# Automated Physical Design Watermarking Leveraging Graph Neural Networks

Ruisi Zhang
ruz032@ucsd.edu
UC San Diego
La Jolla, CA, USA

Rachel Selina Rajarathnam
rachelselina.r@utexas.edu
UT Austin
Austin, TX, USA

David Z. Pan
dpan@ece.utexas.edu
UT Austin
Austin, TX, USA

Farinaz Koushanfar
farinaz@ucsd.edu
UC San Diego
La Jolla, CA, USA

## Abstract

This paper presents AutoMarks, an automated and transferable watermarking framework that leverages graph neural networks to reduce the watermark search overheads during the placement stage. AutoMarks's novel automated watermark search is accomplished by (i) constructing novel graph and node features with physical, semantic, and design constraint-aware representation; (ii) designing a data-efficient sampling strategy for watermarking fidelity label collection; and (iii) leveraging a graph neural network to learn the connectivity between cells and predict the watermarking fidelity on unseen layouts. Extensive evaluations on ISPD'15 and ISPD'19 benchmarks demonstrate that our proposed automated methodology: (i) is capable of finding quality-preserving watermarks in a short time; and (ii) is transferable across various designs, i.e., AutoMarks trained on one layout is generalizable to other benchmark circuits. AutoMarks is also resilient against potential watermark removal and forging attacks.

## 1 Introduction

During the integrated circuit (IC) design flow, the physical design stage [13, 37] transforms the design's logic netlist to a physical layout for fabrication. The placement [18, 28] and routing [8, 17] of the circuit components are optimized while ensuring various design constraints and desired functionalities. However, the IC layouts are susceptible to various security threats from the supply chain [34], such as intellectual property (IP) theft, counterfeiting, and unauthorized overproduction. Watermarking [12, 39, 40] protects the physical design IP by encoding invisible, confidential, and robust signatures onto the IC layouts. These signatures help the design company to claim ownership of the physical design layouts and track the distribution of the design outcomes.

Prior physical design watermarking solutions fall into two categories [40]: (i) invasive watermarking and (ii) constraint-based watermarking. Invasive watermarking [5, 30, 35, 39] adds redundant wires or cells into the layouts as watermarks. Nevertheless, the encoded signatures can be easily forged if the adversary has prior knowledge of the watermarking scheme, making it less secure for real-world applications.

Constraint-based watermarking [7, 12, 40], on the other hand, adds additional positional constraints to watermark selected cells during the placement stage. However, selecting the watermark cells without considering the modern IC design constraints, such as macros and fence regions, would significantly deteriorate quality metrics. More recent work introduced ICMarks [40] to encode a watermark region of cells onto modern IC design without impacting the overall quality metrics. It aims to constrain watermark cells to be in the watermark region and the remaining cells to be out. However, searching for the ideal region constraints that adhere to design constraints and maintain watermarking fidelity consumes significant time.

This paper presents AutoMarks that leverages graph neural networks (GNNs) to automate the search for the best region constraints to watermark, resulting in an efficient and transferable watermarking framework for physical design. By representing the layout as a graph with node features capturing physical, semantic, and design constraint information, we leverage GNNs to learn the overall quality degradation incurred by watermarking a region centered on a specific node.

AutoMarks encompasses three key stages: (i) *watermark search*, (ii) *watermark insertion*, and (iii) *watermark extraction*. During *watermark search*, AutoMarks employs GNN to identify the ideal region on the layout to encode the watermark. The *watermark insertion* stage encodes the watermark to the cells within the watermark region during placement, similar to [40]. The placement is then routed to obtain the watermarked layout. Finally, *watermark extraction* employs reverse-engineering approaches [3, 31] to acquire the cell connections and the positions from the layout. The design company can prove its ownership by decoding watermark cell positions and determining the percentage of cells within the watermark region.

In brief, our contributions are summarized as follows:

- We introduce AutoMarks, a transferable watermarking framework for physical design that leverages graph neural networks to automate the *watermark search*.
- AutoMarks preserves the watermarking fidelity by representing the design layout as a graph and employing GNN to predict the quality degradation due to watermarking at different locations in the layout.
- Experiments on the ISPD'15 [6] and ISPD'19 [20] benchmarks demonstrate AutoMarks's transferability toward unseen layouts while maintaining the watermarking fidelity; It reduces the search time by 50% for large designs ($\geq 500k$ cells) compared to ICMakrs[40].
- Evaluations of AutoMarks under various watermark removal and forging attacks showcase AutoMarks's resiliency; the adversary cannot remove the signatures without significant quality degradations.

## 2 Background and Related Work

## 2.1 Graph Neural Networks for Physical Design

A graph $G = (V, E)$ consists of a set of nodes $V = \{v_1, v_2, ..., v_{|V|}\}$ and a set of edges $E = \{e_1, e_2, ..., e_{|E|}\}$. Each node $v_i \in V$ has a $k$-dimensional feature vector $\mathbf{h}_i \in R^k$. The goal of a GNN [41] is to learn a function $g$ that maps the feature embedding $\mathbf{h}_i$ of node $v_i$ to a new embedding vector $\hat{\mathbf{h}_i} \in R^h$, capturing both local and global information. In a multi-layer GNN, this mapping is performed iteratively, allowing nodes to update their feature embeddings using information from their neighbors via message passing. For each node $v_i$, message passing involves receiving feature embeddings from its neighbors $j \in N_i$ and aggregating the messages using customizable functions to obtain an updated representation $\hat{\mathbf{h}_i}$. The computations are formulated in Equation 1, where $f$, $g$, and $\oplus$ are customizable functions, such as convolution. Here, $\mathbf{h}_v^{(l)}$ and $\mathbf{h}_u^{(l)}$ are the node features at layer $l$, $\mathcal{N}(v)$ is the set of neighbors for node $v$, and $\mathbf{h}_v^{(l+1)}$ is the updated feature of node $v$ at layer $(l + 1)$.

$$\mathbf{h}_v^{(l+1)} = g(\mathbf{h}_v^{(l)} \underset{u \in \mathcal{N}(v)}{\oplus} f(\mathbf{h}_u^{(l)}, \mathbf{h}_v^{(l)})) \tag{1}$$

In physical design, graph neural networks (GNNs) serve as a promising approach during various stages of physical design, such as logic synthesis, floorplanning, partitioning, placement, and routing to improve the overall power-performance-area (PPA) metrics and speed up the chip design process [10, 21–24, 33, 38]. During the placement stage of physical design, GNNs have been employed to optimize for macro locations, overall timing, and power [11, 25, 26, 28]. GNNs have also been employed for hardware reliability, security, and reverse engineering of gate-level netlists at the physical design stage [2, 4].

## 2.2 Physical Design Watermarking

Physical design watermarking encodes imperceptible and unique identifiers onto the physical design layouts to protect the design company's intellectual property (IP). Invasive watermarking schemes add redundant cells or nets onto the layout as watermarks [5, 30, 35, 39], but can be easily counterfeited if the attacker has prior knowledge of the watermarking scheme. Constraint-based watermarking schemes [7, 12, 14, 29] modify the cell row index and cell spacing as watermarks during the placement stage. Nevertheless, the constraint-based watermarking schemes can significantly degrade layout quality if the watermarks are selected without considering modern design constraints such as macros and fence regions.

The recent constraint-based watermarking work, ICMarks [40], employs a scoring function to evaluate each subregion to ensure (i) sufficient cells to accommodate the signature bits; (ii) small total cell area within the watermark region for cell maneuverability; and (iii) minimal cell area overlap on the region boundary to reduce the impact of cell displacement on layout performance. ICMarks avoids regions that overlap with macros and fence regions and selects the region with the lowest score to encode the watermark. Then, the watermark insertion encodes the region of cells during placement, which co-optimizes the placement objective to ensure that only the watermark cells are in the watermark region.

