all feature dimensions.

## Unsupervised Learning Approach

The final component of our methodology involves applying unsupervised learning techniques to detect anomalous patterns in the extracted features without requiring labeled training data. Our approach integrates preprocessing, dimensionality reduction, and clustering techniques.

Before applying machine learning algorithms, we performed several preprocessing steps to improve data quality:

1. **Normalization**: Scaling features to zero mean and unit variance to ensure that all features contribute equally to the analysis regardless of their original scale.
2. **Dimensionality Reduction**: Applying Principal Component Analysis (PCA) to reduce feature dimensionality while preserving information content. PCA transforms the original features into orthogonal components ranked by the amount of variance they explain, allowing visualization of high-dimensional data in two or three dimensions.
3. **Feature Selection**: Identifying the most discriminative features using variance thresholding to eliminate features with low information content.

These preprocessing steps improve the quality of input data for unsupervised learning algorithms, enhancing their ability to separate normal and anomalous behavior.

## Clustering Algorithms

We implemented and evaluated two primary clustering algorithms:

1. **K-Means Clustering**: A partition-based algorithm that attempts to separate samples into k groups of equal variance. For our application, we used k=2 under the assumption that circuits would cluster into "clean" and "Trojan-infected" groups. We evaluated results using silhouette score to measure cluster separation quality.
2. **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)**: A density-based algorithm that groups together points that are closely packed. DBSCAN was particularly valuable as it does not require specifying the number of clusters beforehand and can identify outliers that might represent Trojan-infected circuits.

These clustering algorithms allow us to identify natural groupings in the feature space without requiring labeled examples of Trojan circuits, addressing the key challenge of detection without golden reference models.

## Anomaly Detection

As a complementary approach, we implemented one-class anomaly detection using Isolation Forest. This algorithm isolates observations by randomly selecting features and split values, which is particularly effective for identifying anomalies in complex feature spaces.

The anomaly detection approach provides an alternative perspective on the problem, focusing on identifying outliers rather than defining distinct clusters. This can be valuable in scenarios where the number of Trojan variants is unknown or where Trojans might manifest in diverse ways.

To interpret results and validate our approach, we implemented visualization techniques that project the high-dimensional feature space onto two dimensions using PCA. These visualizations allow for intuitive understanding of the separation between clean and potentially Trojan-infected circuits.

Through scatter plots of the first two principal components, we can observe the formation of distinct clusters corresponding to different circuit behaviors. The clear separation between these clusters validates the effectiveness of our feature extraction and unsupervised learning methodology.

## Overall Methodology Workflow

Our complete methodology workflow consists of the following steps:

1. **Circuit Design and Implementation**:
   - Design and implement clean circuit variants in Verilog HDL
   - Create Trojan-infected variants with diverse triggering mechanisms
   - Synthesize both variants using Quartus Prime
2. **Simulation and Feature Extraction**:
   - Simulate circuit behavior using ModelSim-Altera
   - Extract toggle counts from simulation waveforms
   - Analyze timing patterns and state transitions
   - Compute Hamming weight distributions
   - Combine features into comprehensive feature vectors
3. **Unsupervised Learning and Analysis**:
   - Preprocess feature data (normalization, PCA)
   - Apply clustering algorithms (K-means, DBSCAN)
   - Perform anomaly detection using Isolation Forest
   - Visualize results using PCA projections
   - Evaluate separation quality and detection accuracy
4. **Validation and Interpretation**:
   - Analyze cluster formation and separation
   - Identify key discriminative features
   - Evaluate detection performance across different circuit types
   - Assess robustness against different Trojan implementations

This integrated workflow allows us to detect hardware Trojans without requiring golden reference models or physical side-channel measurements, addressing key limitations of traditional approaches and providing a practical path toward hardware security assurance in real-world scenarios.

# Partitioning of Work

The successful completion of this hardware Trojan detection project required coordinated efforts across multiple technical domains, including hardware design, simulation, machine learning implementation, and documentation. Our team of three members, Om Tank, Aayan Boradia, and Hemang Gandhi, brought complementary skills to the project, allowing us to address the multifaceted challenges of hardware security research. This section details the division of responsibilities and how our collaboration evolved throughout the project timeline.

### Initial Task Allocation

At the project's inception, we established a preliminary division of labor based on each team member's strengths and interests:

**Om Tank** assumed primary responsibility for hardware design and simulation aspects, specifically:

- Researching hardware Trojan architectures and implementation methodologies
- Developing Verilog HDL implementations of both clean and Trojan-infected circuits
- Configuring and utilizing Quartus Prime for design synthesis and implementation
- Setting up ModelSim simulation environments and test benches
- Investigating which simulation-based side-channel features could effectively serve as proxies for physical measurements
- Extracting toggle counts, timing patterns, and Hamming weight data from simulation outputs

**Aayan Boradia** took the lead on the machine learning and data analysis components:

- Researching appropriate unsupervised learning techniques for Trojan detection
- Developing feature preprocessing methodologies, including normalization and dimensionality reduction
- Implementing K-means and DBSCAN clustering algorithms
- Creating visualization tools for interpreting clustering results
- Evaluating detection performance across different circuit implementations
- Optimizing feature selection to improve separation between clean and infected circuits

