# TROJAN-GUARD: Hardware Trojans Detection Using GNN in RTL Designs

Kiran Thorat*, Amit Hasan*, Caiwen Ding†, Zhijie Shi*

*Computer Science and Engineering Department, University of Connecticut
{kiran_gautam.thorat, amit.hasan, zshi}@uconn.edu
†Department of Computer Science & Engineering, University of Minnesota
dingc@umn.edu

*Abstract*—Chip manufacturing is a complex process, and to achieve a faster time to market, an increasing number of untrusted third-party tools and designs from around the world are being utilized. The use of these untrusted third party intellectual properties (IPs) and tools increases the risk of adversaries inserting hardware trojans (HTs). The covert nature of HTs poses significant threats to cyberspace, potentially leading to severe consequences for national security, the economy, and personal privacy. Many graph neural network (GNN)-based HT detection methods have been proposed. However, they perform poorly on larger designs because they rely on training with smaller designs. Additionally, these methods do not explore different GNN models that are well-suited for HT detection or provide efficient training and inference processes. We propose a novel framework that generates graph embeddings for large designs (e.g., RISC-V) and incorporates various GNN models tailored for HT detection. Furthermore, our framework introduces domain-specific techniques for efficient training and inference by implementing model quantization. Model quantization reduces the precision of the weights, lowering the computational requirements, enhancing processing speed without significantly affecting detection accuracy. We evaluate our framework using a custom dataset, and our results demonstrate a precision of 98.66% and a recall (true positive rate) of 92.30%, highlighting the effectiveness and efficiency of our approach in detecting hardware trojans in large-scale chip designs.

## I. INTRODUCTION

The widespread presence of hardware chips in modern computing systems highlights its critical importance. The hardware chips are now integrated into diverse applications such as smartphones, IoT, AI, and autonomous vehicles [8], [24], [25]. As these hardware chips collect, analyze, and store the data, security concerns regarding them are increasing. Furthermore, hardware chips are increasingly becoming larger, and more complex which makes hardware chips vulnerable to a plethora of attacks. For example, hardware trojan [29], side-channel attacks [11], [23], information leakage [4], [21], and fault injection attacks [20]. These advanced attacks can threaten the risk of data confidentiality, integrity, and reliability of computing systems. The vulnerabilities come from various sources such as malicious insider tampering, unintentional design errors, weak testing and verification frameworks, and optimization by computer-aided design (CAD) tools that

overlooked security implications [19]. Unlike software bugs which can be mitigated after the software deployment, the vulnerabilities in hardware chips are often non-mitigable after the chips are manufactured. In addition, to reduce the cost, it is preferable to detect and mitigate vulnerabilities as early as possible, for example, in the register transfer logic (RTL) stage.

Driven by time-to-market demands, hardware chip producers often outsource designs to third parties and also use computer-aided design (CAD) tools from different vendors. However, the globalization of the CAD tool industry introduces the risk of hardware trojan (HT) insertion by malicious entities. Ensuring trustworthy hardware chip designs thus requires reliable HT detection methods. Graph Neural Networks (GNNs) have attracted increasing interest because of their performance in graph-based learning tasks [13]. In hardware chip design, netlists can be represented as graphs, making GNNs a promising choice in several chip design tasks including HT detection.

Existing HT detection studies lacks the coverage of different hardware designs in their datasets. Many use the Trusthub [1] benchmark, but only select a subset of designs (e.g., AES, PIC, and RS232). For instance, GNN4HT [3] uses two-stage GNN approach but evaluates only on 21 hardware designs, while [2] generates additional designs for data balancing yet still evaluates on only AES, PIC, and RS232. Moreover, most existing studies use Graph Convolutional Networks (GCNs) [13] without considering other GNN variants that might be better suited for HT detection. For example, HW2VEC [35] relies on a GCN-based model. Additionally, model compression strategies such as sparsity and quantization, which can improve training and inference efficiency, have not been explored in this domain.

To address these gaps, we propose a novel GNN-based framework for HT detection. Our main contributions include:

- We create a new dataset by injecting multiple hardware trojan variants into a diverse collection of RTL designs, giving broad design coverage.
- We design GNN architectures tailored to HT detection by systematically comparing GCN, GraphSAGE, GAT, and GIN layers for netlist-based EDA graphs.
- We develop a 4-bit post-training quantization scheme for GNN-based HT detectors, trimming model size by $8\times$

and speeding up inference while preserving the close to original detection accuracy.

We evaluate proposed framework on the new dataset and achieve a precision of 98.66 % and a recall of 92.30 %, demonstrating the framework's ability to detect HTs in large-scale chip designs. The remainder of this paper is organized as follows: Section II reviews background and related work; Section III explains the trojan-injection procedure and resulting dataset, details the GNN architecture and the 4-bit quantization pipeline; Section IV presents experimental results and discussion.