## 3 Motivation and Problem Formulation

**Scenario:** The IC design company establishes ownership proof for the final physical design layout by encoding watermarks during the placement stage using AUTOMARKS. The watermarked layout is sent along the supply chain for manufacturing and testing. To prove ownership, the design company reverse-engineers the cell locations from the layout [3, 31] and compares the encoded watermarks with the decoded ones.

**Watermarking Criteria:** An ideal watermarking framework shall meet the following criteria:
- *Fidelity*: The watermarked layout does not impact the functionality and the performance of the IC.
- *Effectiveness*: The encoded watermarks can be successfully extracted along the supply chain for ownership proof.
- *Robustness*: The watermarks can withstand various removal and forging attacks and remain detectable after transformations.

**Threats:** We consider the adversary to be a member of the supply chain, e.g., a fabrication or testing company, with access to the layout and general knowledge of watermarking algorithms. However, he/she cannot access the specific watermarking signatures or hyperparameters used by the design company. Potential threats to the watermarked layout include:
- *Watermark Removal Attacks*: The adversary removes the watermark by perturbing the placement solutions.
- *Watermark Forging Attacks*: The adversary forges the watermark by counterfeiting a different set of watermarks.

## 4 Methodology

This section introduces AUTOMARKS and the stages in its pipeline: *Watermark Search*, *Watermark Insertion*, and *Watermark Extraction*, as depicted in Fig. 1. *Watermark Search* identifies potential regions and cells to encode the signature, and *Watermark Insertion* encodes watermarks onto the IC layout without degrading the quality metrics. Finally, *Watermark Extraction* decodes the watermark signature to verify ownership.

### 4.1 Watermark Search

The objective of the *Watermark Search* is to identify a region of cells that incurs minimal quality degradation after the signature is inserted. It starts by constructing a graph representing the netlist while preserving location, semantic, and design constraint features. Then, AUTOMARKS collects the quality degradation labels on the training design and trains a graph neural network to learn the watermarking fidelity on the training nodes. The pre-trained GNN inferences on an unseen layout and the minimum-scored region node is selected as the target cell to watermark.

*4.1.1 Graph Construction.* The design instances, e.g., standard cells and macros, are connected to others by wired nets through the cell pins. As in Fig. 1(A-1), the design is converted into a directed homogeneous graph $G = (V, E)$, where the source nodes are the cells sending signals and the destination nodes represent the cells loading signals. The edges $E$ are the nets connecting the cells.

For each node in the graph, we construct an 8-dimension feature $D$ as listed in Table 1. The first four dimensions describe the physical information, including cell locations from initial placement and the cell size. To capture the semantic information, we extract the cell name embeddings [27] from S-BERT [32] and use Principal
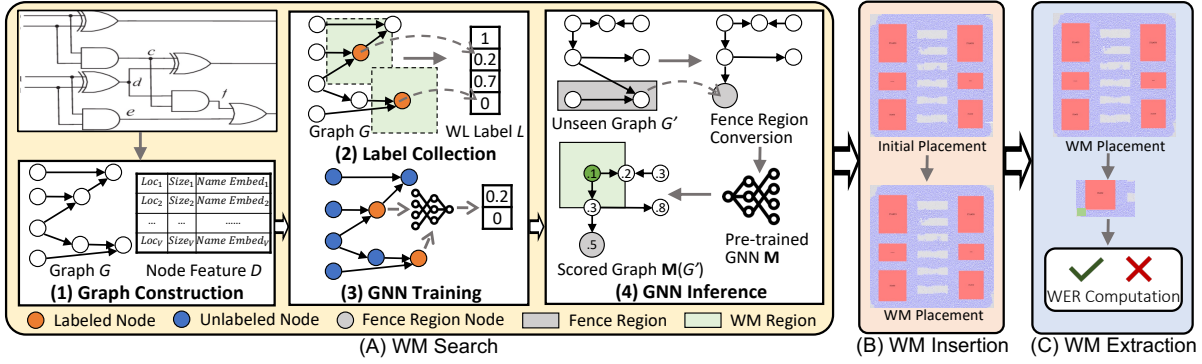
**Figure 1: AutoMarks flow:** The *Watermark Search* leverages GNN to identify the region and cells to watermark. Then, the *Watermark Insertion* encodes the selected watermark region of cells on the layout during the placement stage. The *Watermark Extraction* decodes the watermark and compares it with the encoded ones for ownership proof.

Component Analysis (PCA) [1] to reduce them to four dimensions to form the last four features of $D$.

**Table 1: Node Feature Construction**

| Type | Dim | Description |
|---|---|---|
| Cell Location | 2 | from initial placement $(x, y)$ |
| Cell Size | 2 | cell size (width, height) |
| Cell Name | 4 | S-BERT semantics [32] |

Modern IC designs often have additional design constraints, such as macros and fence regions. We incorporate these design constraints into the graph modeling as follows:

• *Macros*: The macros are treated the same as other standard cells in the graph. They are distinguished from the standard cells by different cell name embeddings.

• *Fence Regions*: As depicted in Fig. 1(A-4), the cells in the fence region are characterized as a single node in the graph, with cell location/size set to the fence region location/size. The name of the fence region cell is assigned after the macros.

*4.1.2* **Label Collection**. We collect the watermarking quality degradation labels on the ispd19test6 [20] design. We chose this design as the training graph for two reasons: (i) the design is of medium size with ~ 180k nodes. It ensures a complex graph structure for the GNN to learn, and the label collection takes one minute per node; (ii) the design has macros, enabling the GNN to learn the graph structure with design constraints. The watermark region has a fixed size of $N \times row\_height$, where $N = 10$.

Watermarking on every node [1] and acquiring its corresponding watermarking quality degradation would take several months for the ispd19test6 design. Therefore, we grid the layout with the size of the watermark region and randomly sample one node from each grid as the training node. The subsampling results in approximately $3.8k$ training nodes and takes only three days to obtain all labels. For the $i$-th node, we denote its improvement over the initial placement wirelengh $WL_{INIT}$ as $\hat{L}_i = \frac{WL_i}{WL_{INIT}}$.

After acquiring the labels of each training node, $\hat{L}_i$ is transformed into the range of $[0, 1]$ for better GNN learning. As specified in Equation 2, if $\hat{L}_i < 1$, it indicates watermarking on the node will not introduce any performance degradation. Such nodes are scored as 0, indicating an ideal position for watermark encoding. If $\hat{L}_i > \beta$,

---

[1]For ease of explanation, we refer to watermarking the node as encoding the watermark region centered on the node throughout the paper.

it suggests significant performance degradation from watermarking. As such, $L_i$ is set to 1, denoting an unsatisfactory position for watermarking. For the $\hat{L}_i$ between 1 and $\beta$, we normalize the scores between $[0, 1]$. For GNN traning, we set $\beta = 1.01$. We provide additional analysis of different $\beta$ choices on AutoMarks performance in Appendix .3.

$$L_i = \begin{cases} 0 & \text{if } \hat{L}_i < 1 \\ \frac{\hat{L}_i - 1}{\beta} & \text{if } 1 \leq \hat{L}_i \leq \beta \\ 1 & \text{if } \hat{L}_i > \beta \end{cases} \quad (2)$$

*4.1.3* **Graph Neural Network Training**. The GNN's objective is to predict the quality metrics degradation from watermarking on the node. The model $\mathbf{M}$ consists of seven-layer graph convolutions with ReLU activation [15]. As in Equation 3, for the $l$-th layer, $\mathcal{N}(i)$ is the set of neighboring nodes of node $i$, $\mathbf{W}^{(l)}$ is the learnable weight matrix at layer $l$, $\mathbf{h}_j^{(l)}$ represents the feature vector of the neighboring node $j$ at layer $l$. The $\mathbf{h}_i^{(l+1)}$ is node $i$'s learned feature at layer $(l+1)$.

