# AttackGNN: Red-Teaming GNNs in Hardware Security Using Reinforcement Learning

**Vasudev Gohil**

Satwik Patnaik

Dileep Kalathil

Jeyavijayan "JV" Rajendran
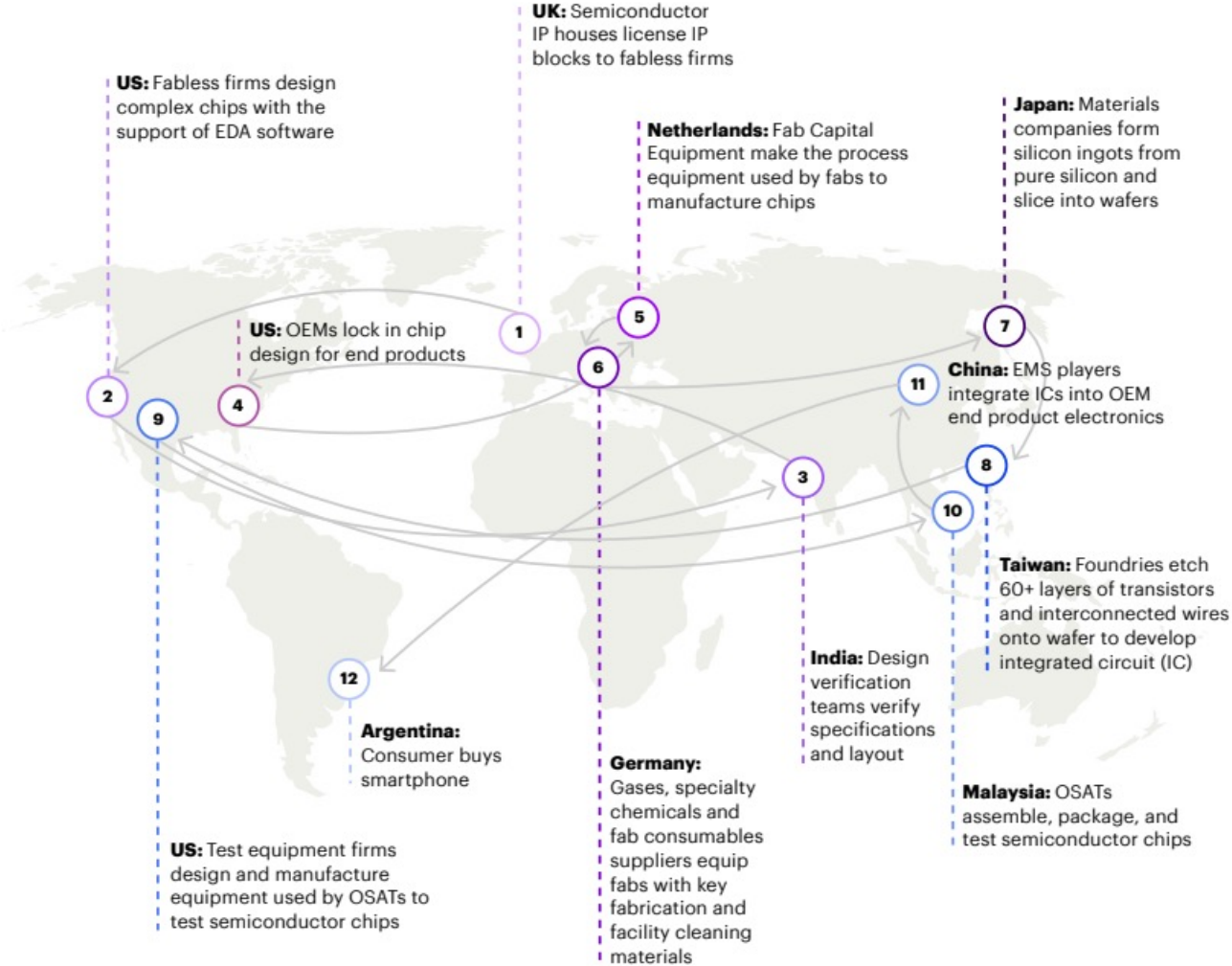
TEXAS A&M UNIVERSITY
Engineering

UNIVERSITY OF DELAWARE

# Hardware-focused Threats to Computing Systems

## Due to Globalized Supply Chain



**UK:** Semiconductor IP houses license IP blocks to fabless firms

**US:** Fabless firms design complex chips with the support of EDA software

**Netherlands:** Fab Capital Equipment make the process equipment used by fabs to manufacture chips

**Japan:** Materials companies form silicon ingots from pure silicon and slice into wafers

**US:** OEMs lock in chip design for end products

**China:** EMS players integrate ICs into OEM end product electronics

**Taiwan:** Foundries etch 60+ layers of transistors and interconnected wires onto wafer to develop integrated circuit (IC)

**India:** Design verification teams verify specifications and layout

**Argentina:** Consumer buys smartphone

**Germany:** Gases, specialty chemicals and fab consumables suppliers equip fabs with key fabrication and facility cleaning materials

**Malaysia:** OSATs assemble, package, and test semiconductor chips

**US:** Test equipment firms design and manufacture equipment used by OSATs to test semiconductor chips

# Hardware-focused Threats to Computing Systems

## Due to Globalized Supply Chain

Real     Fake

Counterfeiting



Overproduction

**US:** Fabless firms design complex chips with the support of EDA software

**UK:** Semiconductor IP houses license IP blocks to fabless firms

**Netherlands:** Fab Capital Equipment make the process equipment used by fabs to manufacture chips

**Japan:** Materials companies form silicon ingots from pure silicon and slice into wafers

**US:** OEMs lock in chip design for end products

1

5

6

2

9

4

7

**China:** EMS players integrate ICs into OEM end product electronics

11

8

3

10

12

**Argentina:** Consumer buys smartphone

**Germany:** Gases, specialty chemicals and fab consumables suppliers equip fabs with key fabrication and facility cleaning materials

**India:** Design verification teams verify specifications and layout

**Taiwan:** Foundries etch 60+ layers of transistors and interconnected wires onto wafer to develop integrated circuit (IC)