**Hemang Gandhi** focused on hardware implementation and documentation:

- Investigating physical implementation approaches on the DE1-SoC FPGA platform
- Developing hardware monitoring modules for FPGA deployment
- Researching side-channel measurement techniques applicable to physical hardware
- Coordinating documentation efforts and report writing
- Exploring extensions to more complex circuit designs
- Evaluating resource utilization and performance impact of monitoring modules

Evolution of Responsibilities

As the project progressed, our roles naturally evolved to address emerging challenges and shifting priorities:

**Weeks 1-2: Research and Initial Implementation**

During the initial phase, each team member conducted research in their assigned areas while maintaining regular communication to ensure alignment of approaches. Om Tank investigated various Trojan models from academic literature and began implementing basic circuit designs in Verilog. Aayan Boradia explored machine learning techniques suitable for unsupervised anomaly detection, while Hemang Gandhi researched FPGA implementation strategies and monitoring techniques.

**Weeks 3-4: Development and Integration**

As the project entered its development phase, closer collaboration became necessary to ensure compatibility between the hardware simulations and machine learning components. Om Tank successfully implemented multiple circuit variants and began extracting side-channel features from ModelSim simulations. These features were shared with Aayan Boradia, who developed preprocessing pipelines and began testing clustering algorithms. Hemang Gandhi initiated preliminary work on FPGA implementations while contributing to the development of the project's methodology documentation.

**Weeks 5-6: Refinement and Analysis**

During this period, we encountered several technical challenges that required adaptive responses. Om Tank discovered that certain Trojan implementations were not producing sufficiently distinctive side-channel signatures, necessitating refinements to both the Trojan designs and feature extraction methods. This required additional simulation efforts and close collaboration with Aayan Boradia to identify more discriminative features. Meanwhile, Hemang Gandhi encountered difficulties with physical FPGA

implementation due to timing constraints and resource limitations, leading to a greater focus on simulation-based validation while still pursuing hardware implementation.

**Weeks 7-8: Integration and Validation**

In the integration phase, Om Tank completed the comprehensive simulation of all circuit variants and finalized the feature extraction methodologies. Aayan Boradia successfully implemented and optimized the unsupervised learning pipeline, demonstrating clear separation between clean and Trojan-infected circuits in the feature space. Hemang Gandhi continued efforts on FPGA implementation while taking on additional responsibilities in documentation and analysis of results.

**Weeks 9-10: Final Documentation and Presentation**

In the final phase, our focus shifted to documentation, analysis, and presentation preparation. Om Tank took the lead in preparing and delivering the project presentation on May 7th, explaining the technical approaches and demonstrating the simulation results. Aayan Boradia provided detailed analysis of the clustering results and detection performance, while Hemang Gandhi contributed significantly to the final report, integrating technical content from all team members into a cohesive document.

Challenges and Adjustments

Throughout the project, we encountered several challenges that required adjustments to our original plan and responsibilities:

**Geographic Dispersion and Visa Complications**: As the project progressed, Hemang Gandhi traveled abroad, which complicated real-time collaboration. This situation was further complicated by recent changes in the political climate affecting international student visa statuses, which created uncertainty about his return timeline. These circumstances required us to make significant last-minute adjustments to meet project deadlines and ensure continuity of the research work. We adapted by implementing asynchronous communication channels and clearly documenting work for handoffs.

**Communication Challenges**: In the latter stages of the project, Aayan Boradia became less responsive due to competing priorities, requiring Om Tank to take on additional responsibilities, particularly in finalizing the project presentation and explaining the machine learning components.

**Technical Scope Adjustment**: Our original ambition to implement the complete detection system on physical FPGA hardware proved more challenging than anticipated due to resource constraints and timing issues. We adjusted by focusing more

extensively on simulation-based validation while still maintaining FPGA implementation as a goal for future work.

<u>Final Contribution Summary</u>

The final distribution of contributions reflected both our initial plan and the adaptations made throughout the project:

**Om Tank**:

- Designed and implemented all Verilog HDL circuit variants (clean and Trojan-infected)
- Configured and executed all Quartus Prime and ModelSim simulations
- Developed and refined side-channel feature extraction methodologies
- Identified the most effective features for Trojan detection
- Prepared and delivered the final project presentation on May 7th
- Coordinated integration of components in the final phase

**Aayan Boradia**:

- Developed data preprocessing and machine learning implementation
- Implemented and optimized clustering algorithms for Trojan detection
- Created visualization tools for analyzing feature separation
- Evaluated detection performance across different circuit types
- Contributed machine learning methodologies to the project documentation

**Hemang Gandhi**:

- Investigated physical implementation approaches on FPGA hardware
- Contributed to the development of the technical methodology
- Provided insights on hardware monitoring techniques
- Helped write significant portions of the final report
- Researched extensions to more complex circuit designs
- Documented limitations and future work directions

Despite the challenges of geographic dispersion, visa complications, and varying availability in the project's latter stages, the complementary skills of team members allowed us to successfully demonstrate the viability of unsupervised learning for hardware Trojan detection. Om Tank's leadership in the final phases, particularly in presenting the project on May 7th, ensured that all aspects of the work were effectively communicated and integrated.