$$\mathbf{h}_i^{(l+1)} = \text{ReLU}\left(\sum j \in \mathcal{N}(i) \frac{1}{c_{ij}} \mathbf{W}^{(l)} \mathbf{h}_j^{(l)}\right) \quad (3)$$

$\mathbf{M}$ is trained to minimize the MSE loss between the predicted label $\mathbf{M}(G)$ and the ground truth label $L$ as in Equation 4.

$$\mathcal{L}_{\text{MSE}} = \frac{1}{|L|} \sum_{i=1}^{|L|} (\mathbf{M}(G_i) - L_i)^2 \quad (4)$$

*4.1.4* **Inference on Unseen Designs**. The unseen design is converted to a graph $G'$ following the **Graph Construction**. The trained GNN $\mathbf{M}$ scores each node in $G'$ and obtains score $L'$. The cell index $c_{wc} = \arg\min(L')$ is selected as the target cell to watermark. As in Fig. 1(A-4), the watermark region $R_w$ is centered at $c_{wc}$. The cells $C_w \in R_w$ are designated as the watermark cells.

When encoding watermarks on larger designs, where the chip canvas is much larger to accommodate more cells than the training design ispd19test6, a larger watermark region $R_w$ is required to ensure the watermark region is hard to remove. To address this, we perform a post-aggregation over the graph and obtain $L'_{agg}$ using Equation 5. The $L'_{agg}$ better estimates the cell neighbors over larger ranges without retraining the GNN on larger designs. The target cell index is obtained as $c_{wc} = argmin(L' + \gamma * L'_{agg})$. We set

$\gamma = 0.2$ here and provide additional analysis of different $\gamma$ choices on AUTOMARKS performance in Appendix .3.

$$L'_{agg,i} = \frac{1}{N} \sum_{k=1}^{N} \left( \frac{1}{|\mathcal{N}_k(i)|} \sum_{j \in \mathcal{N}_k(i)} L'_j \right) \tag{5}$$

## 4.2 Watermark Insertion

After obtaining the watermark region $R_w$ and the cells $C_w$, we encode them as a placement co-optimization objective to ensure only $C_w$ are placed in the region $R_w$. For a design with $K$ fence regions, AUTOMARKS encodes an additional fence region for $R_w$ with $C_w$ as the watermark following Equation 6 [9]. Here, $W$ is the wirelength term, and $D$ is the cell density term with density multiplier $\lambda$. $v$ denotes the $(x, y)$ location of cell and $e \in E$ is the design net.

$$\begin{aligned} \min_v \quad & \sum_{e \in E} W(e; v) + \lambda D(v), \\ \text{s.t.} \quad & v_k = (x_k, y_k) \in R_k, \quad k = 1, \cdots, R_K, R_w, \end{aligned} \tag{6}$$

The watermarked placement is then routed, and the design company obtains the watermarked layout.

## 4.3 Watermark Extraction

To prove ownership, the design company reverse-engineers the suspect layout [3, 31] and obtains netlist connectivities and cell locations. Such methods can recover large layouts containing over $7M$ cells with over 98% accuracy and efficiency. Then, AUTOMARKS queries the locations of $C_w$ and computes the cells $C'_w$ within the watermark region $R_w$. The watermark extraction rate ($\%WER$) is computed as in Equation 7.

$$\%WER = 100 \times \frac{|C'_w|}{|C_w|} \tag{7}$$

# 5 Experiments
## 5.1 Experiment Setups

*5.1.1 Hardware Infrastructure and Implementation Details.* AUTOMARKS is agnostic to the physical design algorithms. As a proof-of-concept, we employ the state-of-the-art physical design framework with DREAMPlace [9] as the placement algorithm and CUGR [19] as the routing algorithm. The graph neural network model is a seven-layer graph convolution network with ReLU activation. It is trained with the SGD optimizer, with the learning rate set to 0.01, weight decay set to 0.001, and the momentum value set to 0.9. It trains 30 epochs with a batch size set to 1280. The GNN is trained to sample 15, 20, 35, 50, 100, 200, 500 nodes at each hop.

*5.1.2 Benchmarks.* We evaluate AUTOMARKS's watermarking performance on the ISPD'2015 [6], and ISPD'2019 [20] benchmarks. Designs with fence region constraints are highlighted in blue in Tables 2 and 3.

*5.1.3 Baseline.* We compare AUTOMARKS with state-of-the-art frameworks employing both invasive and constraint-based watermarking:

- **Row Parity** [12, 14] inserts unique bit sequences as watermarks by shifting cells to different rows in the placement stage. Cells with a 1-bit are moved to an odd row, while cells with a 0-bit are moved to an even row.

- **Cell Scattering** [7] scatters watermark cells on the chip canvas as watermarks using pseudorandom coordinate transformation (PRCT) algorithms. Cells with 1-bit are moved along the y-axis, and cells with 0-bit are moved along the x-axis if they do not overlap with their neighbors.

- **ICMarks** [40] employs a scoring function to evaluate each layout subregion and select the best region of cells that adheres to design constraints and has minimal impact on the layout as the watermark region. The watermark region enforces only the watermark cells to be in the watermark region

- **Buffer Insertion** [39] adds additional buffers as watermarks during the placement stage. Two buffers are inserted to represent 0-bit, and one buffer is inserted to represent 1-bit.

We skip the baselines that: (i) have different watermarking targets, e.g., smaller full-custom IC designs [5, 29] and FPGAs [16, 36]; (ii) have similar watermarking approaches as our baselines, and we use the baselines as a proof-of-concept. For example, flip-flops are encoded as watermarks in [35] instead of buffers in [39].

*5.1.4 Evaluation Metrics.* We evaluate watermarking performance from the following metrics:

- **Placement WireLength Rate (PWLR)**: The rate of estimated half-perimeter wirelength (HPWL) of watermarked layout compared to the original one.

- **Routing WireLength Rate (RWLR)**: The rate of routed wirelength of watermarked layout compared to the original one.

- **Watermark Extraction Rate (WER)**: The percentage of watermark cells successfully extracted.

## 5.2 AUTOMARKS Results

*5.2.1 Fidelity and Transferability.* We compare the watermarking performance of AUTOMARKS and the baselines on the ISPD'2015 [6] benchmarks in Table 2, and the ISPD'2019 [20] benchmarks in Table 3. The encoded watermarks are extracted successfully, i.e., $\%WER = 100$ for all frameworks. We highlight that the AUTOMARKS is only trained on the ispd19test6 design and inference on the rest of the designs, significantly reducing the watermark search time and improving transferability.

**Comparison with Constraint-based WM:** Compared to AUTOMARKS that maintains watermarking fidelity, the baseline Row Parity [12, 14] and Cell Scattering [7] degrades the PWLR by 0.18% and 0.60% and the RWLR by 0.12% and 1.40% over the non-wm designs on ISPD'2015 [6] and ISPD'2019 [20] benchmarks respectively. The Row Parity [12, 14] and Cell Scattering [7] approaches do not consider the design constraints, like fence regions or macros when selecting the watermark cells, resulting in performance degradation.

Compared to AUTOMARKS, ICMarks [40] introduced slightly more degradations on ISPD'2019 [20] benchmarks. It is because AUTOMARKS learns the mapping from layout subgraphs to the actual PWLR improvement, which serves a better quality degradation estimation than the scoring function used in ICMarks [40].