**Malaysia:** OSATs assemble, package, and test semiconductor chips

**US:** Test equipment firms design and manufacture equipment used by OSATs to test semiconductor chips

IP Piracy

Hardware Trojans

Reverse Engineering

# State-of-the-art GNNs in Hardware Security

| Technique Type | Security Problem | Technique | GNN Framework | Claimed Efficacy |
|---|---|---|---|---|
| Defense | Detecting Trojans | GNN4TJ [1] | Attention-based custom GCN | 97% TPR |
| | Locating Trojans | TrojanSAINT [2] | Graph attention network | 98% TPR, 96% TNR |
| | Detecting IP Piracy | GNN4IP [3] | Attention-based custom GCN | 94.61% Accuracy |
| Attack | Reverse Engineering | GNN-RE [4] | Graph attention network | 98.87% Accuracy |
| | Hardware Obfuscation | OMLA [5] | Graph isomorphism network | 89.55% Accuracy |

# State-of-the-art GNNs in Hardware Security

| Technique Type | Security Problem | Technique | GNN Framework | Claimed Efficacy |
|---|---|---|---|---|
| | | | | |
| | Trojans | [2] | network | |
| | Detecting IP Piracy | GNN4IP [3] | Attention-based custom GCN | 94.61% Accuracy |
| Attack | Reverse Engineering | GNN-RE [4] | Graph attention network | 98.87% Accuracy |
| | Hardware Obfuscation | OMLA [5] | Graph isomorphism network | 89.55% Accuracy |

**Are Graph Neural Networks (GNNs) Used To Solve Hardware Security Problems Robust?**

# State-of-the-art GNNs in Hardware Security

| Technique Type | Security Problem | Technique | GNN Framework | Claimed Efficacy |
|---|---|---|---|---|

**Are Graph Neural Networks (GNNs) Used To Solve Hardware Security Problems Robust?**

| | Trojans | [2] | network | |
|---|---|---|---|---|
| | Detecting IP | GNN4IP [3] | Attention-based | 94.61% Accuracy |

**AttackGNN: Red-Teaming GNNs in Hardware Security Using Reinforcement Learning**

| | Hardware Obfuscation | OMLA [5] | Graph isomorphism network | 89.55% Accuracy |

# Threat Model

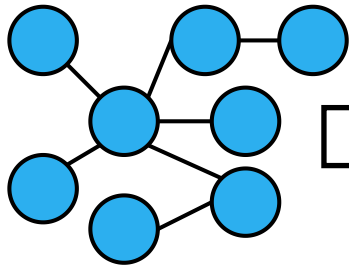## Standard attack model of adversarial attacks



Kevin Eykholt et al., "Robust physical-world attacks on deep learning visual classification," In Proc. of CVPR, 2018

# Threat Model

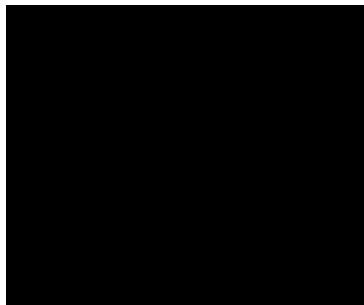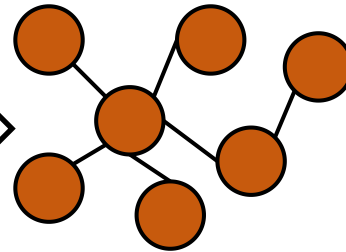## Standard attack model of adversarial attacks



Kevin Eykholt et al., "Robust physical-world attacks on deep learning visual classification," In Proc. of CVPR, 2018



Trained GNN

No Modifications

# Threat Model

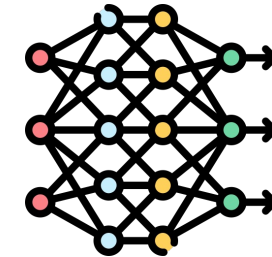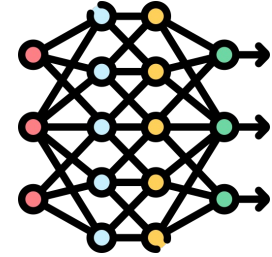## Standard attack model of adversarial attacks

Original Circuit



Perturbed Circuit



Trained GNN



No Modifications
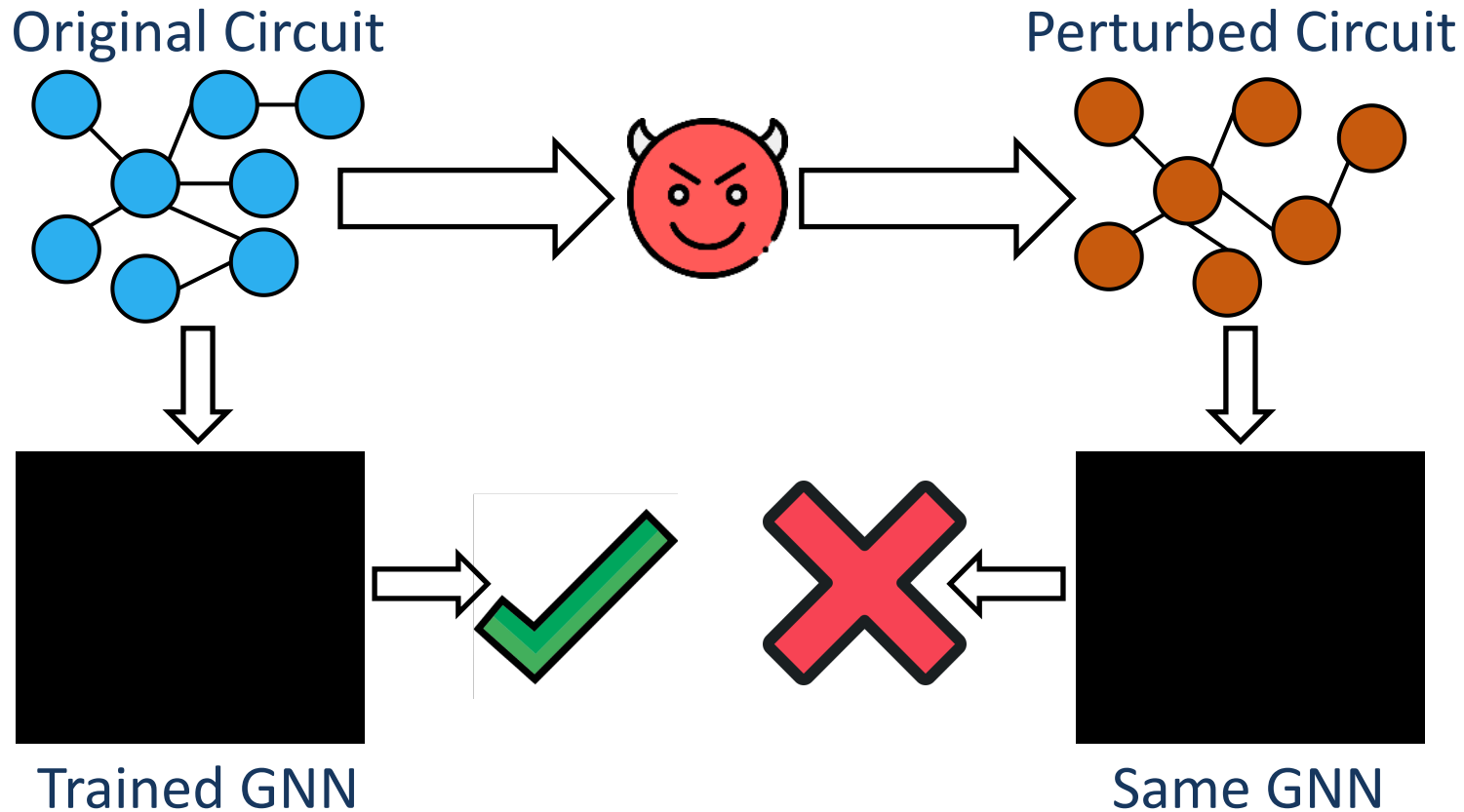
Perturbations Following Circuit Design Rules



