



Risk-Aware and Explainable Framework for Ensuring Guaranteed Coverage in Evolving Hardware Trojan Detection

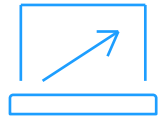
Rahul Vishwakarma

Graduate Research Assistant
California State University Long Beach

Amin Rezaei

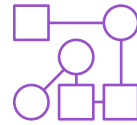
Assistant Professor
California State University Long Beach

Agenda



Introduction

Problem
statement and
preliminaries



PALETTE

Solution and
experimental
results



Inference

Practical use
cases in
decision making



Hardware Trojans

Applied machine learning to detect hardware trojans

Hardware Trojans

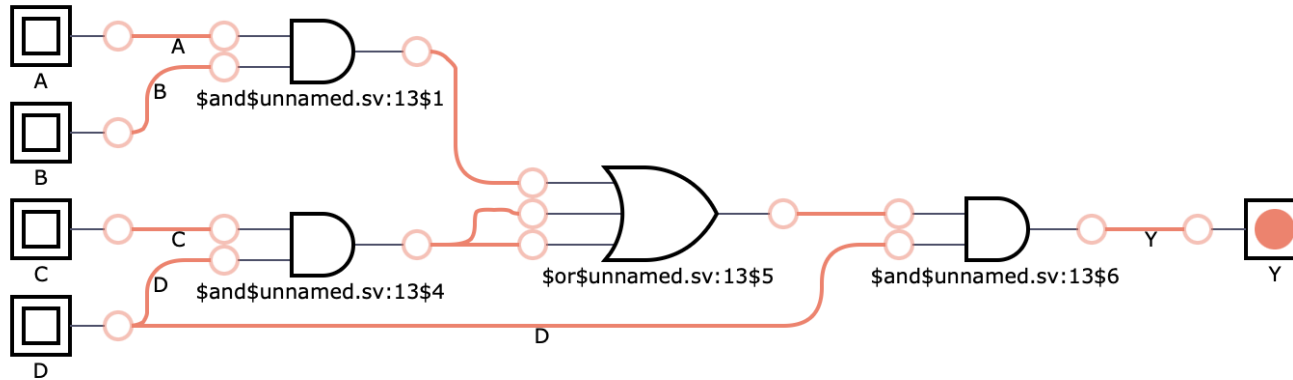


Figure 1: Trojan free circuit