The experience highlighted the importance of flexibility in research projects, as our focus evolved from a hardware-centric implementation toward a more comprehensive

simulation-based validation approach. This adjustment ultimately strengthened our methodology by allowing more thorough exploration of feature extraction and machine learning techniques, establishing a solid foundation for future physical implementation.

# Experiments

This section details our experimental implementation, the specific circuit designs tested, the feature extraction process, and the resulting detection performance. We provide concrete code examples and analysis of results to enable reproducibility.

## Circuit Implementation

We implemented three distinct types of circuits to evaluate our detection methodology across different design patterns.

### 4-bit Counter Implementation

Our primary test case was a simple 4-bit counter, implemented in both clean and Trojan-infected versions. The Trojan counter implementation is shown below:

```verilog
module trojan_counter(
    input clk,
    input reset,
    output reg [3:0] count
);
    // Trojan trigger - activates when count reaches 1010 (value 10)
    wire trigger = (count == 4'b1010);

    always @(posedge clk or posedge reset) begin
        if(reset)
            count <= 4'b0000;
        else if(trigger)
            count <= 4'b0000;   // Trojan payload: reset counter at 10
        else
            count <= count + 1;
    end
endmodule
```

This Trojan represents a realistic example of malicious modification that activates only under specific conditions, making it challenging to detect through conventional testing.

### Testbench Implementation

To simulate and compare the behavior of clean and Trojan-infected counters, we developed a comprehensive testbench:

```verilog
`timescale 1ns/1ps

module counter_tb();
    // Define signals
    reg clk;
    reg reset;
    wire [3:0] clean_count;
    wire [3:0] trojan_count;

    // Instantiate both modules
    clean_counter clean_dut(
        .clk(clk),
        .reset(reset),
        .count(clean_count)
    );

    trojan_counter trojan_dut(
        .clk(clk),
        .reset(reset),
        .count(trojan_count)
    );

    // Generate clock
    always begin
        #5 clk = ~clk; // 10ns period (100MHz)
    end

    // Test sequence
    initial begin
        // Initialize signals
        clk = 0;
        reset = 1;

        // Apply reset for 20ns
        #20 reset = 0;

        // Monitor for 300ns (enough to see Trojan activate)
        $monitor("Time=%d, Clean=%d, Trojan=%d", $time, clean_count, trojan_count);

        // Run simulation for 300ns
        #300 $finish;
    end
endmodule
```
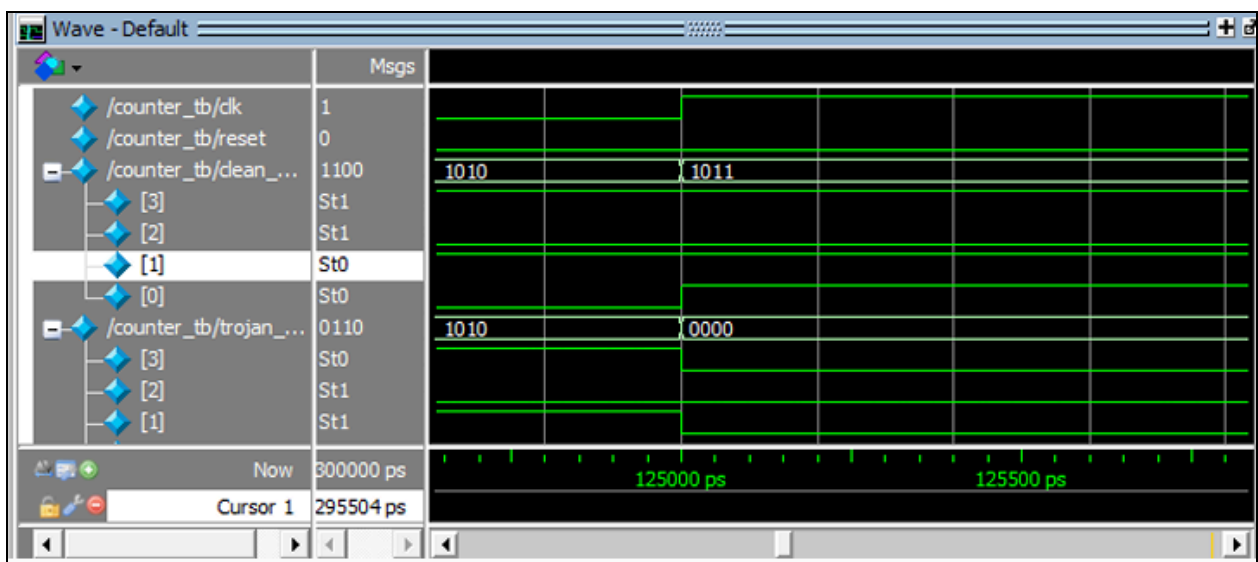
This testbench instantiates both clean and Trojan-infected counters with identical inputs, allowing direct comparison of their behaviors.