Kevin Eykholt et al., "Robust physical-world attacks on deep learning visual classification," In Proc. of CVPR, 2018

# Threat Model
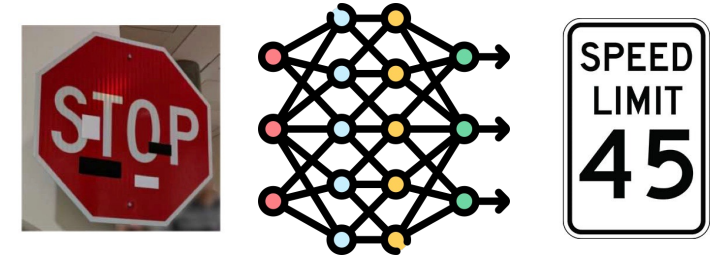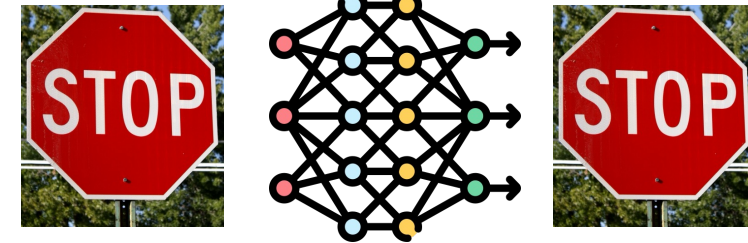
## Standard attack model of adversarial attacks
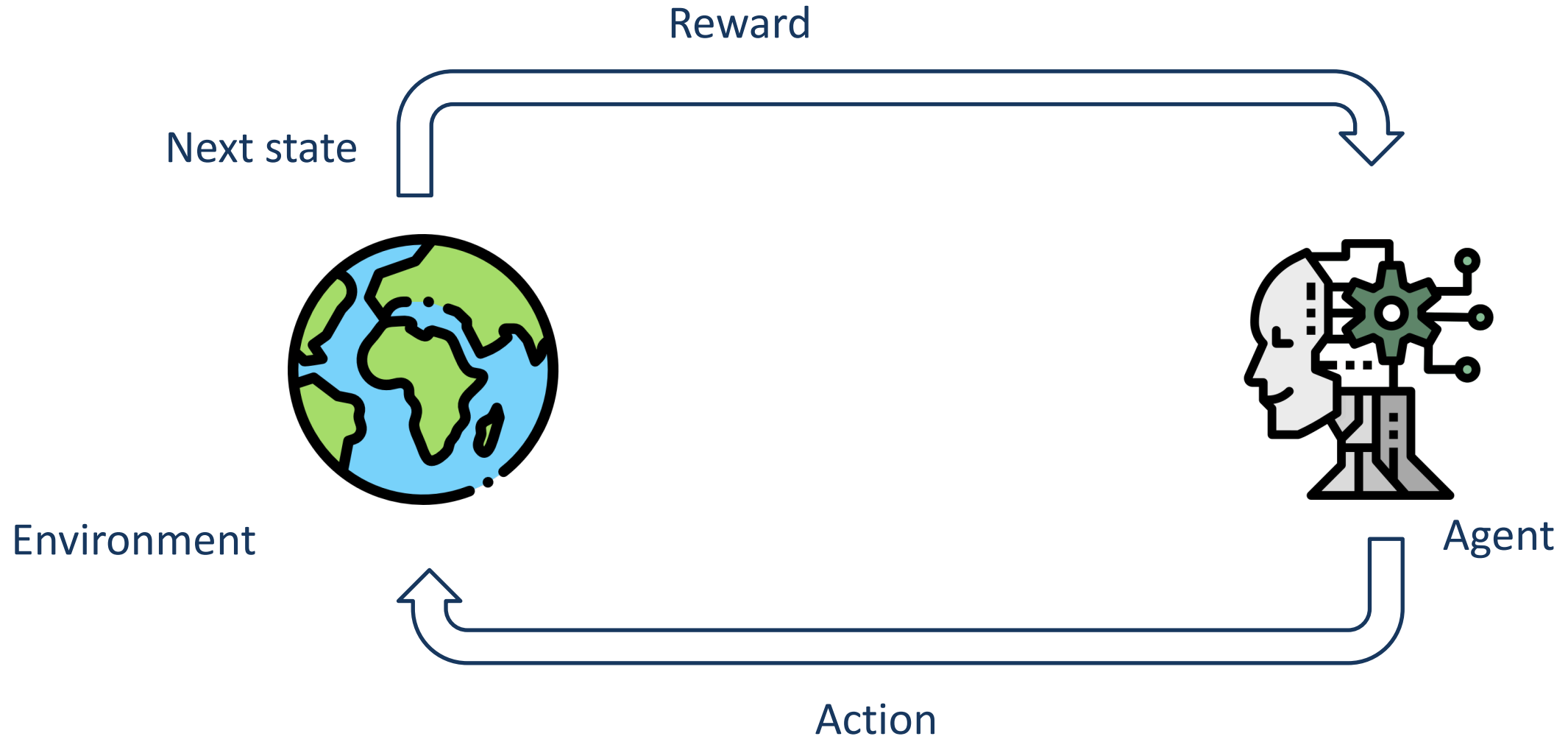
Original Circuit → 😈 → Perturbed Circuit