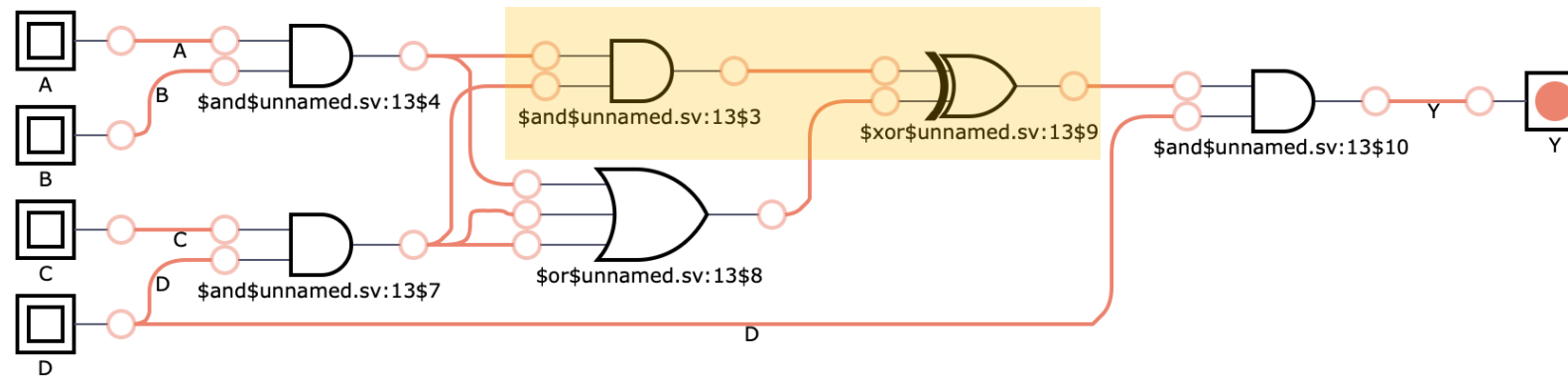


Figure 2: Trojan inserted circuit

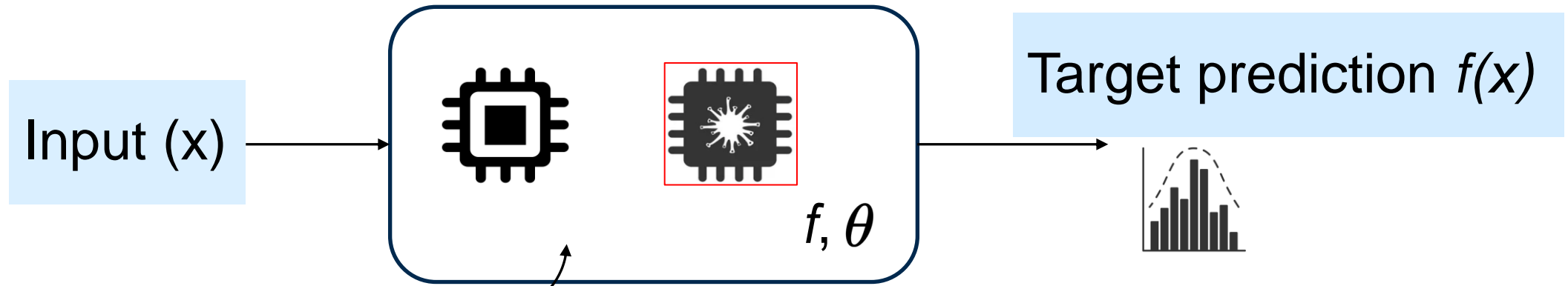
Conventional Approach

- Side-channel Power analysis
- Temperature analysis
- Wireless transmission power analysis
- Regional supply current analysis

Challenges

- Complexity of modern Integrated Circuits
- Challenges in detecting minute Hardware Trojans
- Early-Stage detection significance in design process
- Cost-Efficient Hardware Trojan detection strategies

Machine Learning



Model of problem

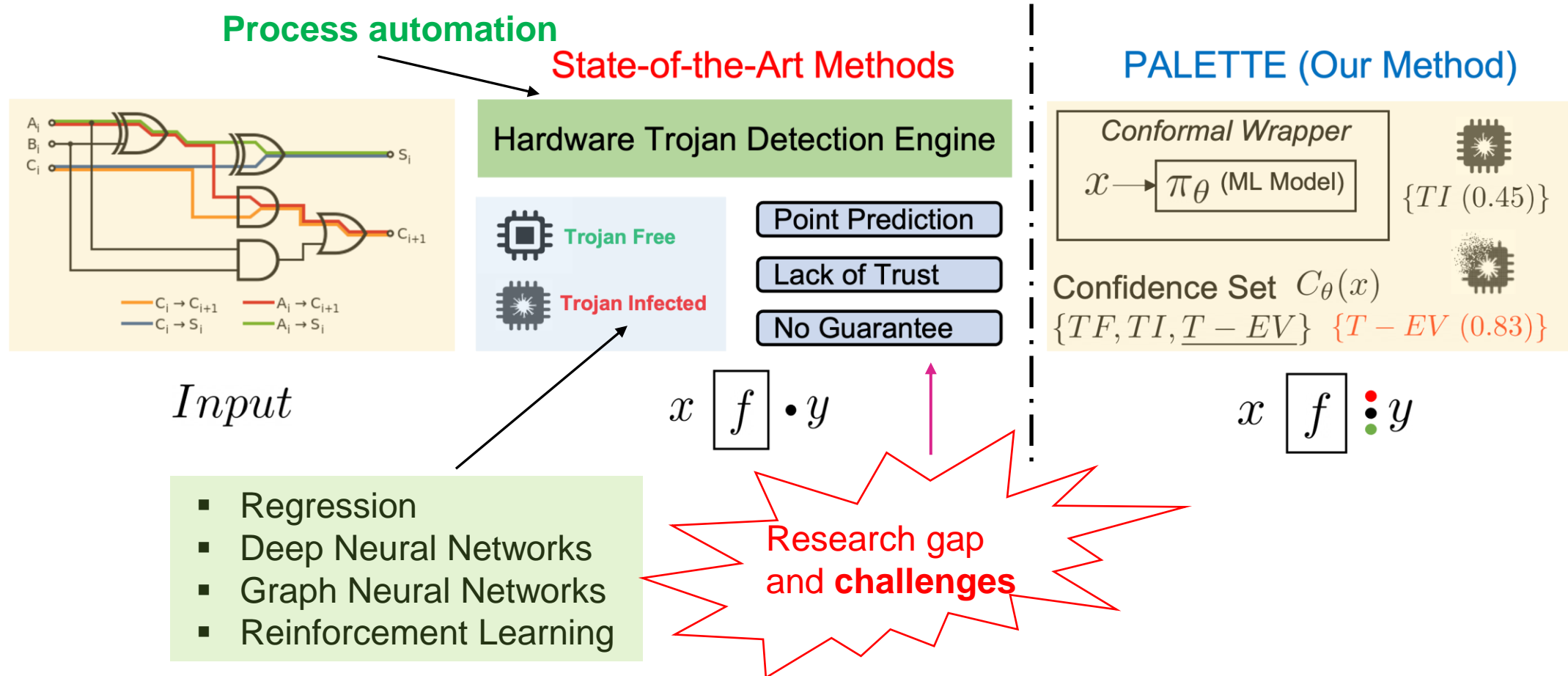
- Given a dataset $\{(x_1, y_1), \dots, (x_N, y_N)\}$
- Pick θ that minimizes the Loss