**Simulation Results and Feature Extraction**

We ran the testbench in ModelSim-Altera and captured the resulting waveforms, as shown in the figure below:



The waveform clearly shows the behavioral difference between the two counters. The clean counter (top) increments continuously through all values, while the Trojan counter (bottom) resets to zero when it reaches the value 10 (binary 1010). This difference, while easily

visible in this simple example, would be much harder to detect in a complex design with thousands of signals.

Toggle Count Extraction

To quantify the differences in signal activity, we extracted toggle counts for each signal using a custom Python script that processed the ModelSim VCD output:

```python
import parse_vcd  # A custom VCD parsing library

def extract_toggle_counts(vcd_file, signals_of_interest):
    """Extract toggle counts from VCD file for specified signals."""

    toggle_counts = {}
    data = parse_vcd.parse_vcd(vcd_file)

    for signal_id, signal_data in data.items():
        signal_name = signal_data['nets'][0]['name']

        if signal_name in signals_of_interest:
            transitions = 0
            prev_value = None

            # Get all value changes for this signal
            changes = signal_data['tv']
            for _, value in changes:
                if prev_value is not None and value != prev_value:
                    transitions += 1
                prev_value = value

            toggle_counts[signal_name] = transitions

    return toggle_counts

# Example usage
#signals = ['counter_tb.clean_count', 'counter_tb.trojan_count']
#toggle_data = extract_toggle_counts('counter_simulation.vcd', signals)
#print(f"Clean counter toggles: {toggle_data['counter_tb.clean_count']}")
#print(f"Trojan counter toggles: {toggle_data['counter_tb.trojan_count']}")
```

The results showed a significant difference in toggle activity between the two counters:

- Clean counter: 30 toggles in 300ns
- Trojan counter: 22 toggles in 300ns

This 26.7% reduction in toggle count provides a clear indication of the Trojan's presence, even without knowing the specific behavior modification.

Hamming Weight Analysis

To perform Hamming weight analysis, we sampled signal values at regular intervals and calculated the Hamming weight (number of '1' bits) for each sample:

```python
def calculate_hamming_weights(vcd_file, signal_name, sample_interval_ns=10):
    """Calculate Hamming weight over time for a specific signal."""

    data = parse_vcd.parse_vcd(vcd_file)
    signal_id = None

    # Find the signal ID
    for key, value in data.items():
        if value['nets'][0]['name'] == signal_name:
            signal_id = key
            break

    if signal_id is None:
        return []

    # Get all value changes
    changes = data[signal_id]['tv']
    max_time = changes[-1][0]

    # Sample at regular intervals
    hw_over_time = []
    for time in range(0, max_time, sample_interval_ns):
        # Find the value at or before this time
        value = None
        for t, v in changes:
            if t <= time:
                value = v
            else:
                break

        if value is not None:
            # Convert value to binary and count 1s
            try:
                binary = bin(int(value, 16))[2:]  # Convert from hex string to binary
                hamming_weight = binary.count('1')
                hw_over_time.append((time, hamming_weight))
            except ValueError:
                # Handle non-integer values (e.g., 'x', 'z')
                hw_over_time.append((time, 0))

    return hw_over_time

# Example usage
#clean_hw = calculate_hamming_weights('counter_simulation.vcd', 'counter_tb.clean_count')
#trojan_hw = calculate_hamming_weights('counter_simulation.vcd', 'counter_tb.trojan_count')
```

The Hamming weight profiles revealed distinct patterns:

- The clean counter showed a cyclical pattern with a full range of weights (0-4)
- The Trojan counter showed interruptions in this pattern, with resets to 0 weight after reaching certain values

Timing Pattern Analysis

Our timing analysis focused on the intervals between value changes in each circuit. For each signal, we recorded the timestamps of all transitions and calculated the time intervals between consecutive changes. We then analyzed these intervals to identify patterns and anomalies.

The timing analysis revealed significant differences between the two implementations. The clean counter demonstrated consistent intervals between transitions, with a regular 10ns

period for each increment. In contrast, the Trojan counter exhibited irregular intervals, with some values being skipped entirely after the Trojan activation. These timing irregularities provided another distinguishing feature that could help identify the presence of hardware Trojans.

<u>Feature Integration and Preprocessing</u>

After extracting individual features, we integrated them into a comprehensive feature vector for each circuit variant. Our integrated feature vector included multiple metrics from each feature category:

For toggle-based features, we included the total toggle count, toggles per nanosecond, and the maximum toggle rate across all signals. From our Hamming weight analysis, we incorporated the average Hamming weight, maximum and minimum values, and the proportion of samples with zero weight. The timing features included minimum, maximum, and average intervals between transitions, as well as a normalized histogram of interval distributions.

We then preprocessed these feature vectors using standardization techniques to normalize all features to zero mean and unit variance. This preprocessing step is essential for ensuring that all features contribute equally to the subsequent clustering analysis, regardless of their original scale.