Kevin Eykholt et al., "Robust physical-world attacks on deep learning visual classification," In Proc. of CVPR, 2018

Trained GNN

No Modifications

Perturbations Following Circuit Design Rules

Black-box Access

# Threat Model

## Standard attack model of adversarial attacks

Original Circuit

Perturbed Circuit

Trained GNN

Same GNN

Goal: Misclassification

Kevin Eykholt et al., "Robust physical-world attacks on deep learning visual classification," In Proc. of CVPR, 2018

No Modifications

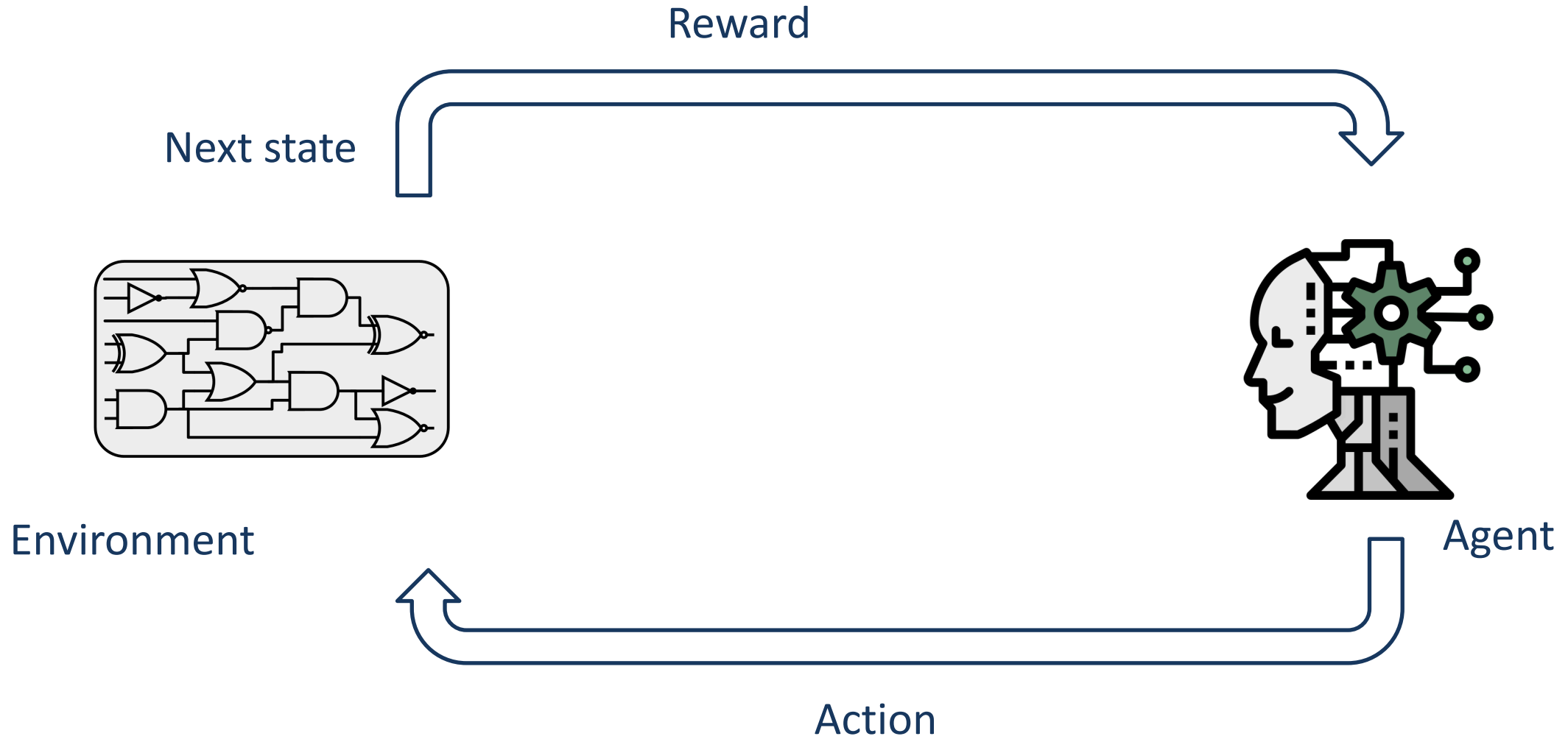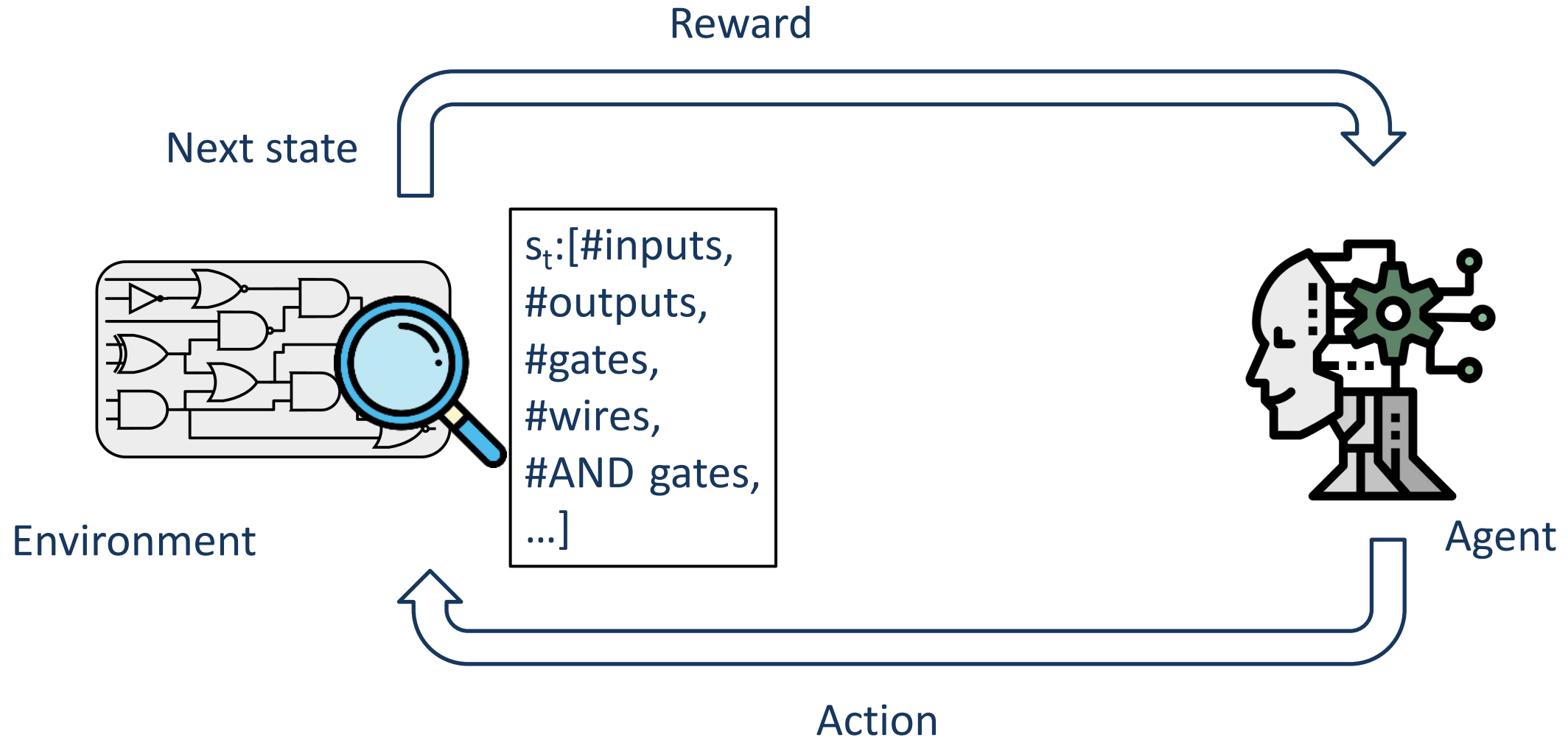Perturbations Following Circuit Design Rules