## II. BACKGROUND AND RELATED WORK

Hardware trojans are malicious modifications made to hardware chips that can threaten national security, the economy, and personal privacy. Detecting HTs early in the design process, especially at the RTL, is important because it is more cost-effective in the early stages. As hardware chips become more complex and are used in various applications such as smartphones, IoT devices, AI systems, and autonomous vehicles, the need for effective HT detection methods has increased. Traditional ML methods have been widely used for HT detection by analyzing structural and functional features extracted from circuit designs. These methods often rely on manually created features from gate-level netlists or other design representations. For example, some studies have used neural networks (NNs) and random forests (RFs) to classify nodes in gate-level netlists, achieving good true positive rates (TPR) and true negative rates (TNR) [15]. Similarly, other research has used models like eXtreme Gradient Boosting (XGBoost) with selected structural features to improve classification performance [22]. However, these traditional ML models face challenges such as the need for extensive feature engineering, difficulties in scaling with increasing circuit complexity, and the lack of a standard feature set for HT detection. GNNs offer a promising alternative to traditional ML models by using the inherent graph structure of circuit designs. GNNs can automatically learn and update node feature representations through message passing and aggregation, reducing the need for manual feature extraction [9]. This ability allows GNNs to generalize better to new circuit designs and adapt to different structural patterns without significant changes. GNNs [13] can perform tasks such as node classification, link prediction, and graph classification, making them suitable for identifying HTs. Several studies have shown the effectiveness of GNNs in HT detection. For instance, some approaches have used Graph Convolutional Networks (GCNs) and Graph Attention Networks (GATs) on netlists, achieving high accuracy in identifying trojan nodes [3], [17], [33], [34]. Other research has used centrality measures from social network analysis, like betweenness centrality (BC) and PageRank (PR), to enhance the detection capabilities of GNN models [10]. These methods show that GNNs can capture both local and global structural properties of circuits, providing a more thorough analysis for HT detection. Despite the progress with GNN-based methods, there are still limitations in current approaches. Many existing studies use limited datasets that do not represent the full range of hardware designs, often testing models on a small set of benchmark circuits from TrustHub [1]. Additionally, most research has focused on specific types of GNN architectures, such as GCNs, without exploring other variants that might offer better performance or efficiency [33]. Furthermore, techniques for model compression, including sparsity and quantization, which could improve the scalability and deployment of GNNs in resource-constrained environments, have not been thoroughly explored in the context of HT detection [3], [17]. To address these challenges, we propose a novel framework that generates graph embeddings for large hardware designs, such as RISC-V cores, and incorporates various GNN models tailored for HT detection. Our framework introduces domain-specific techniques for efficient training and inference by implementing model quantization and inducing sparsity. Model quantization reduces the precision of the weights, lowering computational requirements, while sparsity decreases the number of active connections, enhancing processing speed without significantly affecting detection accuracy. By combining these techniques, our framework achieves both high precision and recall in detecting hardware trojans, making it suitable for large-scale and real-time applications.

## III. METHODOLOGY

### A. Hardware Design

We begin our framework with hardware designs written in RTL Verilog. In chip design, after high-level exploration, the netlist is created using hardware description languages like VHDL or Verilog. We use hardware designs described in RTL Verilog to build the RTL netlist graph. As shown in Fig. 1, we have a full adder example described in an RTL Verilog netlist.

*1) Hardware Trojan Insertion:* Hardware trojans usually have two parts: a trigger and a payload. Both parts are often inserted into low-activity, small modules of the RTL [3].

In the trigger, additional Verilog code triggers signals to activate the payload. As shown in Fig. 2 (a), we create two registers: `Trojan_Counter` (16 bits) and `Trojan_Trigger_Out` in the top module of the RISC-V architecture. If `Trojan_Counter` goes above 100, we set `Trojan_Trigger_Out = 1`.

As shown in Fig. 2 (b), the always block increments `Trojan_Counter` when reset is not active and `IDATA` [6:2] matches certain opcodes (for example, BCC, JALR, RCC).

In the payload, once `Trojan_Trigger_Out` is set, as shown in Fig. 3, the code overrides `IDATA` with `MALICIOUS_INSTR = 32'h0010_0073`, which corresponds to `EBREAK`. This disrupts the CPU (like sleep). We mux the real `IDATA` with the malicious opcode based on `Trojan_Trigger_Out`. This type of trojans is used for denial of service. Other types of the trojans we inserted into our diverse set of hardware designs include information leakage, functionality change, and performance degradation.
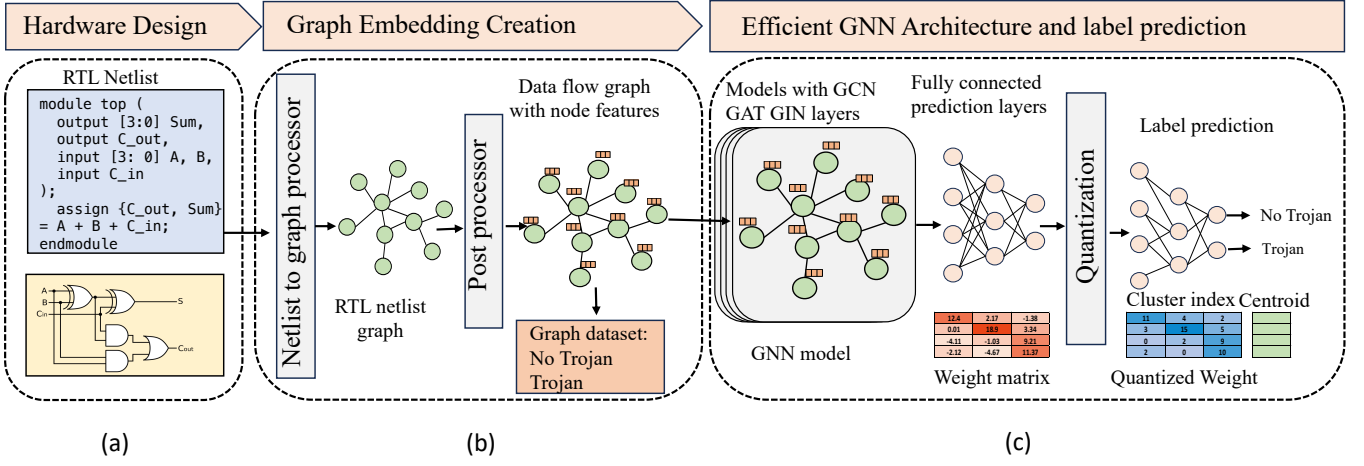
Fig. 1. Framework: (a) RTL hardware design (example, full addder, (b) Graph embedding with node features creation, (c) GNN models and efficient training and inference for prediction.
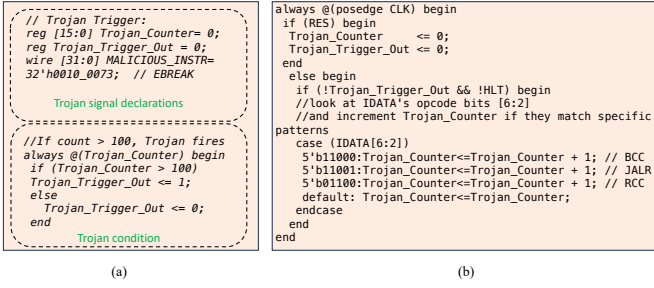