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_i L(f_\theta(x_i), y_i)$$



All models are wrong,
but some are useful!

Machine Learning based solution

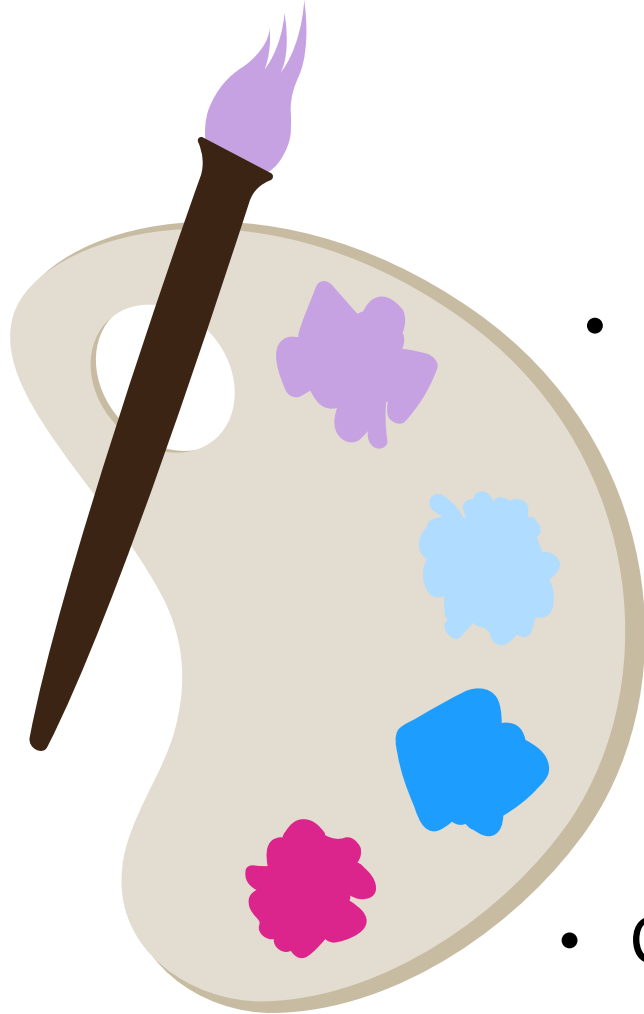




PALETTE

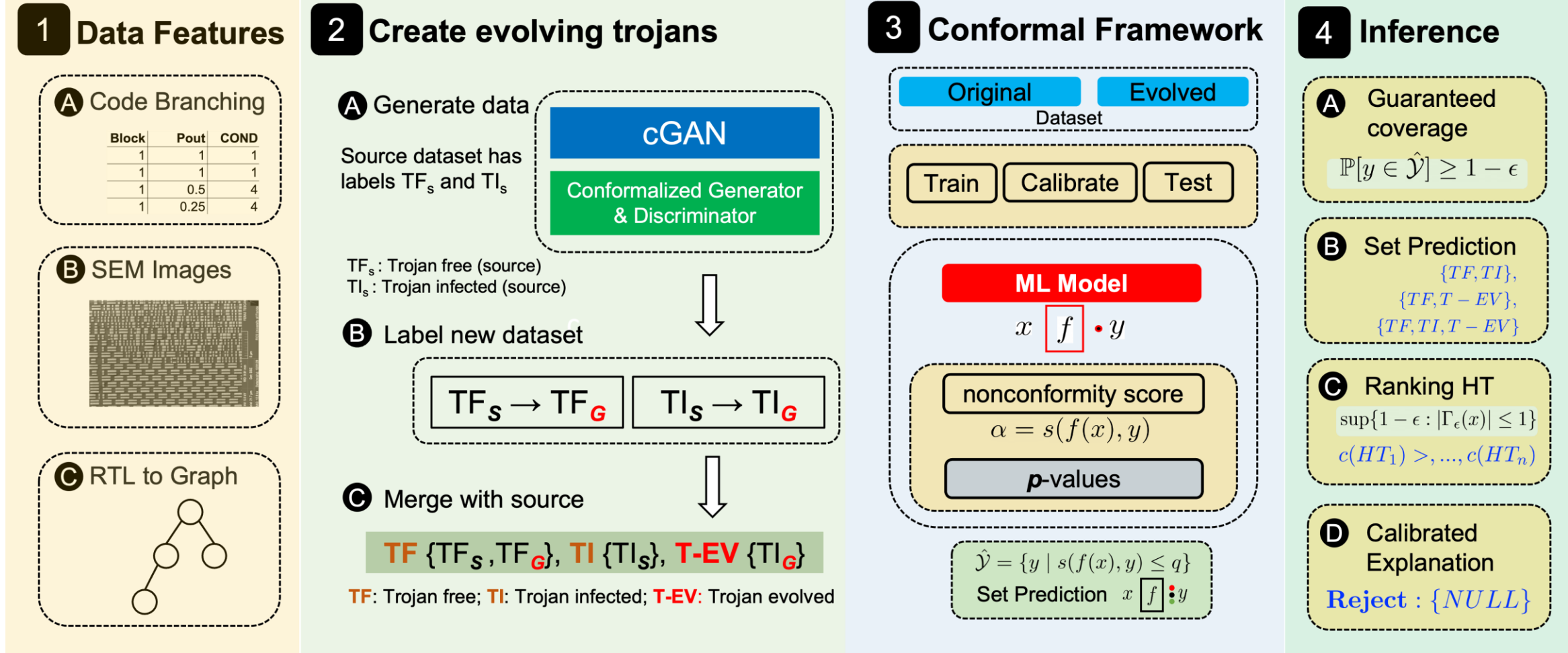
Solution and experimental results

PALETTE: exPLainable frAmework for evoLving hardwarE Trojan deTEction



- Guaranteed coverage of each prediction
 - Model says – “I don’t know”
 - Ranking of Hardware Trojans
- Calibrated explanations for each decision