Black-box Access

# AttackGNN – Preliminary Formulation

Reward

Next state



Environment

Agent

Action

# AttackGNN – Preliminary Formulation

Reward

Next state



Environment

Agent

Action

# AttackGNN – Preliminary Formulation

Reward

Next state

$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, ...]

Environment

Agent

Action

# AttackGNN – Preliminary Formulation



Reward

Next state

$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, …]

Environment

Agent

Action

$a_t$: circuit transformation command

# AttackGNN – Preliminary Formulation

$$r_t = \begin{cases} \boldsymbol{\alpha} \ (> \boldsymbol{0}) & \text{if next state is misclassified} \\ \boldsymbol{0} & \text{else} \end{cases}$$

Reward

Next state

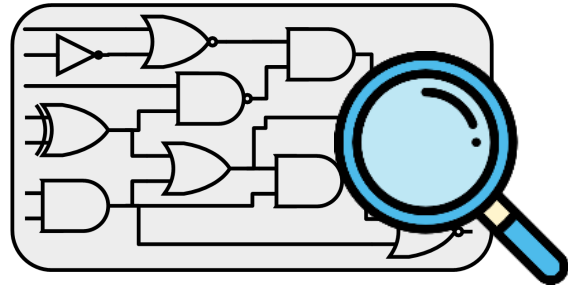$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, …]

Environment

Agent

Action

$a_t$: circuit transformation command
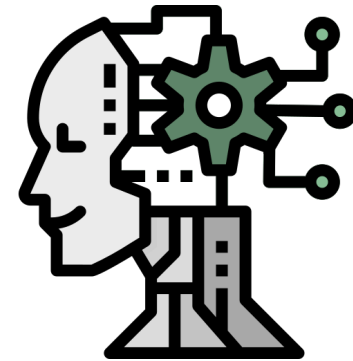
# AttackGNN – Preliminary Formulation

$$r_t = \begin{cases} \boldsymbol{\alpha} \ (> \mathbf{0}) & \text{if next state is misclassified} \\ \mathbf{0} & \text{else} \end{cases}$$

Reward

Next state

$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, ...]
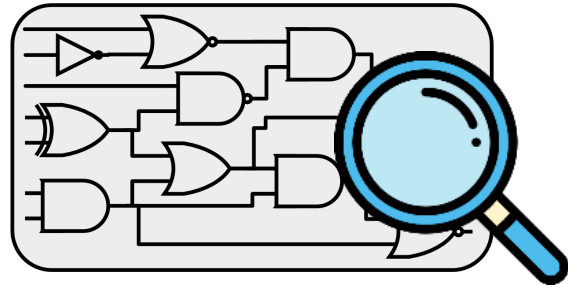
Environment

Agent

Action  "rewrite"  "refactor"

$a_t$: circuit transformation command

# AttackGNN – Challenges

$$r_t = \begin{cases} \alpha \; (> 0) & \text{if next state is misclassified} \\ 0 & \text{else} \end{cases}$$

Reward

Next state

$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, ...]