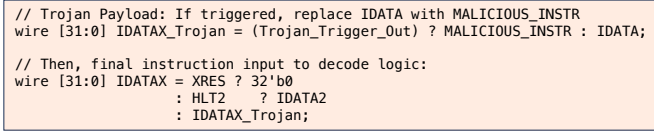Fig. 2. (a) Trojan signal declaration and condition, (b) Trojan trigger

```
// Trojan Payload: If triggered, replace IDATA with MALICIOUS_INSTR
wire [31:0] IDATAX_Trojan = (Trojan_Trigger_Out) ? MALICIOUS_INSTR : IDATA;

// Then, final instruction input to decode logic:
wire [31:0] IDATAX = XRES ? 32'b0
                   : HLT2    ? IDATA2
                   : IDATAX_Trojan;
```

Fig. 3. Trojan payload

## B. Graph Embedding Creation

The hardware design is an RTL netlist which is similar to a graph. As shown in Fig. 1 (b), we use a netlist to graph processor to create RTL netlist graph, which shows the relationships and dependencies among signals and gives a fundamental expression of the code's computational structure. An RTL design may have multiple modules in separate files, so we need a preprocess phase to flatten the design and fix syntax incompatibilities. Next, we use a toolkit called Pyverilog [28], where a parser extracts an abstract syntax tree from Verilog. This tree goes to the data flow analyzer, which builds a tree for each signal in the circuit, with that signal as the root node. We then merge all these signal DFGs into one graph for the entire circuit, and we trim any disconnected sub-graphs and redundant nodes in the merge phase.

As shown in Fig. 1 (b), the final DFG based graph is a created with node features. We construct using a directed graph that shows data dependencies from output signals (root nodes) to input signals (leaf nodes). We define it as $G = (V, E)$, where $V$ is the set of vertices and $E$ is the set of directed edges. Let $V = \{v_1, v_2, \ldots, v_n\}$, where each $v_i$ can be a signal, a constant value, or an operation such as xor, and, concatenation, branch, or branch condition. We define $E = e_{ij}$ for all $i, j$ such that $e_{ij} \in E$ if the value of $v_i$ depends on the value of $v_j$, or if the operation $v_j$ is applied to $v_i$.

## C. Graph Neural Networks

GNNs are used in many applications across various domains by modeling graph structures [13], [14], [26], [27], [36]. They address tasks at the node, edge, or graph level using a consistent two-step procedure: aggregation followed by update. In the aggregation step, the network collects information from a target node's neighbors. Given node representations at the $t$-th layer $h_u^t$ and the set of neighbors of target node $v$, $N_v$, the aggregation function at time $t + 1$ can be written as:

$$a^{(t+1)} = \text{AGGREGATE}^{(t+1)} \left( h_u^{(t)} : u \in N_v \right) \quad (1)$$

Following aggregation, the update step combines the existing node representation with the aggregated information to update the target node's representation

$$h_v^{(t+1)} = \text{UPDATE}^{(t+1)} \left( h^{(t)}, a^{(t+1)} \right) \quad (2)$$

.

Figure 4 shows the process of the "aggregate" and "update" mechanism in GNNs. Figure 4 (a) shows the input graph where the target node $a$ is connected with neighboring nodes $c$, $d$, and $e$. The message-passing happens through the edges that connect the nodes. Figure 4 (b) shows the neighboring nodes pass message (node features) to node $a$. GNN does aggregate operation and node features and updates the node features of node $a$.

Different GNN architectures emerge from the choices made for the functions $\text{AGGREGATE}^{(t+1)}(\cdot)$ and $\text{UPDATE}^{(t+1)}(\cdot)$

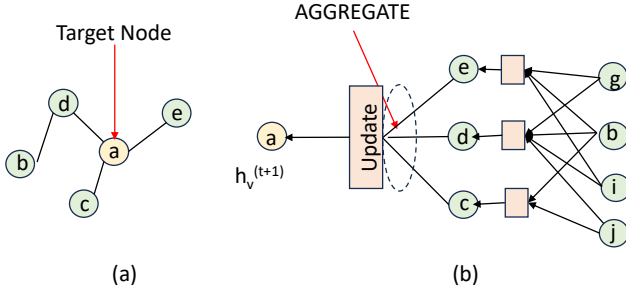[30]. The subsequent sections explore various methods that implement these functions



Fig. 4. (a) Input Graph, (b) GNN: Node aggregate and update

*1) Graph Attention Network (GAT):* GATs [31] replace traditional graph convolutions with masked self-attentional layers. These layers assign different weights to each neighbor during the aggregation of neighboring features. This mechanism allows nodes to focus on more relevant neighbors and helps the model capture complex relationships within the graph. Another advantage is that the importance of neighboring nodes is determined dynamically, without requiring prior knowledge of the entire graph structure.