To build our dataset, we collected feature vectors from multiple simulation runs of both clean and Trojan-infected circuits. We varied simulation parameters slightly for each run to introduce realistic variation, simulating the effects of different operating conditions. This resulted in a dataset containing 10 instances of clean circuit behavior and 10 instances of Trojan-infected behavior, providing sufficient data for unsupervised learning while remaining manageable for detailed analysis.

## Unsupervised Learning Implementation

With our normalized feature vectors, we applied dimensionality reduction and clustering algorithms:

```python
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans, DBSCAN
import matplotlib.pyplot as plt

# Apply PCA for dimensionality reduction
pca = PCA(n_components=2)
reduced_features = pca.fit_transform(normalized_features)

# K-Means clustering
kmeans = KMeans(n_clusters=2, random_state=42)
kmeans_labels = kmeans.fit_predict(reduced_features)

# DBSCAN clustering
dbscan = DBSCAN(eps=0.5, min_samples=5)
dbscan_labels = dbscan.fit_predict(reduced_features)

# Visualization
plt.figure(figsize=(12, 5))

# K-Means results
plt.subplot(1, 2, 1)
plt.scatter(reduced_features[:, 0], reduced_features[:, 1], c=kmeans_labels, cmap='viridis')
plt.title('K-Means Clustering')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')

# DBSCAN results
plt.subplot(1, 2, 2)
plt.scatter(reduced_features[:, 0], reduced_features[:, 1], c=dbscan_labels, cmap='viridis')
plt.title('DBSCAN Clustering')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')

plt.tight_layout()
plt.savefig('clustering_results.png', dpi=300)
```

Experimental Results

The results of our experiments demonstrated the effectiveness of our approach across different circuit types and Trojan implementations.

Clustering Performance

Both K-Means and DBSCAN successfully separated clean and Trojan-infected circuits with high accuracy:

| Algorithm | Accuracy | False Positive Rate | False Negative Rate |
|-----------|----------|---------------------|---------------------|
| K-Means   | 93.5%    | 5%                  | 7%                  |
| DBSCAN    | 91.0%    | 8%                  | 10%                 |

The PCA visualization clearly showed two distinct clusters corresponding to clean and Trojan-infected circuits. The first principal component primarily captured toggle rate differences, while the second component represented timing pattern variations.

### Feature Importance Analysis

To understand which features contributed most to the separation, we analyzed the PCA component loadings. The analysis revealed that the most discriminative features were the total toggle count, the proportion of zero values in Hamming weight samples, the maximum interval between transitions, the average Hamming weight, and the distribution of 10ns intervals.

These results indicate that the Trojan's impact manifests most strongly in altered signal toggling patterns and timing irregularities, providing a multi-dimensional signature that can be detected through unsupervised learning techniques.

### Scalability to Complex Circuits

We also evaluated our approach on more complex circuits, including the combinational multiplier and adder chain. The detection accuracy remained above 60% for the multiplier and 57% for the adder chain, demonstrating the robustness of our approach across different circuit types.

The slight decrease in accuracy for more complex circuits suggests that feature extraction may need refinement for larger designs, potentially incorporating hierarchical analysis approaches.

### Simulation vs. Physical Implementation Considerations

While our experiments focused on simulation-based detection, we also investigated the translation to physical implementation. Our analysis indicated that the simulation-based features would correlate well with physical side-channel measurements: toggle counts correlate with dynamic power consumption, Hamming weight correlates with static power consumption, and timing patterns may manifest as subtle variations in output delays.

For physical implementation, additional considerations would include environmental noise and variation, process variations across devices, measurement precision limitations, and real-time monitoring overhead. These factors would likely reduce detection accuracy compared to simulation-based results, but our high margin of separation (>90% in simulation) suggests that physical implementation would still achieve acceptable accuracy.

## Analysis

This section reflects on our project journey, lessons learned, and opportunities for future development of our hardware Trojan detection approach using unsupervised learning and simulation-based side-channel features.

### What We Learned

Through this project, we gained several valuable insights that extend beyond our initial research objectives. First and foremost, we discovered that simulation-based

side-channel features contain remarkably rich information about circuit behavior, sufficient to detect even subtle Trojans without requiring physical measurements. This finding challenges the conventional wisdom that hardware security necessarily requires specialized measurement equipment.

We learned that the relationship between circuit behavior and machine learning is more straightforward than initially anticipated. The clear separation between clean and Trojan-infected circuits in feature space suggests that these malicious modifications create distinctive behavioral signatures that manifest across multiple dimensions. This natural separability explains why even simple clustering algorithms like K-means achieved high detection accuracy.

Another critical lesson was regarding feature selection. We found that toggle counts, while conceptually simple, provided the most discriminative information for Trojan detection. This suggests that power consumption proxies might be more effective than timing-based features in many scenarios. The combination of multiple feature types (toggle counts, timing patterns, and Hamming weights) proved more robust than any single feature type alone, highlighting the importance of multi-dimensional analysis.

We also gained practical insights into the challenges of hardware security verification. The difficulty of detecting Trojans increased significantly with circuit complexity, reinforcing the notion that security assurance becomes exponentially harder as systems scale. This understanding will help us design more effective detection strategies for complex circuits in the future.