**Comparison with Invasive WM:** Compared to AUTOMARKS, the invasive watermarking Buffer Insertion [39] significantly degrades the PWLR and RWLR by 5.36% and 9.01% on the ISPD'2015 [6] designs; and by 3.04% and 8.71% on the ISPD'2019 [20] benchmarks

| Design | Cells | Nets | Row Parity [12, 14] | | Cell Scattering [7] | | Buffer Insertion [39] | | ICMarks [40] | | AutoMarks | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | PWLR ↓ | RWLR ↓ | PWLR ↓ | RWLR ↓ | PWLR ↓ | RWLR ↓ | PWLR ↓ | RWLR ↓ | PWLR ↓ | RWLR ↓ |
| perf_a ★ | 108K | 115K | 1.0045 | 1.0155 | 0.9978 | 0.9967 | 1.5289 | 2.3270 | 0.9972 | **0.9873** | **0.9956** | 0.9929 |
| perf_b | 113K | 113K | 1.0058 | 1.0267 | 1.0020 | 1.0001 | 1.0176 | 1.0745 | **0.9890** | 0.9901 | 0.9898 | **0.9824** |
| dist_a | 127K | 134K | 1.0015 | 0.9999 | **0.9984** | **0.9965** | 1.0995 | 1.1072 | 1.0004 | 1.0052 | 1.0001 | 1.0040 |
| mult_b ★ | 146K | 152K | 1.0047 | 1.0199 | **0.9991** | 0.9923 | 1.8966 | 2.9374 | 0.9994 | 1.0028 | 1.0021 | **0.9922** |
| mult_c ★ | 146K | 152K | 1.0031 | 1.0252 | 0.9980 | 1.0023 | 1.6003 | 2.7459 | **0.9963** | **0.9979** | 1.0010 | 1.0022 |
| pci_a ★ | 30K | 34K | 1.0080 | 1.0340 | **0.9931** | **0.9846** | 1.6269 | 1.6849 | 0.9997 | 0.9950 | 0.9996 | 0.9974 |
| pci_b ★ | 29K | 33K | 1.0069 | 1.1173 | **0.9912** | 0.9977 | 1.2239 | 1.8207 | 0.9951 | 1.0026 | 0.9973 | **0.9885** |
| superblue11 ★ | 926K | 936K | 1.0052 | 1.0344 | 1.0278 | 1.0358 | 1.6930 | 3.7086 | **0.9986** | **0.9992** | 1.0073 | 1.0136 |
| superblue16 | 680K | 697K | 1.0030 | 1.0210 | **0.9988** | **0.9989** | 1.1023 | 1.1377 | 1.0017 | 1.0204 | 1.0025 | 1.0129 |
| perf_1 | 113K | 113K | 1.0035 | 1.0199 | 0.9985 | 0.9983 | 1.0150 | 1.0642 | 0.9964 | 0.9973 | **0.9909** | **0.9956** |
| fft_1 | 35K | 33K | 1.0017 | 1.0260 | 1.0167 | 1.0156 | 1.0680 | 1.1457 | **0.9671** | 0.9673 | 0.9711 | **0.9660** |
| fft_2 ★ | 35K | 33K | 1.0050 | 1.0260 | 1.0148 | 1.0114 | 1.4603 | 1.4476 | **0.9767** | **0.9770** | 0.9854 | 0.9823 |
| fft_a ★ | 34K | 32K | 1.0121 | 1.0200 | 1.0024 | 0.9933 | 1.8203 | 2.1923 | 0.9939 | 0.9895 | **0.9839** | **0.9773** |
| fft_b ★ | 34K | 32K | 1.0018 | 1.0075 | 0.9961 | 1.0030 | 1.9859 | 2.5615 | 0.9909 | **0.9890** | **0.9869** | 1.0039 |
| mult_1 | 160K | 159K | 1.0050 | 1.0238 | 1.0077 | 1.0065 | 1.0464 | 1.0631 | **0.9753** | **0.9744** | 0.9773 | 0.9769 |
| mult_2 | 160K | 159K | 1.0033 | 1.0199 | 0.9984 | 0.9967 | 1.0736 | 1.1147 | 0.9852 | 0.9900 | **0.9831** | **0.9872** |
| mult_a ★ | 154K | 154K | 1.0037 | 1.0105 | 0.9995 | 0.9968 | 1.3862 | 1.6738 | 0.9973 | **0.9916** | **0.9965** | 0.9934 |
| superblue12 | 1293K | 1293K | 1.0031 | 1.0067 | 0.9979 | 0.9956 | 1.0683 | 1.1044 | **0.9854** | **0.9732** | 0.9869 | 0.9783 |
| superblue14 | 634K | 620K | 1.0020 | 1.0057 | 0.9991 | 0.9981 | 1.0212 | 1.0286 | **0.9887** | **0.9867** | 1.0083 | 1.0037 |
| superblue19 | 522K | 512K | 1.0025 | 1.0077 | 1.0001 | 1.0005 | 1.0295 | 1.0672 | **0.9814** | **0.9809** | 1.0100 | 1.0563 |
| Average | - | - | 1.0043 | 1.0231 | 1.0018 | 1.0012 | 1.0536 | 1.0901 | **0.9908** | **0.9908** | 0.9937 | 0.9966 |

**Table 2: Performance on the ISPD'2015 benchmarks [6]. All the design watermarks are successfully extracted, i.e., WER = 100%. The PWLR and RWLR are the placement and routed wirelength rates over the original designs. The results in gray fail buffer insertion WM with significant degradation on the high-utilized designs (denoted with ★).**

| Design | Cells | Nets | Row Parity [12, 14] | | Cell Scattering [7] | | Buffer Insertion [39] | | ICMarks [40] | | AutoMarks | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | PWLR ↓ | RWLR ↓ | PWLR ↓ | RWLR ↓ | PWLR ↓ | RWLR ↓ | PWLR ↓ | RWLR ↓ | PWLR ↓ | RWLR ↓ |
| ispd19test1 | 9K | 3K | 1.0060 | 1.0129 | 0.9978 | 0.9998 | 1.0394 | 1.0619 | **0.9955** | 1.0015 | 1.0026 | **0.9989** |
| ispd19test2 | 73K | 72K | 1.0184 | 1.0201 | 1.0096 | 1.0016 | 1.0127 | 1.0408 | 0.9988 | 0.9999 | **0.9921** | **0.9927** |
| ispd19test3 | 8K | 9K | 1.0130 | 1.0475 | 1.0180 | 1.0193 | 1.0582 | 1.0801 | 1.0059 | 1.0045 | **0.9875** | **0.9865** |
| ispd19test4 | 151K | 146K | 1.0017 | 1.0050 | 0.9999 | 0.9998 | 1.1452 | 1.2494 | **0.9957** | **0.9907** | 0.9970 | 1.0001 |
| ispd19test5 | 29K | 29K | 1.0102 | 1.0556 | 1.0998 | 1.0975 | 0.9859 | 1.0121 | 1.0013 | 0.9998 | **0.9982** | **0.9953** |
| ispd19test6 | 181K | 180K | 1.0032 | 1.0106 | **0.9994** | 0.9993 | 1.0044 | 1.1108 | 1.0023 | 1.0026 | 1.0000 | **0.9971** |
| ispd19test7 | 362K | 359K | 1.0028 | 1.0160 | 1.0136 | 1.0109 | **1.0003** | 1.0717 | 1.0050 | 1.0054 | 1.0150 | 1.0139 |
| ispd19test8 | 543K | 538K | 1.0001 | 1.0082 | 0.9966 | 0.9958 | 1.0118 | 1.0806 | 0.9961 | **0.9929** | 0.9945 | 0.9996 |
| ispd19test9 | 903K | 895K | 1.0072 | 1.0108 | 1.0108 | 1.0095 | 1.0164 | 1.0814 | 1.0023 | 1.0023 | **0.9973** | **1.0000** |
| ispd19test10 | 903K | 895K | 1.0025 | 1.0090 | 1.0043 | 1.0108 | 1.0378 | 1.0982 | 0.9972 | **0.9966** | 0.9972 | 1.0002 |
| Average | - | - | 1.0065 | 1.0194 | 1.0060 | 1.0140 | 1.0304 | 1.0871 | 1.0000 | 0.9996 | **0.9981** | **0.9982** |

**Table 3: Performance on the ISPD'2019 benchmarks [20]. All the design watermarks are successfully extracted, i.e., WER = 100%. The PWLR and RWLR are the placement and routed wirelength rates over the original designs.**

over the non-watermarked designs. As the layout is highly utilized, adding redundant buffers results in significant cell placement reordering to accommodate the watermark cells, requiring more routing resources to connect the additional components than constraint-based AutoMarks.