The updated node representations are obtained as follows:

$$h_v^{(t+1)} = \sigma \left( \sum_{k=1}^{K} \sum_{u \in N_v} \alpha_{vu}^k W_k h_u^{(t)} \right) \quad (3)$$

Here, $\alpha_{vu}^k$ represents the normalized attention coefficient for the $k$-th attention head, and $W_k$ is the corresponding weight matrix. By using attention, the model effectively detects both local and global dependencies.

*2) Graph Convolutional Network (GCN):* A commonly used GNN architecture is the Graph Convolutional Network (GCN) [13], [18]. Inspired by image convolutions, GCNs use convolution filters that operate directly on the graph structure. Unlike images, the neighborhood size of a node in a graph varies. A parameter matrix transforms the node representations obtained from the previous layer, and the transformed representations are weighted according to the graph's adjacency matrix. The update step using GCN is defined as

$$H^{(t+1)} = \sigma \left( \hat{A} H^{(t)} W \right) \quad (4)$$

where $H^{(t+1)}$ denotes the matrix of stacked node representations at layer $t + 1$, $H^{(t)}$ represents the node representations at layer $t$, $\sigma$ is an activation function (e.g., ReLU), $\hat{A}$ is the normalized adjacency matrix, and $W$ is the parameter matrix. GCN uses a shared weight for all edges, making the model relatively simple. For more complex graph structures, this approach may offer less expressiveness.

*3) Graph Isomporhism Network (GIN):* Another variant is GIN [32], which uses multi-layer perceptrons (MLPs) to learn the parameters of the update function. It draws inspiration from the Weisfeiler-Lehman (WL) graph isomorphism test

[16], which determines the similarity between two graphs. The node update is defined by

$$h_v^{(t+1)} = \text{MLP}^{(t+1)} \left( \left(1 + \epsilon^{(t+1)}\right) \cdot h_v^{(t)} + \sum_{u \in N(v)} h_u^{(t)} \right) \quad (5)$$

where $\epsilon^{(t+1)}$ is a learnable parameter. GIN uses a simple sum operator to aggregate features, which makes it computationally efficient. The learnable parameter $\epsilon$ allows the model to adapt to various graph structures and effectively capture graph-level features, making GIN a common choice for graph-level tasks.

*D. Efficient training and inference*

Quantization is a commonly used technique to reduce the precision of neural network weights and activations. Thus, the model consumes less memory and with reduced precision, computations become faster. Given a full-precision floating-point weight matrix $W \in \mathbb{R}^{m \times n}$, we create a quantized version $\hat{W}$ by mapping each element to a lower-bit representation:

$$\hat{W} = \text{round}(W/S) \times S, \quad (6)$$

where $S$ is a scaling factor. $S$ adjusts the range of the quantized values to match the distribution of the original weights by determining an appropriate step size for rounding, minimizing the error introduced by quantization while ensuring that the values remain representative of the full precision data.

For hardware Trojan detection using GNNs (as shown in Figure 1 (a)), quantization makes our workflow more efficient. We apply 4-bit quantization using the `bitsandbytes` library [7], which reduces our 32-bit model weights to 4-bit precision while keeping accuracy nearly the same. This process reduces the memory footprint, speeds up inference, makes it easier to deploy the model on hardware with limited resources. In 4-bit quantization, weights are restricted to $2^4 = 16$ levels and we can cut storage needs by a factor of 8 compared to 32-bit floating-point weights [6], [12].

In our work, applying 4-bit quantization helps process and store the model efficiently while keeping accuracy high. Given a graph representation with node features $X$ and an adjacency matrix $A$, a graph convolution operation is performed as follows:

$$H^{(l+1)} = \sigma(A H^{(l)} W^{(l)}), \quad (7)$$

where $W^{(l)}$ is the weight matrix at layer $l$. After quantization, this becomes:

$$\hat{H}^{(l+1)} = \sigma(A \hat{H}^{(l)} \hat{W}^{(l)}), \quad (8)$$

This allows faster calculations due to the lower precision. The balance between accuracy and efficiency makes low-bit quantization a viable option for faster inference.

## IV. EVALUATION

*A. Experimental Setup*

In evaluating our framework, we utilize an AMD EPYC 7763 64-Core Processor and Ubuntu 22.04.5 LTS system. We conduct experiments on NVIDIA RTX A6000 graphics card

| Hyperparameter | Value |
|---|---|
| Learning Rate | 0.001 |
| Epochs | 200 |
| Hidden Units | 200 |
| Dropout Rate | 0.5 |
| Batch Size | 4 |
| Pooling Ratio | 0.8 |
| Embedding Dimension | 2 |

equipped with 48 GB of GDDR6 memory, and CUDA version 12.6.

Table I presents the hyperparameters used in training our GNN models. We set the learning rate to 0.001 and trained the model for 200 epochs with 200 hidden units and a dropout rate of 0.5 to prevent overfitting. A batch size of 4 was used, and we employed a pooling ratio of 0.8. Additionally, we evaluated the model every 10 epochs and used an embedding dimension of 2.