Proposed Solution



Dataset

- GAINESIS¹
- Chip-level Trojan: Trust-Hub²
 - Creating Evolved Hardware Trojans

$$Label = \{TF, TI, T - EV\}$$

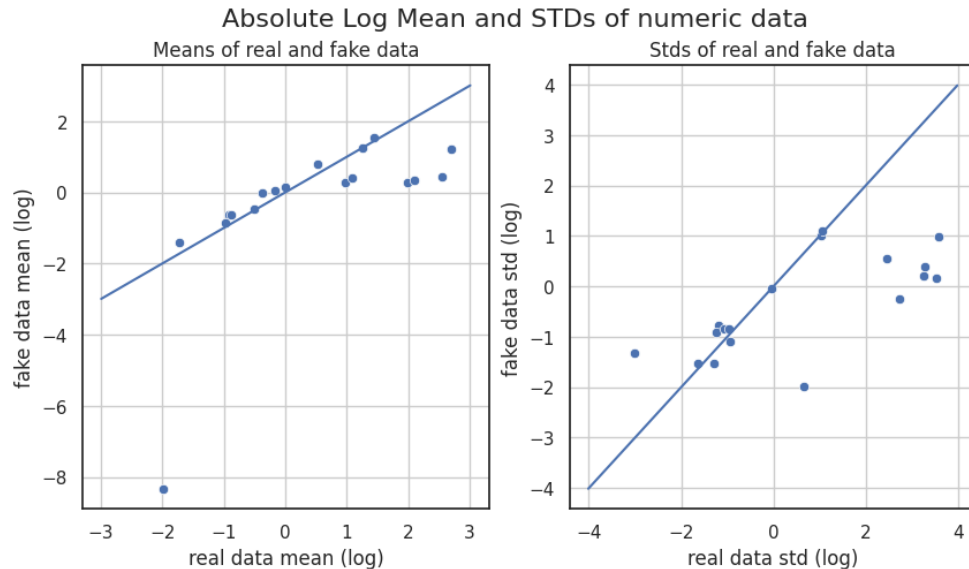
¹ K. G. Liakos, G. K. Georgakilas, F. C. Plessas, and P. Kitsos, "Gainesis: Generative artificial intelligence netlists synthesis," *Electronics*, vol. 11, no. 2, p. 245, 2022.

² H. Salmani, M. Tehranipoor, S. Sutikno, and F. Wijitrnanto, "Trust-hub trojan benchmark for hardware trojan detection model creation using machine learning," 2022. [Online]. Available: <https://dx.doi.org/10.21227/px6s-sm21>

Evolving Hardware Trojans

- Notion of Evolution:

$$E_{HT} \rightarrow HT \blacksquare HT_{structural_changes}$$



Comparison of real and synthetically generated dataset on trust-hub chip-level trojan dataset.

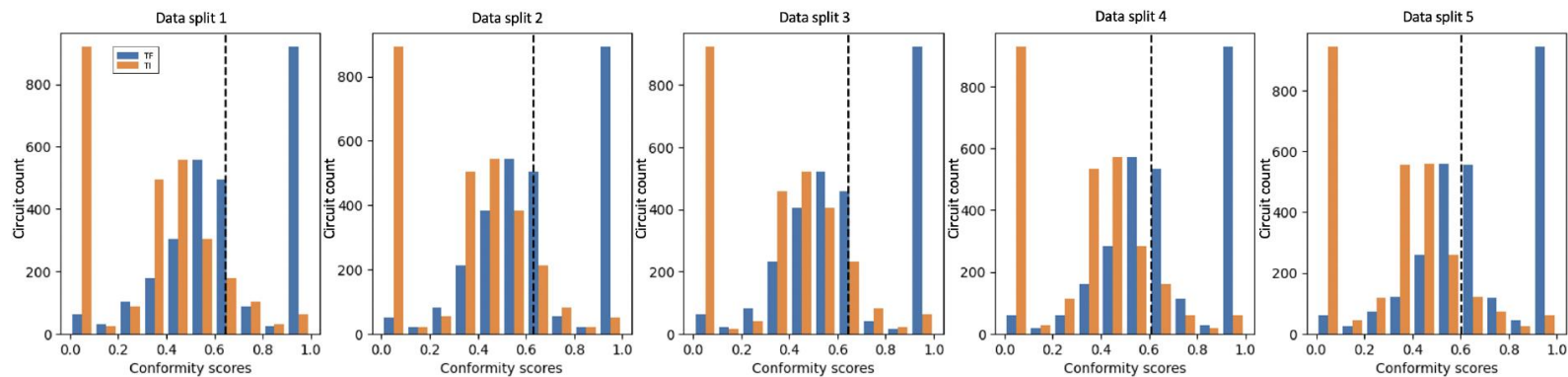
	<i>Train</i>	<i>Calibration</i>	<i>Test</i>
<i>TF</i>	1436	470	471
<i>TI</i>	114	33	44
<i>T-EV</i>	308	117	105
Total count	1858	620	620
T-EV	16.50%	18.87%	16.93%