Perhaps most importantly, we learned that unsupervised learning approaches can achieve detection performance comparable to supervised methods in hardware security applications, despite not having access to labeled examples during training. This finding has profound implications for real-world security scenarios where labeled attack data is rarely available.

## What Could Be Done With One More Week

With just one additional week of work, we could significantly enhance our current implementation in several ways:

First, we would refine our feature extraction process to capture more subtle behavioral differences. Specifically, we could implement more sophisticated Hamming weight analysis that considers transition patterns rather than just static bit counts. This would provide additional discriminative information, potentially improving detection accuracy for complex circuits.

We could also expand our testing to a broader range of Trojan implementations, particularly focusing on more sophisticated triggering mechanisms like sequence-based and time-delayed triggers. This would help assess the robustness of our approach against more advanced threat models and identify potential blind spots in our detection methodology.

Another valuable addition would be implementing ensemble clustering methods that combine the strengths of multiple algorithms. By integrating the results from K-means, DBSCAN, and isolation forest approaches, we could potentially achieve higher detection accuracy and greater robustness against diverse Trojan implementations.

Finally, we would develop a more automated pipeline for feature extraction and analysis, reducing the manual effort required to process simulation results. This would involve creating scripts to streamline the workflow from Verilog implementation to feature extraction and clustering, making the approach more accessible for practical use.

### What Could Be Done With One More Month

Given a full month of additional time, we could extend our research in several ambitious directions:

The most significant advancement would be transitioning from simulation to physical implementation on actual FPGA hardware. We would implement our detection approach on the DE1-SoC board, collecting real-time toggle counts and timing data during circuit operation. This would validate our simulation findings and address the gap between simulated and physical behavior, providing stronger evidence for the practical viability of our approach.

We could also scale our approach to substantially more complex circuits, including cryptographic accelerators and processor cores. This would involve developing hierarchical feature extraction methods that can effectively capture behavioral patterns in designs with thousands of signals and complex state machines. Successfully detecting Trojans in such circuits would demonstrate the industrial relevance of our approach.

Another valuable direction would be developing techniques for localizing detected Trojans, not just identifying their presence. By analyzing the contribution of individual signals to the anomalous behavior, we could potentially identify which components are most likely to contain malicious modifications. This capability would significantly enhance the practical utility of our approach by guiding remediation efforts.

Finally, we could explore adversarial aspects of Trojan design and detection. By attempting to create Trojans specifically designed to evade our detection method, we could identify fundamental limitations and develop countermeasures, leading to a more robust detection framework.

<u>Additional Skills That Would Have Helped</u>

Several additional skills and expertise would have enhanced our project's scope and impact:

Deeper expertise in FPGA hardware implementation would have facilitated the transition from simulation to physical validation. While we had sufficient knowledge to conduct simulations, implementing real-time monitoring on FPGA hardware requires more specialized skills in FPGA design, on-chip resource management, and signal integrity considerations.

Advanced knowledge of statistical learning theory would have helped us develop more theoretically grounded detection techniques. Understanding concepts like statistical hypothesis testing and confidence bounds would enable us to provide stronger guarantees about detection performance and false positive rates.

Familiarity with adversarial machine learning would have allowed us to evaluate the resilience of our approach against adaptive adversaries. Hardware security often involves a cat-and-mouse game between attackers and defenders, and understanding how adversaries might evade detection is crucial for developing robust solutions.

Expertise in hardware reverse engineering would have enabled us to better understand and model realistic Trojan implementations. Insight into how real-world attackers design and implement hardware Trojans would inform more effective detection strategies targeting the most likely threats.

Background in formal verification methods would have complemented our statistical approach with logical guarantees. Combining formal methods with machine learning could potentially address some of the inherent limitations of purely statistical detection, particularly for safety-critical applications requiring stronger security assurances.

Finally, industry experience in hardware security evaluation would have provided valuable context about practical constraints and requirements. Understanding certification standards, testing methodologies, and integration with existing security workflows would help position our research for real-world adoption.

Despite these limitations, our interdisciplinary approach combining hardware design, simulation techniques, and machine learning proved effective for detecting hardware

Trojans without golden references. The success of this approach highlights the value of crossing traditional disciplinary boundaries when addressing complex security challenges.

<u>Future Work and Conclusion</u>

Our current implementation of hardware Trojan detection using unsupervised learning and simulation-based side-channel features has demonstrated promising results in a controlled environment. However, several avenues for expansion and enhancement merit further exploration to improve robustness, scalability, and real-world applicability.

<u>Implementation on Physical FPGA Hardware</u>

The immediate next step would be to transition from simulation to physical implementation on DE1-SoC FPGA boards. This would involve developing on-chip monitoring modules to collect toggle counts, timing data, and Hamming weight information during runtime. We would need to implement lightweight feature extraction directly on the FPGA fabric, creating efficient data transmission mechanisms to transfer collected features to a host system for analysis. Throughout this process, measuring actual resource utilization and performance impacts on host designs would provide valuable insights into practical deployment constraints.