*B. Dataset*

We create a dataset from the Trusthub benchmark [1] and the open-source GitHub repository [5]. The dataset consists of 51 different RTL designs with inserted trojans that cover denial of service, information leakage, functionality change, and performance degradation. As shown in Figure 5, the dataset includes 40 designs with trojans (class 1) and 11 designs without trojans (class 0). We split the dataset into 70% for training, 10% for validation, and 20% for testing the GNN model. The training set consists of 27 class 1 designs and 8 class 0 designs. The validation set includes 5 class 1 designs and 1 class 0 design. The testing set contains 9 class 1 designs and 2 class 0 designs. This distribution ensures that the dataset is diverse and balanced for training and testing the GNN models.

TABLE II
GRAPH STATS FOR SELECTED DESIGNS

| Design | Type | Nodes | Edges | Time (s) |
|---|---|---|---|---|
| AES-T1900 | TjIn | 20446 | 25036 | 24.5754 |
| Darkriscv-T200 | TjFree | 1587 | 2175 | 1.5074 |
| Darkriscv-T100 | TjIn | 1538 | 2098 | 1.2373 |
| AES-T1500 | TjIn | 21100 | 25845 | 24.9904 |
| PIC16F84-T300 | TjIn | 2636 | 3694 | 3.2588 |

*C. Framework evaluation*

We evaluate our framework using a diverse dataset to demonstrate its effectiveness in detecting hardware trojans. To address the class imbalance, where class 0 (designs without trojans) is underrepresented, we employ a 5-fold cross-validation method in our experiments. By averaging the results over five runs, we follow standard practices to prevent overfitting and ensure that our reported metrics accurately reflect the model's performance. We tested various GNN architectures with different layer configurations, including 2-layer, 3-layer,
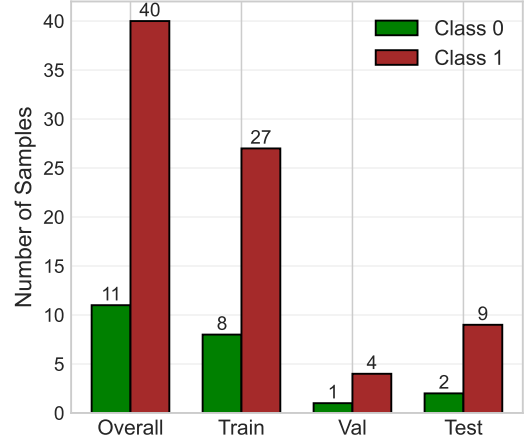


Fig. 5. Dataset distribution

and 5-layer models for GCN, GAT, and GIN. Figure 6 shows the evaluation results, where the x-axis represents the different GNN types and their respective layers, and the y-axis shows the percentage values of each performance metric.

Figure 6 (a) shows the precision which measures the proportion of correctly identified trojans out of all instances flagged by the model. the 2-layer GCN model achieves the highest precision of 98.66%, indicating it is highly effective in minimizing false positives. Additionally, the 2-layer GIN model also demonstrates high precision at 98.5%. As the number of layers increases, precision decreases in each model. For example, GCN precision drops from 98.66% with 2 layers (GCN-2) to 93.9% with 3 layers (GCN-3), and further to 84.8% with 5 layers (GCN-5). Similarly, GAT and GIN models show a decline in precision with deeper layers, highlighting that deeper networks tend to have lower precision due to increased complexity and potential overfitting. Accuracy measures the overall correctness of the model by calculating the proportion of true results (both true positives and true negatives) among the total number of cases examined. Figure 6 (b) shows that the 2-layer GCN model also achieves the highest accuracy of 92.4%. The 5-layer GIN model follows with an accuracy of 89.8%. As the number of layers increases, accuracy decreases across GCN and GAT types. For instance, GCN-3 and GCN-5 achieve accuracies of 81.33% and 69.4%, respectively. The GIN-5 shows 89.8% which due to potential overfitting. This decline indicates that deeper GNN models may struggle to maintain performance as model complexity increases. Figure 6 (c) shows F1-score is the harmonic mean of precision and recall, providing a balance between the two metrics. The 2-layer GCN model achieves the highest F1-score of 94.6%, indicating a strong balance between precision and recall. The 5-layer GIN model also achieves a notable F1-score of 91.9%. As the number of layers increases, the F1-score decreases for all GNN types. For example, GCN-3 and GCN-5 have F1-scores of 87.3% and 73.5%, respectively. This trend highlights that deeper models may compromise the
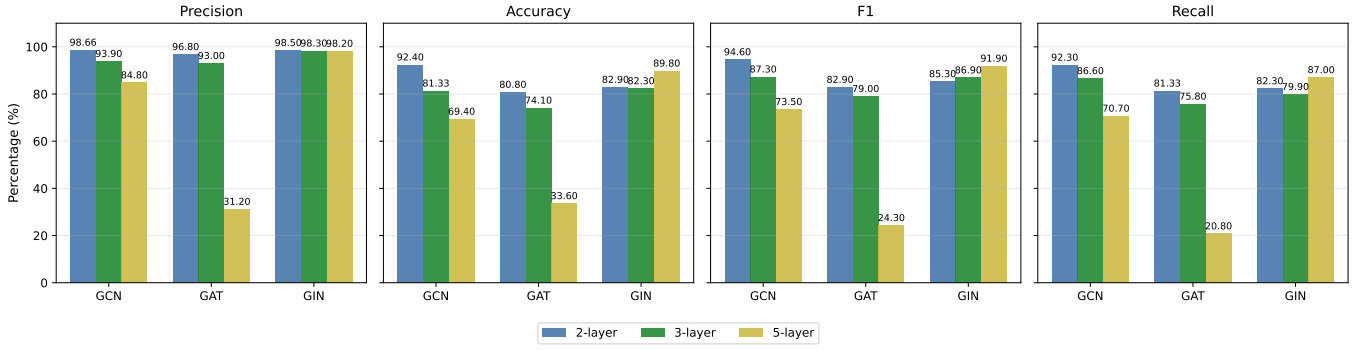
Fig. 6. Evaluation results in comparison with different models: (a) Accuracy, (b) Precision , (c) F1-Scores, (d) Recall