Dataset split for model input

Experimental Results

circuit	TI	TF	y-pred	Conf
1	FALSE	TRUE	TF	0.891
2	FALSE	TRUE	TF	0.796
3	FALSE	TRUE	TF	0.996
4	FALSE	TRUE	TF	0.997
...
4596	FALSE	TRUE	TF	1
4597	FALSE	TRUE	TF	0.991
4598	TRUE	FALSE	TI	0.995
4599	FALSE	TRUE	TF	0.989
4600	FALSE	TRUE	TF	0.992

Conformal inference for GAINESIS dataset



Distribution of scores on each five of the calibration fold for GAINESIS dataset

Experimental Results

	TF	TI	T-EV	pTF	pTI	pT-EV	y_pred	Conf
1	T	F	F	0.319	0	0.003	TF	0.997
2	T	F	F	0.243	0.002	0.006	TF	0.994
3	T	T	F	0.161	0.078	0.016	TF	0.992
4	T	T	T	0.114	0.053	0.119	T-EV	0.886
5	T	F	F	0.645	0.001	0.004	TF	0.996
6	F	F	T	0.653	0	0.971	T-EV	0.365
7	T	F	F	0.3	0	0.002	TF	0.998

Conformal inference and associated p-values for trust-hub chip-level trojan dataset.

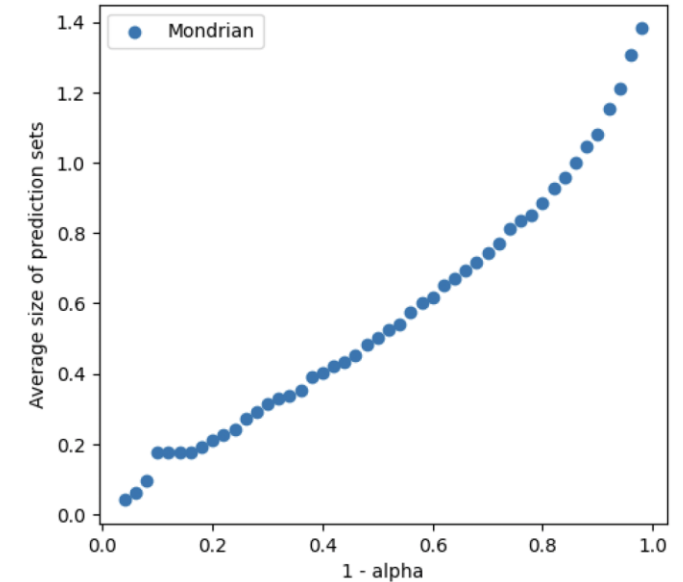
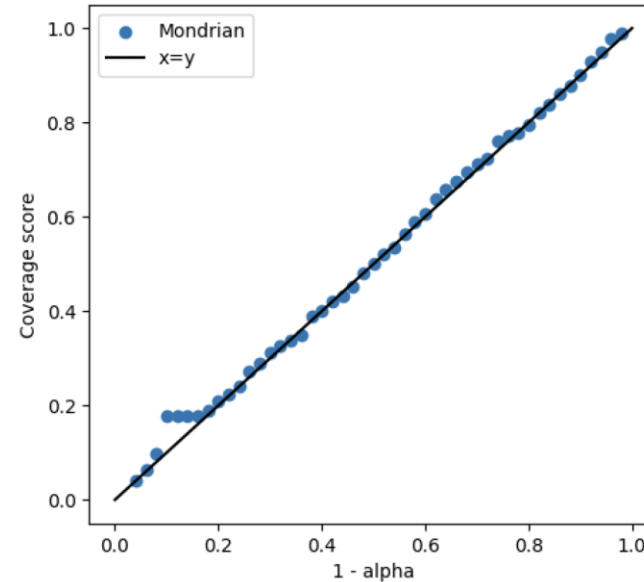
alpha	mondrian	raps	naïve	top_k
0.05	10	37	35	0
0.5	45	57	57	61
0.9	45	61	61	61

Comparison of conformal predictors with corresponding significance level on trust-hub chip-level trojan dataset

Performance Metrics

Coverage: Proportion of true target values that fall within the pred intervals.

Efficiency: How tight the prediction intervals are.



Effective coverage and average prediction set size for Trust-Hub chip dataset.