**Figure 2: The watermark search time for different designs. The dotted line is the average search time of the large designs (≥ 500k cells).**

*5.2.2* **Efficiency**. Fig. 2 compares the watermark search time for ICMarks [40] and AutoMarks. For smaller designs, the search time of AutoMarks is very close to ICMarks [40]. However, for the large designs(≥ 500k cells), the average search time of AutoMarks is 176.38s, whereas ICMarks requires 320.93s. As ICMarks uses a fixed window size to traverse and score the layout, while AutoMarks employs GNN to batch-predict the region node scores, AutoMarks reduces the watermark search time for large designs by 45.04%.
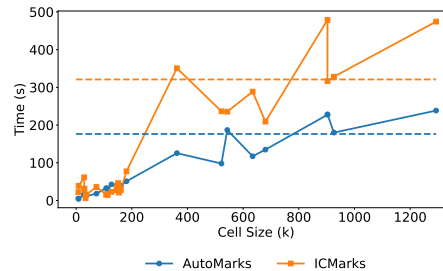
## 5.3 AutoMarks's Robustness

*5.3.1* **Watermark Removal Attacks**. Fig. 3 evaluates the impact of watermark removal attacks on AutoMarks and the baseline approaches. Due to space limitation, more details of the attacks under different parameter settings are in Appendix .1. The **Location swap attack** targeting Row Parity [12, 14] selects 0.1% of the total cells randomly and swaps their locations. The **Constraint**

(a) Location swap attack (0.1%)

(b) Constraint perturbation attack (10%)

(c) Optimization attack

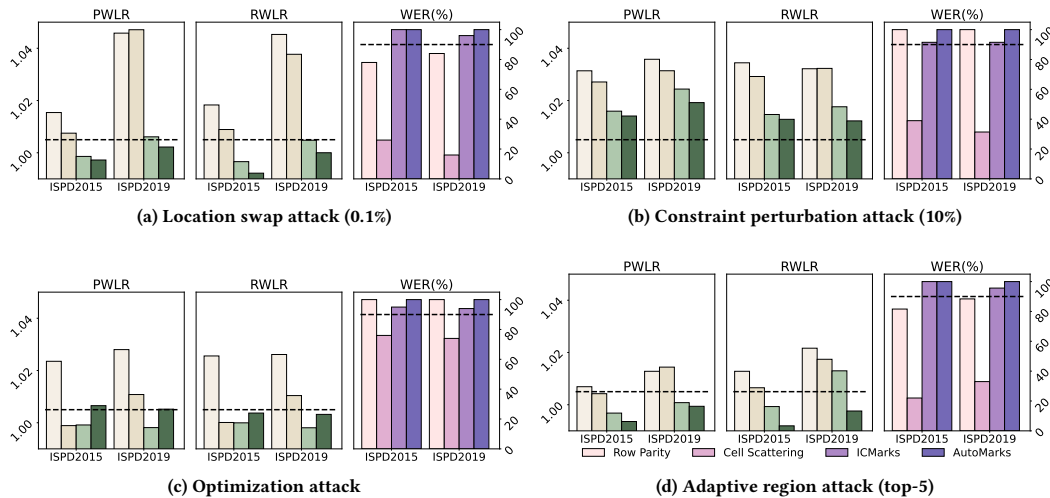(d) Adaptive region attack (top-5)

**Figure 3: Watermarking performance under different attacks for wirelength-driven placement on the ISPD'2015 [6] and ISPD'2019 [20] benchmarks. The black dotted line in the two left subfigures denotes the quality degradation threshold of 1.005, and the black dotted line in the rightmost subfigure denotes the watermark extraction threshold of 90%.**

**perturbation attack** targeting Cell Scattering [7] moves 10% of the cells randomly along the x-axis for one unit or the y-axis for one-row height if the cells do not overlap with the neighboring cells. The **Optimization attack** targeting all frameworks applies another round of detailed placement on the watermarked layout to remove signatures. The **Adaptive region attack** targeting IC-Marks [40] and AᴜᴛᴏMᴀʀᴋs, assumes the adversary knows the size of the watermark region and uses the evaluation function of IC-Marks [40] to identify the watermark regions. Then, the adversary perturbs cells in the top-5 regions to remove watermarks. We do not consider attacks on Buffer Insertion as it incurs significantly more performance degradation for watermark insertion than the constraint-based watermarking approaches.

AᴜᴛᴏMᴀʀᴋs is resilient to all attacks and maintains the watermark extraction rates of 100% even when the quality metrics PWLR and RWLR are greatly compromised; because as long as the watermark cells are within the watermark region, it is more robust toward such slight perturbation of cell locations. In contrast, Row Parity is vulnerable to Location swap attacks; Cell Scattering is vulnerable to Location swap attacks and Constraint perturbation attacks. While ICMarks employs a two-level watermarking framework to strengthen the watermarking strength, the detailed watermarking in ICMarks is vulnerable to watermark removal attacks. As a result, the overall WER for ICMarks slightly degrades.

*5.3.2* **Watermark Forging Attacks**. To forge the signature, the adversary needs to provide the watermark region by reproducing the GNN inference results. However, the adversary does not have access to the design used for GNN training, making it difficult to counterfeit the signature. Therefore, AᴜᴛᴏMᴀʀᴋs is resilient to the forging attacks.

## 5.4   Ablation Study and Analysis

This subsection provides ablation studies over the different hyperparameter choices in AᴜᴛᴏMᴀʀᴋs's performance using the ISPD'2019 [20] as benchmarks. Additional visualizations are in Appendix .2, and ablation studies are in Appendix .3.

*5.4.1* **The Effectiveness of Node Feature**. We analyze how different node features impact the GNN performance. In Table 4, we skip the cell location, size, and name as the node feature and report their PWLR and RWLR performance respectively. Other settings follow the default ones in AᴜᴛᴏMᴀʀᴋs.

Cell location and name have a more significant impact on the performance of the AᴜᴛᴏMᴀʀᴋs compared to cell size. Excluding the cell location and name into the node feature construction results in 2.9% and 3.2% PWLR and 3.4% and 3.3% RWLR degradation, respectively. In contrast, excluding the cell size only results in 0.2% degradation. This is mainly because the cell name indicates the type of cells, e.g., standard cells and macros, and cell location indicates the position of the cells. Both attributes are essential in helping AᴜᴛᴏMᴀʀᴋs learn which node is the ideal candidate for watermark insertion.

| Node Feature | PWLR | RWLR |
|---|---|---|
| All | 0.9981 | 0.9982 |
| w/o cell location | 1.0280 | 1.0312 |
| w/o cell size | 1.0002 | 0.9977 |
| w/o cell name | 1.0303 | 1.0304 |

**Table 4: The effectiveness of different node feature constructions on AᴜᴛᴏMᴀʀᴋs's performance.**

## 6   Conclusion

This work presents AᴜᴛᴏMᴀʀᴋs, an automated and highly transferable watermarking framework for physical design. Leveraging graph neural networks for the watermark region search, AᴜᴛᴏMᴀʀᴋs significantly reduces the search time while preserving watermarking fidelity. Extensive evaluations on the ISPD'15 and ISPD'19 benchmarks demonstrate the effectiveness and transferability of our proposed framework compared to existing physical design watermarking approaches. AᴜᴛᴏMᴀʀᴋs is resilient against watermarking removal and forging attacks with 100% watermark extraction rate for proof of ownership.