Environment

Agent

Action "rewrite" "refactor"

$a_t$: circuit transformation command
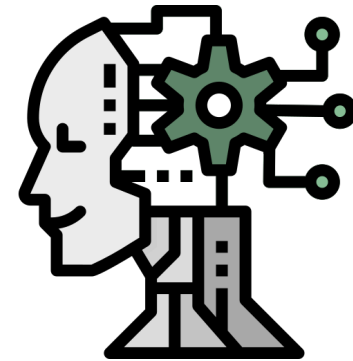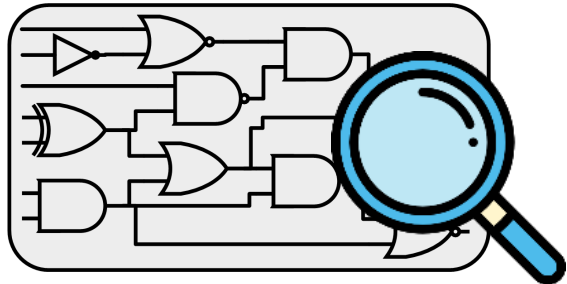
# AttackGNN – Challenges

$$r_t = \begin{cases} \boldsymbol{\alpha} \, (> \mathbf{0}) & \text{if next state is misclassified} \\ \mathbf{0} & \text{else} \end{cases}$$

Reward

② ⚠

Unnecessary Reward Computations

Next state

③ ⚠

MDP Specific to One GNN

$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, …]

Environment

Agent

Ineffective and Specific Actions

① ⚠

Action "rewrite" "refactor"

$a_t$: circuit transformation command

# AttackGNN – Solutions

$$r_t = \begin{cases} \boldsymbol{\alpha} \ (> \mathbf{0}) & \text{if next state is misclassified} \\ \mathbf{0} & \text{else} \end{cases}$$

Reward

②

**Unnecessary Reward Computations**

Next state

③

**MDP Specific to One GNN**

$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, …]

Environment

Agent

**Ineffective and Specific Actions**

①

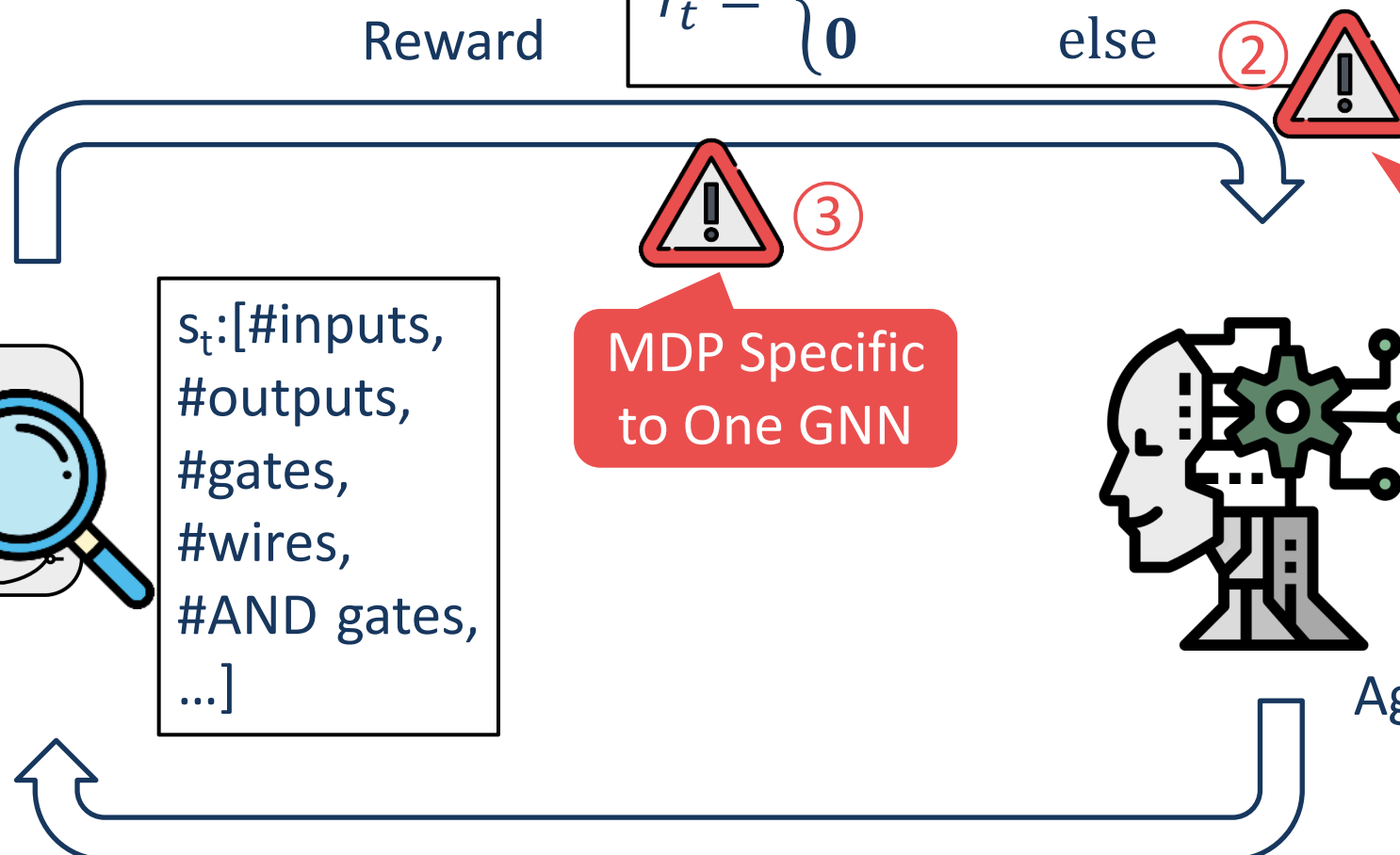Action "rewrite" "refactor"

$a_t$: circuit transformation command

# AttackGNN – Solutions
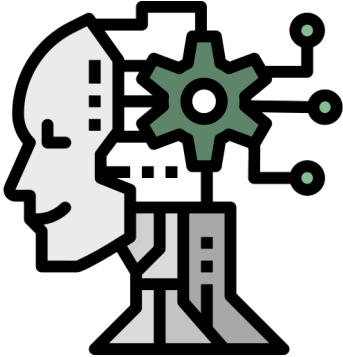
$$r_t = \begin{cases} \boldsymbol{\alpha} \; (> \mathbf{0}) & \text{if next state is misclassified} \\ \mathbf{0} & \text{else} \end{cases}$$

Reward

② ⚠

Unnecessary Reward Computations

Next state

③ ⚠

MDP Specific to One GNN

$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, ...]