Performance Metrics

sig	mean_err	avg_c	n_correct	mean_T-EV
0.05	0.049	1.040	589	0.012
0.1	0.102	0.941	556	0.045
0.2	0.204	0.812	493	0.133
0.3	0.303	0.701	431	0.220
0.4	0.406	0.596	367	0.319
0.5	0.504	0.497	307	0.423
0.6	0.604	0.397	245	0.536
0.7	0.702	0.298	184	0.650
0.8	0.798	0.202	125	0.764
0.9	0.900	0.100	61	0.884

Performance metrics of conformal inference on trust-hub chip-level trojan dataset



Inference

Risk-Aware decision making

Risk-Aware Ranking

- Confidence score

$$\text{Confidence}(x) = \sup\{1 - \epsilon : |\Gamma_\epsilon(x)| \leq 1\}$$

- Assign confidence score

$$\alpha_{0.05}(\text{circuit 12}) = \{T - EV\}_{C=0.88}$$

$$\alpha_{0.05}(\text{circuit 13}) = \{T - EV\}_{C=0.81}$$

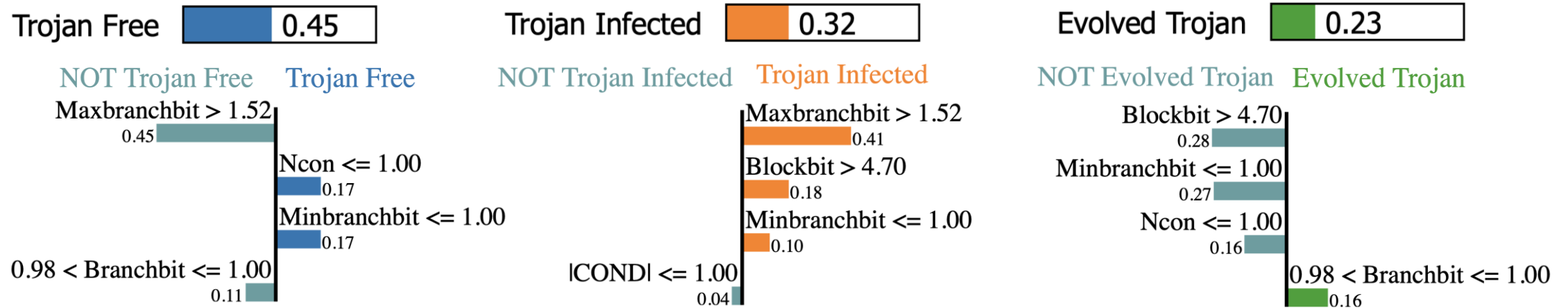
$$\alpha_{0.05}(\text{circuit 14}) = \{T - EV\}_{C=0.61}$$

	confidence	credibility	y_pred
1	0.997	0.319	TF
2	0.994	0.242	TF
3	0.922	0.162	TF
4	0.886	0.119	T-EV
5	1	0.645	TF
6	0.999	0.97	T-EV
7	0.998	0.301	TF

Adoption of confidence for risk-aware ranking on trust-hub chip-level trojan dataset

Calibrated Explanation for Reject

- Explanation for each decisions are calibrated



Calibrated explanation for rejecting a decision

Model says – “I don’t know”

Algorithm 1: Prediction with Reject Option

Input: model, instance

Output: prediction

```
1 confidence_scores = model.predict(instance);
2 if  $\max(\text{confidence\_scores}) < \text{threshold}$  then
3   | return "I don't know";
4 return class_with_highest_confidence(confidence_scores);
```

Algorithm 2: Prediction with Non-Conformity Measure Threshold

Input: model, instance

Output: prediction

```
1 non_conformity = calculate_non_conformity(model, instance);
2 if  $\text{non\_conformity} > \text{threshold}$  then
3   | return "I don't know";
4 confidence_scores = model.predict(instance);
5 return class_with_highest_confidence(confidence_scores);
```



Conclusion

Key Takeaways

- Evolving hardware Trojan (HT) evolution.
- Guaranteed coverage and tunable significance levels.
- Algorithm-agnostic and explainability-aware rejection of predictions.
- Ranking mechanism for evolved Trojans.

While there's no silver bullet for zero-day attacks, adopting a proactive risk-aware defense strategy significantly reduces the attack.



Amin Rezaei



Rahul Vishwakarma

Computer Architecture, Reliability, and Security Laboratory (CARS-Lab)



Scan me!