## Acknowledgments

# References

[1] Hervé Abdi and Lynne J Williams. 2010. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics* 2, 4 (2010), 433–459.

[2] Lilas Alrahis, Johann Knechtel, and Ozgur Sinanoglu. 2023. Graph Neural Networks: A Powerful and Versatile Tool for Advancing Design, Reliability, and Security of ICs. In *Proceedings of the 28th Asia and South Pacific Design Automation Conference (ASPDAC '23)*. Association for Computing Machinery, New York, NY, USA, 83–90. https://doi.org/10.1145/3566097.3568345

[3] Lilas Alrahis, Abhrajit Sengupta, Johann Knechtel, Satwik Patnaik, Hani Saleh, Baker Mohammad, Mahmoud Al-Qutayri, and Ozgur Sinanoglu. 2021. GNN-RE: Graph neural networks for reverse engineering of gate-level netlists. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 41, 8 (2021), 2435–2448.

[4] Lilas Alrahis, Abhrajit Sengupta, Johann Knechtel, Satwik Patnaik, Hani Saleh, Baker Mohammad, Mahmoud Al-Qutayri, and Ozgur Sinanoglu. 2022. GNN-RE: Graph Neural Networks for Reverse Engineering of Gate-Level Netlists. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 41, 8 (2022), 2435–2448. https://doi.org/10.1109/TCAD.2021.3110807

[5] Fujun Bai, Zhiqiang Gao, Yi Xu, and Xueyu Cai. 2007. A watermarking technique for hard IP protection in full-custom IC design. In *2007 International Conference on Communications, Circuits and Systems*. IEEE, 1177–1180.

[6] Ismail S Bustany, David Chinnery, Joseph R Shinnerl, and Vladimir Yutsis. 2015. ISPD 2015 benchmarks with fence regions and routing blockages for detailed-routing-driven placement. In *Proceedings of the 2015 Symposium on International Symposium on Physical Design*. 157–164.

[7] Xueyu Cai, Zhiqiang Gao, Fujun Bai, and Yi Xu. 2007. A watermarking technique for hard IP protection in post-layout design level. In *2007 7th International Conference on ASIC*. 1317–1320. https://doi.org/10.1109/ICASIC.2007.4415879

[8] Upma Gandhi, Ismail Bustany, William Swartz, and Laleh Behjat. 2019. A reinforcement learning-based framework for solving physical design routing problem in the absence of large test sets. In *2019 ACM/IEEE 1st Workshop on Machine Learning for CAD (MLCAD)*. IEEE, 1–6.

[9] Jiaqi Gu, Zixuan Jiang, Yibo Lin, and David Z Pan. 2020. DREAMPlace 3.0: Multi-electrostatics based robust VLSI placement with region constraints. In *Proceedings of the 39th International Conference on Computer-Aided Design*. 1–9.

[10] Hao-Hsiang Hsiao, Yi-Chen Lu, Pruek Vanna-Iampikul, and Sung Kyu Lim. 2024. FastTuner: Transferable Physical Design Parameter Optimization using Fast Reinforcement Learning. In *Proceedings of the 2024 International Symposium on Physical Design (ISPD '24)*. Association for Computing Machinery, New York, NY, USA, 93–101. https://doi.org/10.1145/3626184.3633328

[11] Xun Jiang, Zizheng Guo, Zhuomin Chai, Yuxiang Zhao, Yibo Lin, Runsheng Wang, and Ru Huang. 2023. Invited Paper: Accelerating Routability and Timing Optimization with Open-Source AI4EDA Dataset CircuitNet and Heterogeneous Platforms. In *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*. 1–9. https://doi.org/10.1109/ICCAD57390.2023.10323938

[12] Andrew B Kahng, John Lach, William H Mangione-Smith, Stefanus Mantik, Igor L Markov, Miodrag Potkonjak, Paul Tucker, Huijuan Wang, and Gregory Wolfe. 2001. Constraint-based watermarking techniques for design IP protection. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 20, 10 (2001), 1236–1252.

[13] Andrew B Kahng, Jens Lienig, Igor L Markov, and Jin Hu. 2011. *VLSI physical design: from graph partitioning to timing closure*. Vol. 312. Springer.

[14] Andrew B Kahng, Stefanus Mantik, Igor L Markov, Miodrag Potkonjak, Paul Tucker, Huijuan Wang, and Gregory Wolfe. 1998. Robust IP watermarking methodologies for physical design. In *Proceedings of the 35th annual Design Automation Conference*. 782–787.

[15] Yuanzhi Li and Yang Yuan. 2017. Convergence analysis of two-layer neural networks with relu activation. *Advances in neural information processing systems* 30 (2017).

[16] Wei Liang, Xingming Sun, Zhihua Xia, Decai Sun, and Jing Long. 2011. A chaotic IP watermarking in physical layout level based on FPGA. *Radioengineering* 20, 1 (2011), 118–125.

[17] Jai-Ming Lin, Chung-Wei Huang, Liang-Chi Zane, Min-Chia Tsai, Che-Li Lin, and Chen-Fa Tsai. 2021. Routability-driven global placer target on removing global and local congestion for VLSI designs. In *2021 IEEE/ACM International Conference On Computer Aided Design (ICCAD)*. IEEE, 1–8.

[18] Yibo Lin, Shounak Dhar, Wuxi Li, Haoxing Ren, Brucek Khailany, and David Z Pan. 2019. DREAMPlace: Deep learning toolkit-enabled gpu acceleration for modern vlsi placement. In *Proceedings of the 56th Annual Design Automation Conference 2019*. 1–6.

[19] Jinwei Liu, Chak-Wa Pui, Fangzhou Wang, and Evangeline F. Y. Young. 2020. CUGR: Detailed-Routability-Driven 3D Global Routing with Probabilistic Resource Model. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*. 1–6. https://doi.org/10.1109/DAC18072.2020.9218646

[20] Wen-Hao Liu, Stefanus Mantik, Wing-Kai Chow, Yixiao Ding, Amin Farshidi, and Gracieli Posser. 2019. ISPD 2019 initial detailed routing contest and benchmark with advanced routing rules. In *Proceedings of the 2019 International Symposium on Physical Design*. 147–151.

[21] Daniela Sánchez Lopera, Lorenzo Servadei, Gamze Naz Kiprit, Souvik Hazra, Robert Wille, and Wolfgang Ecker. 2021. A Survey of Graph Neural Networks for Electronic Design Automation. In *2021 ACM/IEEE 3rd Workshop on Machine Learning for CAD (MLCAD)*. 1–6. https://doi.org/10.1109/MLCAD52597.2021.9531070

[22] Yi-Chen Lu, Wei-Ting Chan, Vishal Khandelwal, and Sung Kyu Lim. 2022. Driving Early Physical Synthesis Exploration through End-of-Flow Total Power Prediction. In *Proceedings of the 2022 ACM/IEEE Workshop on Machine Learning for CAD* (Virtual Event, China) *(MLCAD '22)*. Association for Computing Machinery, New York, NY, USA, 97–102. https://doi.org/10.1145/3551901.3556476

[23] Yi-Chen Lu, Sai Surya Kiran Pentapati, Lingjun Zhu, Kambiz Samadi, and Sung Kyu Lim. 2020. TP-GNN: A Graph Neural Network Framework for Tier Partitioning in Monolithic 3D ICs. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*. 1–6. https://doi.org/10.1109/DAC18072.2020.9218582