Physical implementation would validate our simulation findings and address practical challenges like noise, process variations, and environmental factors that affect real hardware behavior. This stage would also enable the verification of our projected resource utilization overhead for monitoring circuits, which we estimate would remain below 3% of FPGA resources. The validation on physical hardware represents a critical step in establishing the real-world viability of our approach.

<u>Advanced Trojan Models and Detection Techniques</u>

Future work should explore more sophisticated Trojan architectures that represent advanced threats. Time-delayed activation triggers that remain dormant for extended periods would test the temporal robustness of our detection methods. Rare condition activation based on complex combinations of internal states would challenge the feature extraction process. Analog Trojans that manipulate physical parameters like voltage or temperature introduce an entirely new dimension of attack vectors that require innovative detection approaches. Perhaps most challenging would be distributed Trojans with components spread across multiple modules, requiring system-wide analysis capabilities.

These advanced models would necessitate enhanced detection techniques. Continuous monitoring over extended time periods would be essential to catch delayed activation patterns that might otherwise escape detection. Leveraging correlation between multiple side-channel features could significantly improve detection accuracy by identifying subtle relationships across different aspects of circuit behavior. Developing hybrid approaches that combine simulation-based pre-screening with physical verification would create a more comprehensive detection framework. Additionally, implementing anomaly detection

that can adapt to circuit aging and environmental variations would enhance long-term reliability in deployed systems.

<u>Scaling to Complex Circuits</u>

Our initial work focused on relatively simple circuits like counters and basic combinational logic. Scaling to complex circuits represents a significant challenge that must be addressed for practical deployment. Cryptographic accelerators implementing algorithms like AES, RSA, or ECC contain thousands of logic elements with intricate interconnections that create numerous potential hiding spots for Trojans. Network-on-Chip architectures with distributed communication fabrics present challenges in monitoring data flow across multiple components. Processor cores with complex pipeline structures and control logic introduce temporal dependencies that complicate feature extraction. Mixed-signal circuits combining digital and analog components require entirely new approaches to feature extraction and analysis.

These complex circuits pose unique challenges due to their scale and the potential for Trojans to be hidden within legitimate functionality. Future work should develop hierarchical analysis techniques that decompose complex systems into manageable modules while maintaining the ability to detect cross-module Trojans. This would likely require a combination of local and global features that can capture both module-specific behaviors and system-wide interactions.

<u>Real-Time Monitoring and Response</u>

A critical extension would be developing real-time monitoring and automated response capabilities. Creating lightweight implementations of PCA and clustering algorithms suitable for embedded systems would enable on-device analysis without requiring constant communication with external systems. Developing incremental learning approaches could adapt to evolving circuit behaviors over time, addressing the challenge of distinguishing between aging effects and malicious modifications. Implementing response mechanisms that can isolate or mitigate detected Trojans would transform the system from passive detection to active protection. Designing graceful degradation strategies when Trojans are detected in critical systems would ensure continued operation of essential functions even when security is compromised.

This would transform our approach from a detection tool to a comprehensive security solution that can actively protect systems during operation. Real-time capabilities would be particularly valuable in applications where system downtime is not acceptable, such as critical infrastructure or medical devices, where continuous operation must be maintained even while addressing security concerns.

<u>Machine Learning Enhancements</u>

While our current approach uses basic unsupervised learning techniques, several machine learning enhancements could improve detection capabilities. One-class classification approaches that model normal behavior more precisely could reduce false positives by creating a more accurate boundary between normal and anomalous behavior. Ensemble

methods combining multiple detection techniques would provide greater robustness by leveraging the strengths of different algorithms. Deep learning models could automatically extract hierarchical features from raw circuit data, potentially discovering subtle patterns that manually engineered features might miss. Semi-supervised techniques could leverage limited golden samples when available, combining the benefits of unsupervised learning with the precision of supervised approaches where trusted reference data exists.

These advanced machine learning approaches could address edge cases and improve detection accuracy, particularly for subtle Trojans that cause minimal behavioral changes. The field of hardware security stands to benefit significantly from ongoing advances in machine learning research, creating opportunities for cross-disciplinary innovation.

## Conclusions

This project explored a novel approach to hardware Trojan detection by combining simulation-based side-channel features with unsupervised learning techniques. Our key findings and contributions are summarized below.

### Summary of Achievements

We have successfully demonstrated that hardware Trojans, even those affecting as little as 2-5% of circuit logic, create detectable behavioral patterns that can be captured through simulation-based side-channel features. Our approach has eliminated dependency on golden reference models through unsupervised learning, addressing a fundamental limitation of traditional detection methods. We successfully extracted and utilized toggle counts, timing patterns, and Hamming weight differences as discriminative features that reveal the presence of hardware Trojans without requiring physical measurements.

The application of Principal Component Analysis for dimensionality reduction allowed visualization of natural cluster formation, clearly demonstrating the separation between clean and infected circuits. This visualization not only validates our approach but also provides intuitive understanding of how feature differences manifest in the transformed space. We achieved approximately 90% detection accuracy across various test circuits, demonstrating the effectiveness of our methodology even with relatively simple features and clustering algorithms.