balance between correctly identifying trojans and minimizing false positives. Figure 6 (d) shows Recall, also known as the true positive rate, measures the proportion of actual trojans correctly identified by the model. The 2-layer GCN model achieves the highest recall of 92.30%, ensuring that most trojans are detected. The 5-layer GIN model also shows a strong recall of 87.0%. Increasing the number of layers leads to a decrease in recall, with GCN-3 and GCN-5 showing recall rates of 86.6% and 70.7%, respectively. Similar decreases are observed in GAT and GIN models, indicating that deeper layers reduce the model's ability to identify all true trojans.

Considering all metrics accuracy, precision, recall, and F1-score, GCN-2 consistently outperforms the other models. GCN-2 achieves the highest precision (98.66%), accuracy (92.4%), recall (92.30%), and F1-score (94.6%), demonstrating its superior capability in effectively detecting hardware trojans while maintaining a low rate of false positives. Additionally, the 5-layer GIN model shows strong performance in accuracy and F1-score, making it a competitive alternative. The decline in performance with increased layers across all GNN types suggests that shallower models are more suitable for this task, likely due to better generalization and reduced complexity in training.

### D. Memory Utilization

TABLE III
GPU MEMORY USAGE BY DIFFERENT GNN MODELS

| GNN Type | Layers | GPU Memory (MiB) |
| --- | --- | --- |
| GCN | 2 | 1,688 |
| GCN | 3 | 1,880 |
| GCN | 5 | 2,322 |
| GIN | 2 | 2,742 |
| GIN | 3 | 3,472 |
| GIN | 5 | 4,940 |
| GAT | 2 | 3,310 |
| GAT | 3 | 4,000 |
| GAT | 5 | 5,296 |

Table III presents the GPU memory usage for different GNN models with varying layer configurations. The results show that GPU memory consumption increases with the number of layers for all GNN types. Specifically, the GAT model has

the highest memory usage, starting at 3,310 MiB for 2 layers and escalating to 4,000 MiB with 3 layers and 5,296 MiB with 5 layers. the GCN model uses 1,688 MiB with 2 layers, which increases to 1,880 MiB with 3 layers and 2,322 MiB with 5 layers. The GIN model starts at 2,742 MiB for 2 layers and rises to 3,472 MiB for 3 layers and 4,940 MiB for 5 layers. These findings indicate that deeper GNN architectures require significantly more GPU memory, which is an important consideration for scalability and deployment in environments with limited resources. In subsection Quantized Model Results IV-F, we discuss memory saving.

### E. Comparison with SOTAs

To demonstrate our framework's performance, we compared it with two GNN-based state-of-the-art (SOTA) methods: SOTA1 [33] and SOTA2 [2]. For a fair comparison, we set the same hyperparameters (learning rate: 0.001, epochs: 200, hidden units: 200) and used the same GNN layers (i.e., a 2-layer GNN model). In these comparison experiments, we did not employ the validation dataset split and k-fold validation technique for fair comparison. This might lead to overfitted models and we tried our best to provide non-overfit results for a fair comparison with SOTAs. As shown in Fig. 7, we present the metrics of interest on the y-axis and the designs (AES, RS232, and PIC) on the x-axis. Note that these papers only used a subset of our dataset, so we created the same dataset as the SOTAs to show the results. We show precision comparison in Fig. 7 (a) with the SOTAs. Our framework outperforms the SOTAs with 98.1% precision for the AES and for RS232 is slightly lower, which can be attributed to fewer training samples for the RS232 designs in our dataset. The PIC dataset is too small models give the overfitted for multiple ran result as evidenced by SOTA 2 and our results. Fig. 7 (b) shows recall comparison. As we can see, our framework achieves the best recall 93.3% for AES designs. For other designs model is overfitting the results. Fig. 7 (c) shows F1 scores, where our scores are better than SOTA1 for AES and RS232 dataset which non a overfitted model. Overall, our framework outperforms the SOTA methods in most of the evaluation metrics.
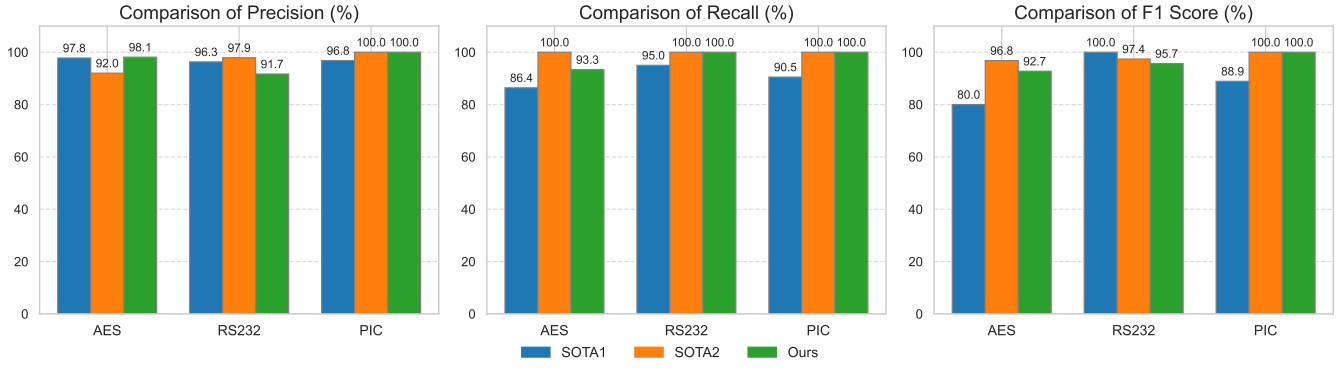
Fig. 7. Evaluation results in comparison with SOTAs: (a) Precision, (b) Recall, (c) F1-Scores; SOTA1= [33] and SOTA2= [2]