[24] Yi-Chen Lu and Sung Kyu Lim. 2022. On Advancing Physical Design Using Graph Neural Networks. In *Proceedings of the 41st IEEE/ACM International Conference on Computer-Aided Design (ICCAD '22)*. Association for Computing Machinery, New York, NY, USA, Article 2, 7 pages. https://doi.org/10.1145/3508352.3561094

[25] Yi-Chen Lu, Sai Pentapati, and Sung Kyu Lim. 2021. The Law of Attraction: Affinity-Aware Placement Optimization using Graph Neural Networks. In *Proceedings of the 2021 International Symposium on Physical Design* (Virtual Event, USA) *(ISPD '21)*. Association for Computing Machinery, New York, NY, USA, 7–14. https://doi.org/10.1145/3439706.3447045

[26] Yi-Chen Lu, Tian Yang, Sung Kyu Lim, and Haoxing Ren. 2022. Placement Optimization via PPA-Directed Graph Clustering. In *Proceedings of the 2022 ACM/IEEE Workshop on Machine Learning for CAD* (Virtual Event, China) *(MLCAD '22)*. Association for Computing Machinery, New York, NY, USA, 1–6. https://doi.org/10.1145/3551901.3556482

[27] Yi-Chen Lu, Tian Yang, Sung Kyu Lim, and Haoxing Ren. 2022. Placement optimization via ppa-directed graph clustering. In *Proceedings of the 2022 ACM/IEEE Workshop on Machine Learning for CAD*. 1–6.

[28] Azalia Mirhoseini, Anna Goldie, Mustafa Yazgan, Joe Wenjie Jiang, Ebrahim Songhori, Shen Wang, Young-Joon Lee, Eric Johnson, Omkar Pathak, Azade Nazi, et al. 2021. A graph placement methodology for fast chip design. *Nature* 594, 7862 (2021), 207–212.

[29] Min Ni and Zhiqiang Gao. 2004. Watermarking system for IC design IP protection. In *2004 International Conference on Communications, Circuits and Systems (IEEE Cat. No. 04EX914)*, Vol. 2. IEEE, 1186–1190.

[30] Tingyuan Nie, Tomoo Kisaka, and Masahiko Toyonaga. 2005. A watermarking system for IP protection by a post layout incremental router. In *Proceedings of the 42nd annual Design Automation Conference*. 218–221.

[31] Rachel Selina Rajarathnam, Yibo Lin, Yier Jin, and David Z. Pan. 2020. ReGDS: A Reverse Engineering Framework from GDSII to Gate-level Netlist. In *2020 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. 154–163. https://doi.org/10.1109/HOST45689.2020.9300272

[32] Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084* (2019).

[33] Haoxing Ren and Jiang Hu. 2022. *Machine Learning Applications in Electronic Design Automation*. Springer.

[34] Masoud Rostami, Farinaz Koushanfar, and Ramesh Karri. 2014. A Primer on Hardware Security: Models, Methods, and Metrics. *Proc. IEEE* 102, 8 (2014), 1283–1295. https://doi.org/10.1109/JPROC.2014.2335155

[35] Debasri Saha, Parthasarathi Dasgupta, Susmita Sur-Kolay, and Samar Sen-Sarma. 2007. A novel scheme for encoding and watermark embedding in VLSI physical design for IP protection. In *2007 International Conference on Computing: Theory and Applications (ICCTA'07)*. IEEE, 111–116.

[36] Debasri Saha and Susmita Sur-Kolay. 2007. Fast robust intellectual property protection for VLSI physical design. In *10th International Conference on Information Technology (ICIT 2007)*. IEEE, 1–6.

[37] Majid Sarrafzadeh and CK Wong. 1996. *An introduction to VLSI physical design*. McGraw-Hill Higher Education.

[38] Atefeh Sohrabizadeh, Yunsheng Bai, Yizhou Sun, and Jason Cong. 2023. Robust GNN-Based Representation Learning for HLS. In *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*. 1–9. https://doi.org/10.1109/ICCAD57390.2023.10323853

[39] Guangyu Sun, Zhiqiang Gao, and Yi Xu. 2006. A watermarking system for ip protection by buffer insertion technique. In *7th International Symposium on Quality Electronic Design (ISQED'06)*. IEEE, 5–pp.

[40] Ruisi Zhang, Rachel Selina Rajarathnam, David Z Pan, and Farinaz Koushanfar. 2024. ICMarks: A Robust Watermarking Framework for Integrated Circuit Physical Design IP Protection. *arXiv preprint arXiv:2404.18407* (2024).

[41] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2020. Graph neural networks: A review of methods and applications. *AI open* (2020).

## .1 Additional Attack Evaluations

In this subsection, we include the watermark removal attack evaluations with different parameter settings, including location swap attack, constraint perturbation attack, and adaptive region attack.
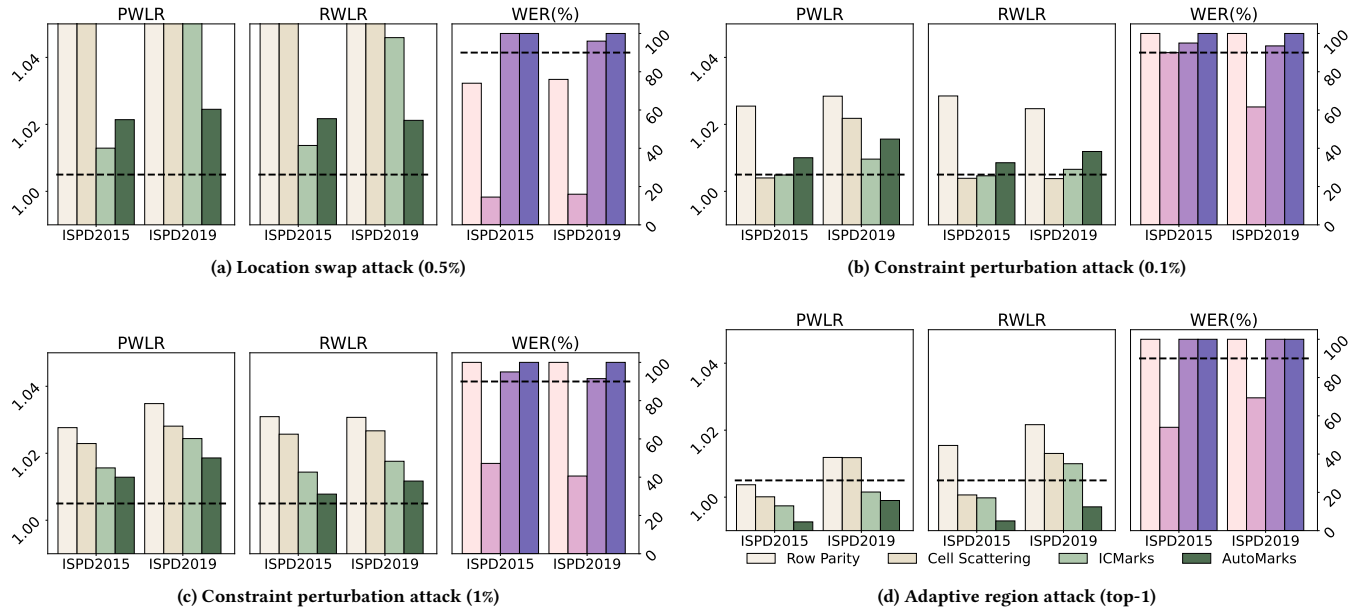


(a) Location swap attack (0.5%)

(b) Constraint perturbation attack (0.1%)

(c) Constraint perturbation attack (1%)

(d) Adaptive region attack (top-1)

**Figure 4: Watermarking performance under different attacks for wirelength-driven placement on the ISPD'2015 [6] and ISPD'2019 [20] benchmarks. The black dotted line in the two left subfigures denotes the quality degradation threshold of 1.005, and the black dotted line in the rightmost subfigure denotes the watermark extraction threshold of 90%.**

## .2 Additional Visualizations

In this section, we provide additional visualizations of AUTOMARKS's label and the training procedure.