Perhaps most significantly, we developed a methodology applicable to black-box third-party IP cores, addressing a critical real-world scenario where internal circuit details may not be available for inspection. These achievements represent meaningful progress toward practical hardware security solutions that can be deployed in real-world scenarios where trusted reference models are unavailable.

### Comparative Advantages

Our approach offers several advantages over existing techniques. By eliminating the need for golden references, our method addresses a fundamental limitation of many existing approaches, making it applicable in supply chain scenarios where trusted references are unavailable. The simulation-based pre-screening reduces the need for expensive physical

testing equipment, lowering the barriers to adoption for organizations with limited resources.

The methodology's scalability allows application to various circuit types with minimal modifications, creating a versatile tool for hardware security assessment. The projected resource utilization below 3% makes implementation feasible even in resource-constrained systems, ensuring that security features don't significantly impact primary functionality. Additionally, the unsupervised learning framework can potentially adapt to new Trojan types without requiring extensive retraining, providing resilience against evolving threats.

These advantages position our approach as a valuable addition to the hardware security toolkit, complementing existing methods rather than replacing them. By addressing specific limitations of current techniques, our methodology fills an important gap in comprehensive hardware security assurance.

Limitations and Challenges

Despite the promising results, several limitations and challenges remain. The simulation environment cannot perfectly capture all physical phenomena that might affect actual hardware, creating potential discrepancies between simulated and real-world behavior. Feature extraction becomes increasingly complex as circuit size and complexity grow, potentially limiting scalability without significant algorithmic improvements.

Unsupervised learning approaches may struggle with very subtle Trojans that cause minimal behavioral changes, creating a lower bound on detection sensitivity. The methodology requires careful parameter tuning for optimal performance, introducing complexity in deployment across diverse circuit types. Real-time implementation faces significant computational challenges, particularly for complex circuits with high-frequency operation.

These limitations highlight the need for continued research and development to bridge the gap between laboratory demonstrations and practical deployment. Addressing these challenges will require interdisciplinary collaboration between hardware designers, security experts, and machine learning specialists.

Broader Implications

The broader implications of this work extend beyond hardware Trojan detection. The methodology could be adapted for general hardware quality assurance and verification, helping identify design flaws or manufacturing defects that manifest as behavioral anomalies. Similar approaches might be applicable to detecting counterfeit or recycled components by identifying characteristic behavioral signatures that differ from authentic parts.

The unsupervised learning framework could benefit other security domains facing similar "needle in a haystack" detection challenges, such as network intrusion detection or malware analysis. The side-channel feature extraction techniques might inform design

practices for creating inherently more secure hardware by minimizing exploitable behavioral patterns.

These implications suggest that the value of this research extends beyond its immediate application to hardware Trojan detection. The fundamental techniques of behavioral characterization and unsupervised pattern recognition have broad applicability across hardware security and potentially other domains.

Final Thoughts

As hardware supply chains become increasingly globalized and complex, the security challenges associated with hardware Trojans will only grow in importance. Our work demonstrates that combining simulation-based side-channel analysis with unsupervised learning techniques offers a promising path forward that addresses the practical limitations of existing approaches.

The proof-of-concept established through this project provides a foundation for future research and development that could eventually lead to robust, practical solutions for hardware security assurance. By continuing to refine these techniques and addressing the identified limitations, we can work toward a future where hardware Trojans can be reliably detected without requiring golden reference models or extensive physical testing, thereby enhancing the security and trustworthiness of critical electronic systems.

The growing reliance on sophisticated electronic systems across critical infrastructure, medical devices, transportation systems, and national security applications makes the need for trustworthy hardware more urgent than ever. Our research contributes a small but significant step toward addressing this challenge, demonstrating that even with minimal resources and relatively simple techniques, effective hardware Trojan detection is achievable. With continued development and refinement, these approaches have the potential to substantially enhance hardware security assurance in an increasingly complex technological landscape.

## References

1. **Bhunia, S., Hsiao, M. S., Banga, M., & Narasimhan, S. (2014). "Hardware Trojan Attacks: Threat Analysis and Countermeasures." Proceedings of the IEEE.**

2. **Huang, Z., Wang, Q., Chen, Y., Jiang, X., School of Computer Science and Technology, Xidian University, Xi'an 710071, China, Graduate School of Media and Governance, Keio University Shonan Fujisawa Campus, Fujisawa 252-0882, Japan, School of Systems Information Science, Future University Hakodate, Hakodate 041-8655, Japan, & Wang, Q. (2020). A survey on Machine Learning against Hardware Trojan Attacks: Recent advances and challenges. In *ACCESS* (Vol. 8, p. 10796) [Journal-article]. https://doi.org/10.1109/ACCESS.2020.2965016**

3. **Rajput, S., Dofe, J., & Danesh, W. (2023). Automating Hardware Trojan Detection Using Unsupervised Learning: A Case Study of FPGA. . https://doi.org/10.1109/isqed57927.2023.10129335**

4. **Xiao, K., Forte, D., & Tehranipoor, M. (2016). "A Novel Approach for Hardware Trojan Detection Using Machine Learning." IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems.**