Environment

Agent

Ⓐ Effective, Generalizable Actions

Ineffective and Specific Actions

① ⚠

Action   Don't use 3-input AND gates

$a_t$: allowed/unallowed gate types

# AttackGNN – Solutions

Reward

$$r_t = \begin{cases} \boldsymbol{\alpha} \ (> \boldsymbol{0}) & \text{if next state is misclassified} \\ \boldsymbol{0} & \text{else} \end{cases}$$

② ⚠

Unnecessary Reward Computations

Next state

③ ⚠

$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, …]

MDP Specific to One GNN

Ⓑ Sparse Rewards

Environment

Agent

Ⓐ Effective, Generalizable Actions

Ineffective and Specific Actions

① ⚠

Action   Don't use 3-input AND gates
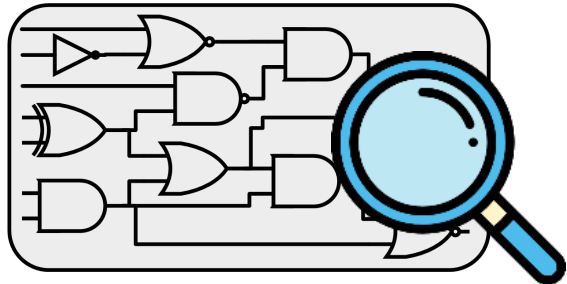
$a_t$: allowed/unallowed gate types

# AttackGNN – Solutions

$$r_t = \begin{cases} \boldsymbol{\alpha} \, (> \mathbf{0}) & \text{if next state is misclassified} \\ \mathbf{0} & \text{else} \end{cases}$$

Reward

② Unnecessary Reward Computations

Next state

③ MDP Specific to One GNN

$s_t$:[#inputs, #outputs, #gates, #wires, #AND gates, …]

© Multi-task Learning: Contextual MDP

Ⓑ Sparse Rewards

Environment

Agent

Ⓐ Effective, Generalizable Actions

① Ineffective and Specific Actions

Action  Don't use 3-input AND gates
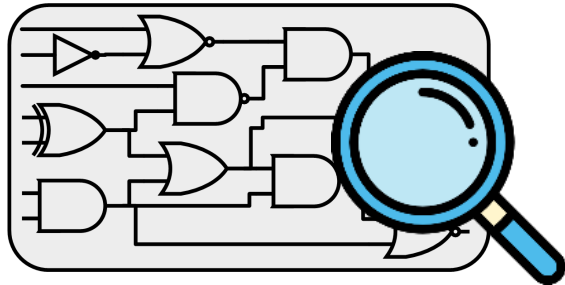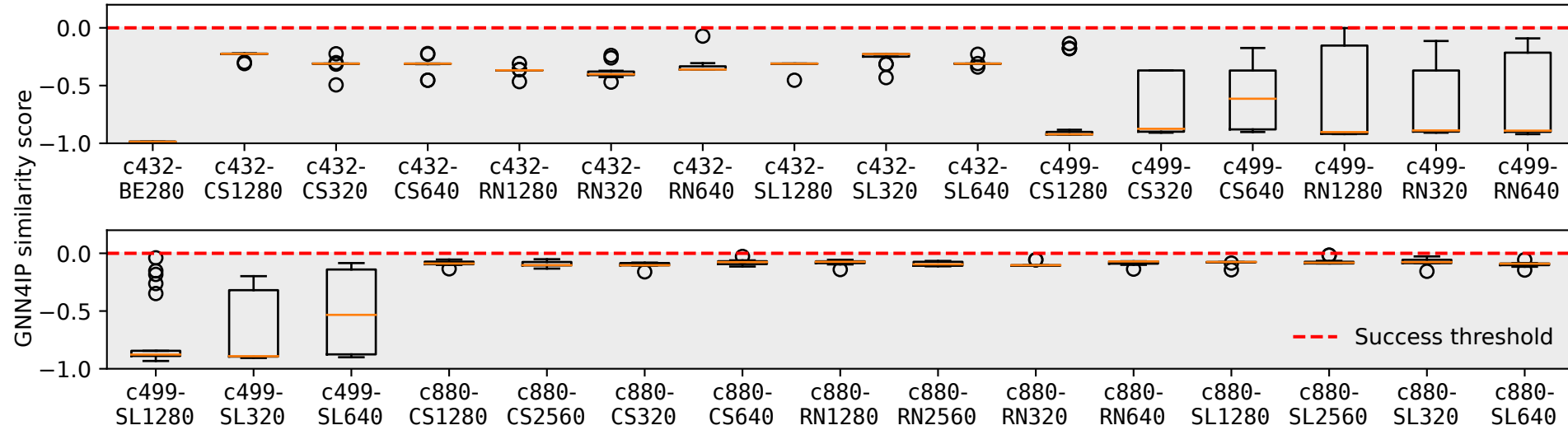
$a_t$: allowed/unallowed gate types