## F. Quantized Model Results

To enhance the efficiency of our framework for training and inference, we quantize our model. The quantized model uses significantly less GPU memory to process the GNN. We apply 4-bit quantization using the `bitsandbytes` library [7], which reduces our 32-bit model weights to 4-bit precision while maintaining similar accuracy. This process decreases the memory footprint, speeds up inference, and facilitates deployment on hardware with limited resources. In 4-bit quantization, weights are limited to $2^4 = 16$ levels, reducing storage needs by a factor of 8 compared to 32-bit floating-point weights [6], [12]. To demonstrate the framework's performance with quantization, we selected the two best-performing models from Figure 6: the 2-layer GCN and the 5-layer GIN models. Figure 8 (a) shows the quantized 2-layer GCN model. In the quantized version, precision is 92.2%, compared to 98.66% in the non-quantized model. Similarly, accuracy decreases from 92.4% to 88.3%. These slight reductions in performance are outweighed by the substantial decrease in GPU memory usage. Figure 8 (b) shows the quantized 5-layer GIN model, which uses 4,940 MiB of GPU memory as shown in Table III. In the quantized version, precision is 96.3%, compared to 98.2% in the non-quantized model. Although there is a minor decrease in precision, the quantized model benefits from significantly reduced GPU memory consumption. These results indicate that quantization slightly lowers some performance metrics but offers substantial benefits in terms of memory efficiency. This trade-off makes our framework suitable for scaling to larger graphs with complex features, enabling deployment in environments with limited GPU resources without a significant loss in detection performance.

## V. CONCLUSION

In this paper, we developed a novel framework for detecting hardware trojans in large-scale chip designs, such as RISC-V cores. Our framework includes graph extraction to enhance coverage and explores various GNN models to identify the most effective ones for HT detection. To improve training and inference efficiency, we applied 4-bit quantization to the models, significantly reducing GPU memory usage while
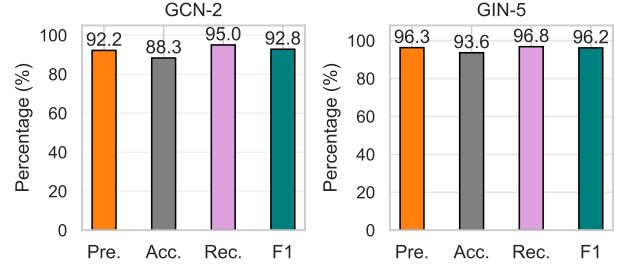


Fig. 8. Quantized GNN model results with: (a) GCN-2 , (b) GIN-5

maintaining high accuracy. We evaluated our framework using a custom dataset, achieving a precision of 98.66% and a recall of 92.30%. These results demonstrate the effectiveness and efficiency of our approach in accurately identifying hardware trojans in complex chip designs, making it suitable for deployment in resource-constrained environments.

## REFERENCES

[1] Chip-level trojan benchmarks. https://trust-hub.org/#/benchmarks/chip-level-trojan.

[2] Anindita Chattopadhyay, Siddharth Bisariya, Anantram Patel, Sowmya Sunkara, and Vijay Kumar Sutrakar. Deep learning based graph neural network technique for hardware trojan detection at register transfer level. In *2024 IEEE 4th International Conference on VLSI Systems, Architecture, Technology and Applications (VLSI SATA)*, pages 1–5, 2024.

[3] Lihan Chen, Chen Dong, Qiaowen Wu, Ximeng Liu, Xiaodong Guo, Zhenyi Chen, Hao Zhang, and Yang Yang. Gnn4ht: A two-stage gnn based approach for hardware trojan multifunctional classification. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2024.

[4] Gustavo K Contreras, Adib Nahiyan, Swarup Bhunia, Domenic Forte, and Mark Tehranipoor. Security vulnerability analysis of design-for-test exploits for asset protection in socs. In *2017 22nd Asia and South Pacific Design Automation Conference (ASP-DAC)*, pages 617–622. IEEE, 2017.

[5] Darklife. Darkriscv. https://github.com/darklife/darkriscv.

[6] Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *ArXiv*, abs/2305.14314, 2023.

[7] Tim Dettmers, Ruslan Svirschevski, Vage Egiazarian, Denis Kuznedelev, Elias Frantar, Saleh Ashkboos, Alexander Borzunov, Torsten Hoefler,

and Dan Alistarh. Spqr: A sparse-quantized representation for near-lossless llm weight compression. *arXiv preprint arXiv:2306.03078*, 2023.

[8] Weimin Fu, Kaichen Yang, Raj Gautam Dutta, Xiaolong Guo, and Gang Qu. Llm4sechw: Leveraging domain-specific large language model for hardware debugging. In *2023 Asian Hardware Oriented Security and Trust Symposium (AsianHOST)*, pages 1–6. IEEE, 2023.

[9] Kento Hasegawa, Kazuki Yamashita, Seira Hidano, Kazuhide Fukushima, Kazuo Hashimoto, and Nozomu Togawa. Node-wise hardware trojan detection based on graph learning. *IEEE Transactions on Computers*, 2023.