***Visualization of the labels***. We visualize the training dataset `ispd19test6`'s labels $L$ distribution in Fig 5. As seen, watermarking on more than half of the nodes yields good watermarking performance. The critical step is to ensure the graph neural network learns the subgraph structure of the good nodes and can thus ensure the watermarking fidelity.

***Loss curve during GNN training***. We show the loss curve in Fig. 6, which demonstrates the graph neural network gradually learns to predict quality degradations of a given subgraph. The loss converges during the GNN training.
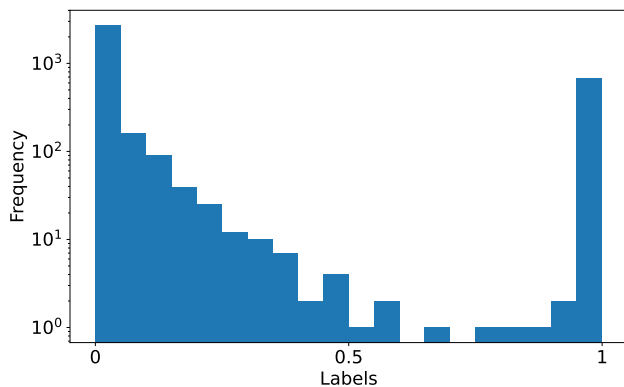


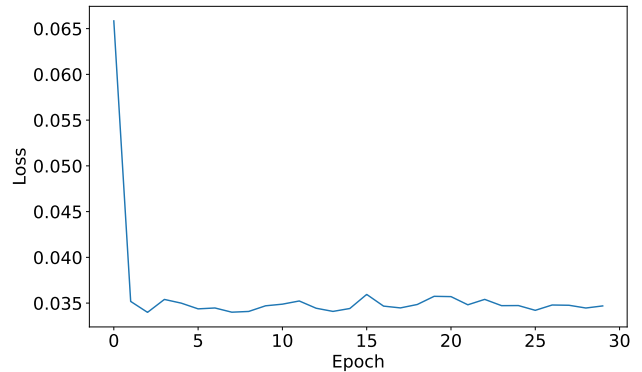**Figure 5: Histogram distribution of AUTOMARKS's labels.**



**Figure 6: Loss curve during the AUTOMARKS's GNN training procedure.**

***Stealthiness****.* To verify the stealthiness of AUTOMARKS, we include the watermarked layout and the non-watermarked layout from both ISPD'2015 [6] and ISPD'2019 [20] benchmarks in Fig. 7. As seen, the watermarks are invisible upon inspection and the adversary cannot differentiate between a non-watermarked layout and a layout watermarked by AUTOMARKS.
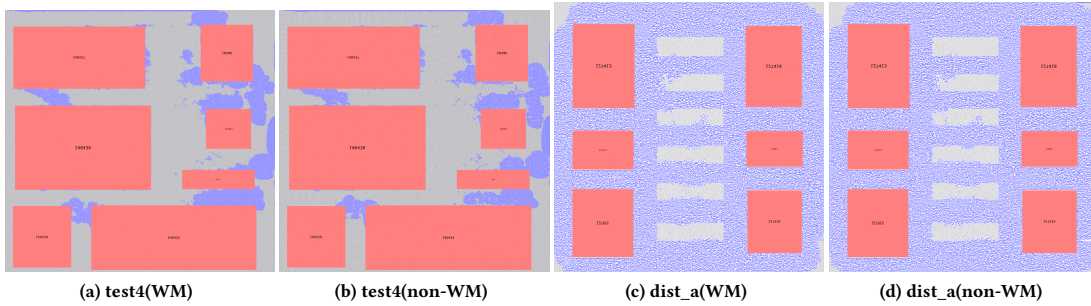


(a) test4(WM)　　　　　　(b) test4(non-WM)　　　　　　(c) dist_a(WM)　　　　　　(d) dist_a(non-WM)

**Figure 7: Watermarked and Non-watermarked examples.**

## .3　Additional Ablation Studies

***The Effectiveness of learning rates****.* We analyze how different GNN learning rates (LR) would affect AUTOMARKS's performance in Table 5. Following Sec. 5.2, the AUTOMARKS is trained on `ispd19test6` design and inference on the rest of the benchmarks. Here, we change the learning rate from 0.01 to 0.001 and 0.1. When the learning rate is increased from 0.001 to 0.01, the GNN learns to capture the node connections and predict the watermarking fidelity better. However, increasing the learning rate from 0.01 to 0.1 introduced marginal quality changes.

| LR | PWLR | RWLR |
|---|---|---|
| 0.001 | 1.0004 | 1.0009 |
| 0.01 | 0.9981 | 0.9982 |
| 0.1 | 0.9990 | 0.9998 |

**Table 5: The effectiveness of different learning rate choices on AUTOMARKS's performance.**

***The Effectiveness of $\beta$ Threshold****.* We analyze how different label threshold $\beta$ choice would impact the watermarking fidelity performance. In Table 6, we change the $\beta$ from 0.01 to 0.005 and 0.2, and show the PWLR and RWLR performance. Other settings follow the default ones in AUTOMARKS.

As seen, increasing the $\beta$ from 0.005 to 0.01 introduced significant PWLR and RWLR performance improvemnt. It is because more training nodes are normalized into the range of [0,1] that help AUTOMARKS to learn more diverse score predictions. However, increasing the $\beta$ from 0.01 to 0.02 does not introduce significant watermarking fidelity changes.

| $\beta$ | PWLR | RWLR |
|---|---|---|
| 0.005 | 1.0078 | 1.0089 |
| 0.01 | 0.9981 | 0.9982 |
| 0.02 | 1.0013 | 0.9990 |

**Table 6: The effectiveness of different $\beta$ choices on AUTOMARKS's performance.**

***The effectiveness of $\gamma$****.* We analyze how different $\gamma$ choices would affect the larger designs' (i.e. cell size larger than 900k) performance in Table 7. Other settings follow the default ones in AUTOMARKS. As seen, adjusting the $\gamma$ from 0.2 to 0 and 0.5 will change the AUTOMARKS performance on the large layout slightly.

| $\gamma$ | PWLR | RWLR |
|---|---|---|
| 0 | 0.9974 | 1.0002 |
| 0.2 | 0.9972 | 1.0001 |
| 0.5 | 0.9974 | 1.0002 |

**Table 7: The effectiveness of different $\gamma$ choices on AUTOMARKS's performance.**

*.3.1　**The effectiveness of GNN layers****.* We analyze how different numbers of the GNN layers would affect the watermarking fidelity performance in Table 8. Here, we fix the watermark region size $N = 10$ to be the same as Section 5.2 and change the number of GNN layers used for training and inference from $5 \rightarrow 10$.

Increasing the layer from 5 to 7 significantly improved the watermarking performance. When the layer number is set to 5, most of the cells are within the watermark region, and the GNN fails to learn the node features on the region boundary, which will be expelled outside the region and play an essential role in quality degradation. Increasing the number of layers results in better watermarking performance. However, increasing the number from 7 to 10 introduced marginal quality improvement and more computational overheads. Therefore, AUTOMARKS employs a graph neural network with 7 layers.

| Layer Num. | PWLR | RWLR |
|:---:|:---:|:---:|
| 5 | 1.1033 | 1.1131 |
| 7 | 0.9981 | 0.9982 |
| 10 | 0.9997 | 0.9999 |

**Table 8: The effectiveness of different GNN layers on AutoMarks's performance.**