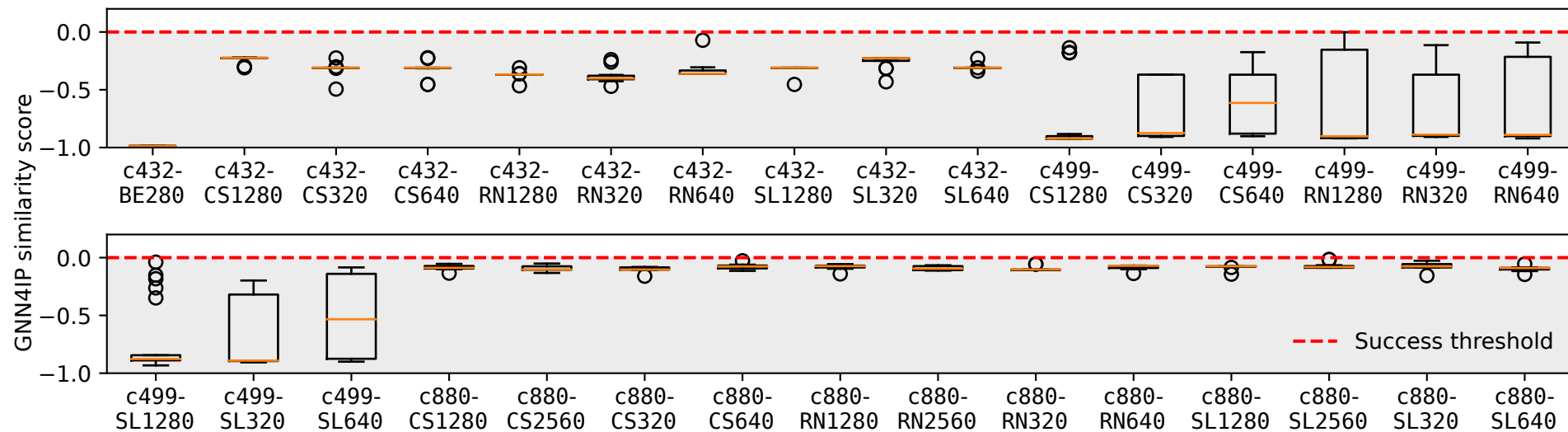# AttackGNN – Results

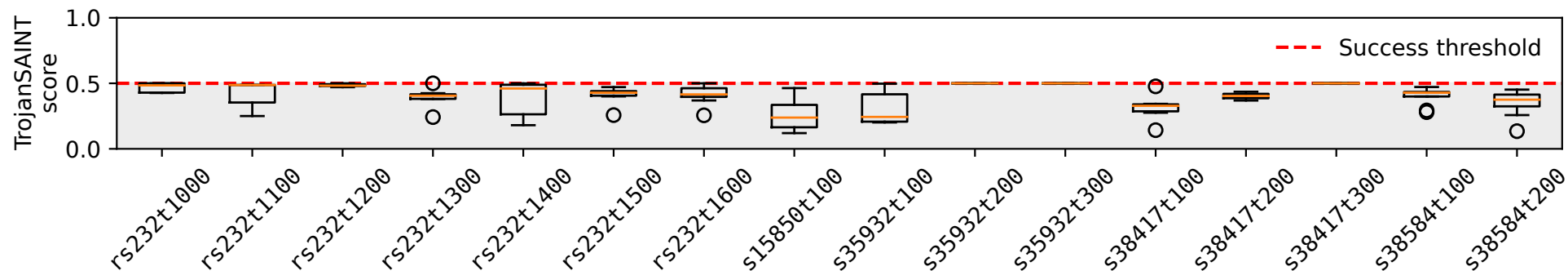## Against GNN4IP (IP Piracy Detection GNN)

# AttackGNN – Results
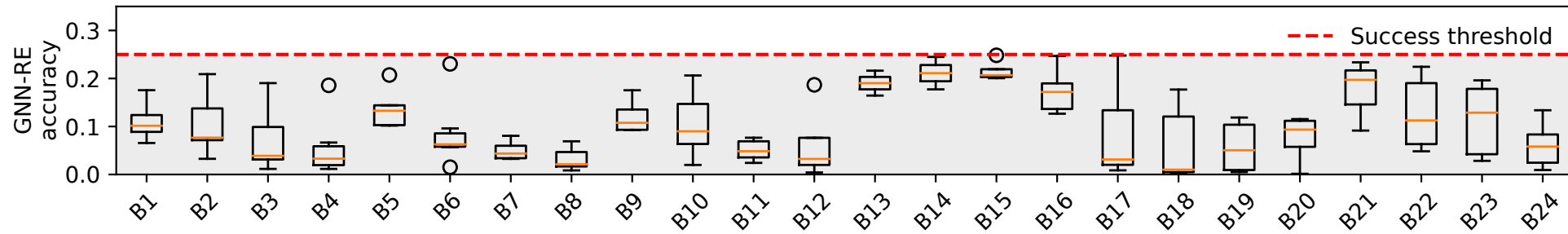
**Against GNN4IP (IP Piracy Detection GNN)**

**Against TrojanSAINT (Trojan Locator GNN)**

# AttackGNN – Results

## Against GNN-RE
## (Reverse Eng. GNN)



GNN-RE accuracy box plot with y-axis from 0.0 to 0.3, x-axis labels B1–B24, and a red dashed "Success threshold" line at 0.25.

## Against GNN4TJ
## (Trojan Detector GNN)

**GNN4TJ predictions**

| | | HT-infested | HT-free |
|---|---|---|---|
| **True labels** | HT-infested | 19 | 0 |
| | HT-free | 15 | 0 |

## Against OMLA
## (De-obfuscation GNN)



OMLA KPA box plot with y-axis from 0.35 to 0.65, x-axis labels c1355, c1908, c2670, c3540, and red dashed "Success threshold" lines at 0.45 and 0.55.

# AttackGNN – Results

## Against GNN-RE (Reverse Eng. GNN)



## Against GNN4TJ (Trojan Detector GNN)

**GNN4TJ predictions**

| | | HT-infested | HT-free |
|---|---|---|---|
| **True labels** | HT-infested | 19 | 0 |
| | HT-free | 15 | 0 |

## Against OMLA (De-obfuscation GNN)



Success rate of all GNNs against AttackGNN-generated adversarial circuits: **0%**

GNNs used in hardware security are **not robust**!

# Thank You

Vasudev Gohil
vasudevgohil.com

Secure and Trustworthy Hardware (SETH) Lab
https://seth.engr.tamu.edu
Texas A&M University

# References

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

[2] Lashen, Hazem, Lilas Alrahis, Johann Knechtel, and Ozgur Sinanoglu. "TrojanSAINT: Gate-level netlist sampling-based inductive learning for hardware Trojan detection."arXiv preprint arXiv:2301.11804, 2023.

[3] Yasaei, Rozhin, Shih-Yuan Yu, Emad Kasaeyan Naeini, and Mohammad Abdullah Al Faruque. "GNN4IP: Graph neural network for hardware intellectual property piracy detection." In 58th ACM/IEEE Design Automation Conference (DAC), pp. 217-222, IEEE, 2021.

[4] Alrahis, Lilas, Abhrajit Sengupta, Johann Knechtel, Satwik Patnaik, Hani Saleh, Baker Mohammad, Mahmoud Al-Qutayri, and Ozgur Sinanoglu. "GNN-RE: Graph neural networks for reverse engineering of gate-level netlists." IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems 41, no. 8: 2435-2448, 2021.

[5] Alrahis, Lilas, Satwik Patnaik, Muhammad Shafique, and Ozgur Sinanoglu. "OMLA: An oracle-less machine learning-based attack on logic locking." IEEE Transactions on Circuits and Systems II: Express Briefs 69, no. 3: 1602-1606, 2021.