[10] Mona Hashemi, Amirabbas Momeni, A Pashrashid, and Siamak Mohammadi. Graph centrality algorithms for hardware trojan detection at gate-level netlists. *International Journal of Engineering*, 35(7):1375–1387, 2022.

[11] Zecheng He and Ruby B Lee. How secure is your cache against side-channel attacks? In *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 341–353, 2017.

[12] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew G. Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2704–2713, 2017.

[13] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2017.

[14] Kishor Kunal, Tonmoy Dhar, Meghna Madhusudan, Jitesh Poojary, Arvind K Sharma, Wenbin Xu, Steven M Burns, Jiang Hu, Ramesh Harjani, and Sachin S Sapatnekar. Gnn-based hierarchical annotation for analog circuits. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 42(9):2801–2814, 2023.

[15] Tatsuki Kurihara, Kento Hasegawa, and Nozomu Togawa. Evaluation on hardware-trojan detection at gate-level ip cores utilizing machine learning methods. In *2020 IEEE 26th International Symposium on On-Line Testing and Robust System Design (IOLTS)*, pages 1–4. IEEE, 2020.

[16] Andrei Leman and Boris Weisfeiler. A reduction of a graph to a canonical form and an algebra arising during this reduction. *Nauchno-Technicheskaya Informatsiya*, 2(9):12–16, 1968.

[17] Peijun Ma, Ge Shang, Hongjin Liu, Jiangyi Shi, Weitao Pan, Yan Zhang, and Yue Hao. Gnn-based hardware trojan detection at register transfer level leveraging multiple-category features. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 2024.

[18] Tengfei Ma, Patrick Ferber, Siyu Huo, Jie Chen, and Michael Katz. Online planner selection with graph neural networks and adaptive scheduling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5077–5084, 2020.

[19] Prabhat Mishra, Mark Tehranipoor, and Swarup Bhunia. Security and trust vulnerabilities in third-party ips. *Hardware IP Security and Trust*, pages 3–14, 2017.

[20] Adib Nahiyan, Farimah Farahmandi, Prabhat Mishra, Domenic Forte, and Mark Tehranipoor. Security-aware fsm design flow for identifying and mitigating vulnerabilities to fault attacks. *IEEE Transactions on Computer-aided design of integrated circuits and systems*, 38(6):1003–1016, 2018.

[21] Adib Nahiyan, Kan Xiao, Kun Yang, Yeir Jin, Domenic Forte, and Mark Tehranipoor. Avfsm: A framework for identifying and mitigating vulnerabilities in fsms. In *Proceedings of the 53rd Annual Design Automation Conference*, pages 1–6, 2016.

[22] Ryotaro Negishi, Tatsuki Kurihara, and Nozomu Togawa. Hardware-trojan detection at gate-level netlists using a gradient boosting decision tree model and its extension using trojan probability propagation. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, 107(1):63–74, 2024.

[23] Nitin Pundir, Jungmin Park, Farimah Farahmandi, and Mark Tehranipoor. Power side-channel leakage assessment framework at register-transfer level. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 30(9):1207–1218, 2022.

[24] Sandip Ray, Eric Peeters, Mark M Tehranipoor, and Swarup Bhunia. System-on-chip platform security assurance: Architecture and validation. *Proceedings of the IEEE*, 106(1):21–37, 2017.

[25] Dipayan Saha, Shams Tarek, Katayoon Yahyaei, Sujan Kumar Saha, Jingbo Zhou, Mark Tehranipoor, and Farimah Farahmandi. Llm for soc security: A paradigm shift. *IEEE Access*, 2024.

[26] Aditya Shahane, Saripilli Swapna Manjiri, Ankesh Jain, and Sandeep Kumar. Graph of circuits with gnn for exploring the optimal design

space. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 6014–6025. Curran Associates, Inc., 2023.

[27] Zixing Song, Yifei Zhang, and Irwin King. No change, no gain: empowering graph neural networks with expected model change maximization for active learning. *Advances in Neural Information Processing Systems*, 36, 2024.

[28] Shinya Takamaeda-Yamazaki. Pyverilog: A python-based hardware design processing toolkit for verilog hdl. In *International Workshop on Applied Reconfigurable Computing*, 2015.

[29] Mohammad Tehranipoor and Farinaz Koushanfar. A survey of hardware trojan taxonomy and detection. *IEEE design & test of computers*, 27(1):10–25, 2010.

[30] Jana Vatter, Ruben Mayer, Hans-Arno Jacobsen, Horst Samulowitz, and Michael Katz. Choosing a classical planner with graph neural networks. *arXiv preprint arXiv:2402.04874*, 2024.

[31] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.

[32] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.

[33] Rozhin Yasaei, Luke Chen, Shih-Yuan Yu, and Mohammad Abdullah Al Faruque. Hardware trojan detection using graph neural networks. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2022.

[34] Rozhin Yasaei, Shih-Yuan Yu, and Mohammad Abdullah Al Faruque. Gnn4tj: Graph neural networks for hardware trojan detection at register transfer level. In *2021 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 1504–1509. IEEE, 2021.

[35] Shih-Yuan Yu, Rozhin Yasaei, Qingrong Zhou, Tommy Nguyen, and Mohammad Abdullah Al Faruque. Hw2vec: A graph learning tool for automating hardware security, 2021.

[36] Xinqi Zhang, Jihao Shi, Junjie Li, Xinyan Huang, Fu Xiao, Qiliang Wang, Asif Sohail Usmani, and Guoming Chen. Hydrogen jet and diffusion modeling by physics-informed graph neural network. *Renewable and Sustainable Energy Reviews*, 207:114898, 2025.