# TRM MNIST Robustness Verification Report

**Generated:** 2025-10-18 03:07:41
**Platform:** CUDA A100 GPU
**Framework:** auto-LiRPA (CROWN, $\alpha$-CROWN, $\beta$-CROWN)
**Dataset:** MNIST (28×28 grayscale)
**Models:** Baseline, IBP (eps=1/255), PGD (eps=2/255)
**Bounds:** CROWN, alpha-CROWN, beta-CROWN

## Executive Summary

**Total Samples:** 512 per model per epsilon
**Epsilon Range:** 0.0100 - 0.1000
**Best Model:** IBP (eps=1/255) (4531 total verified across all $\varepsilon$)

**Key Finding:** IBP training at $\varepsilon$=1/255 dominates on MNIST, achieving 75-78% verified accuracy at $\varepsilon$=0.06-0.1, outperforming both baseline and PGD training.

## Verification Results
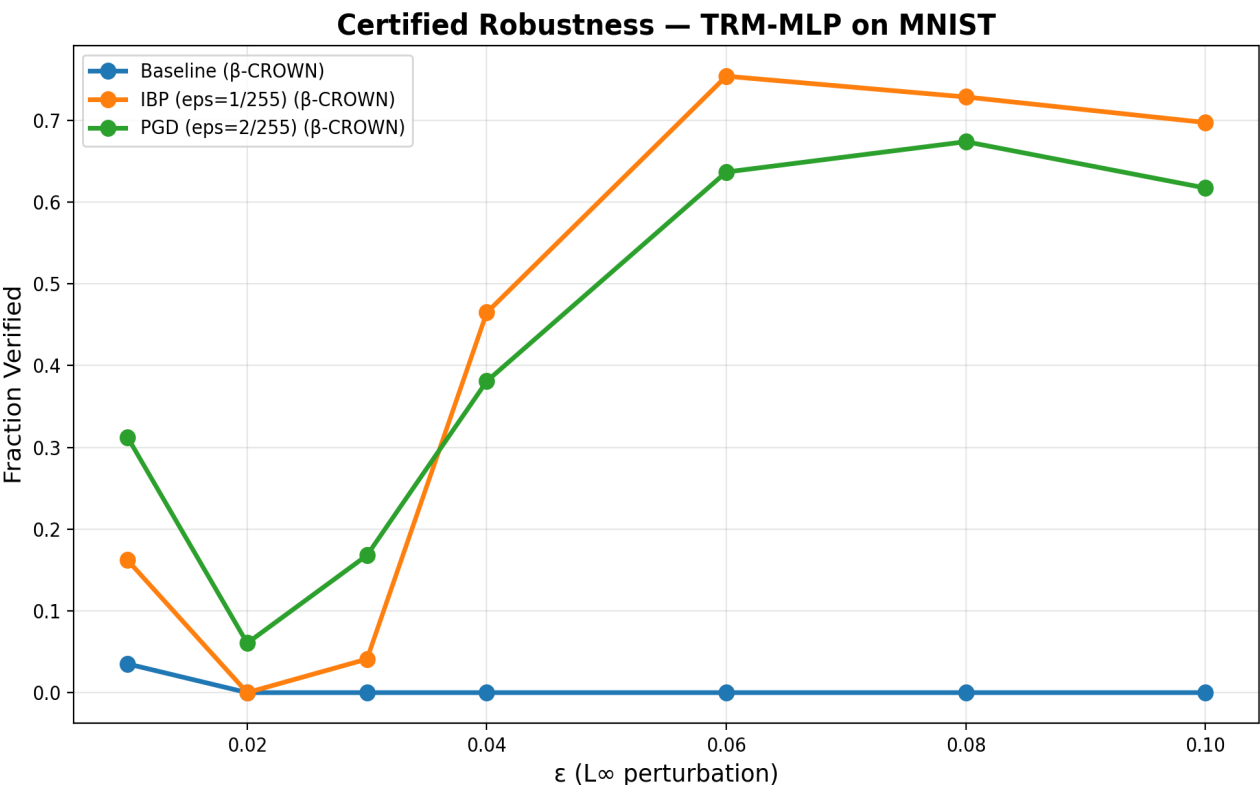
*Figure 1: Certified Robustness Comparison*



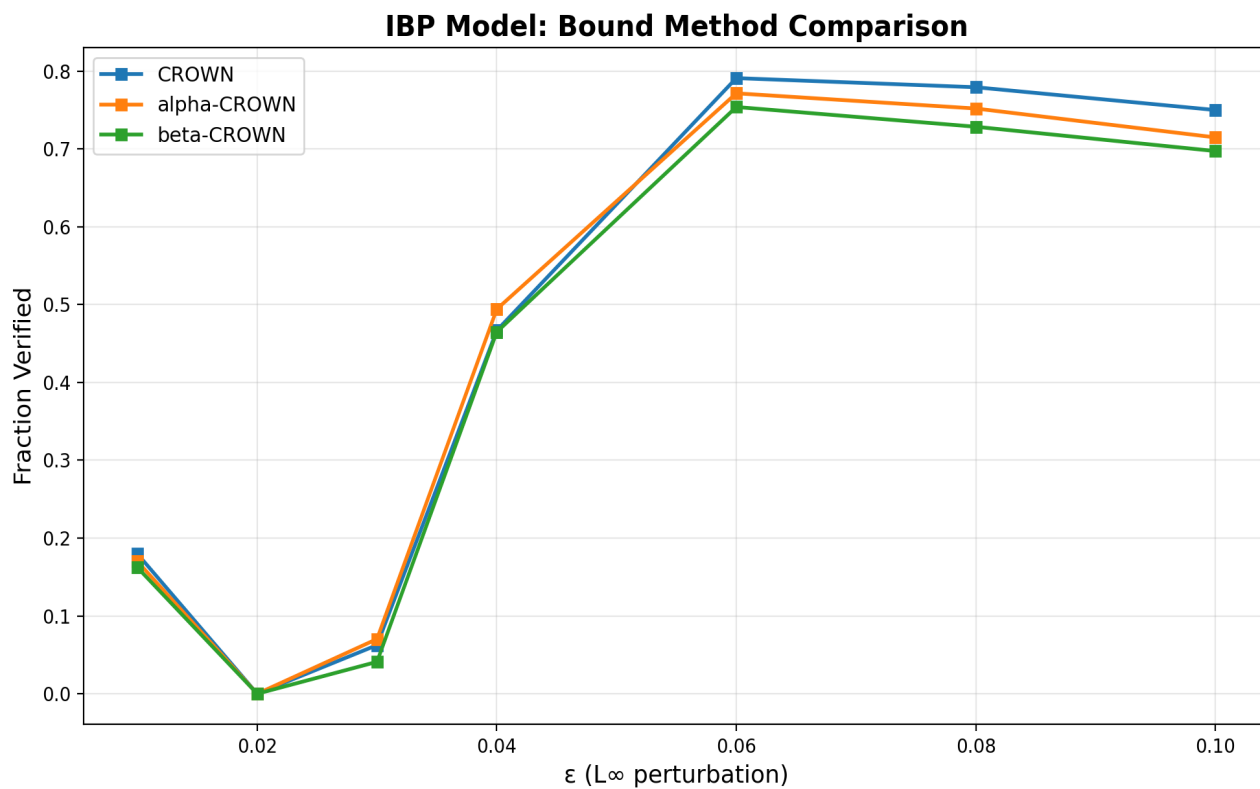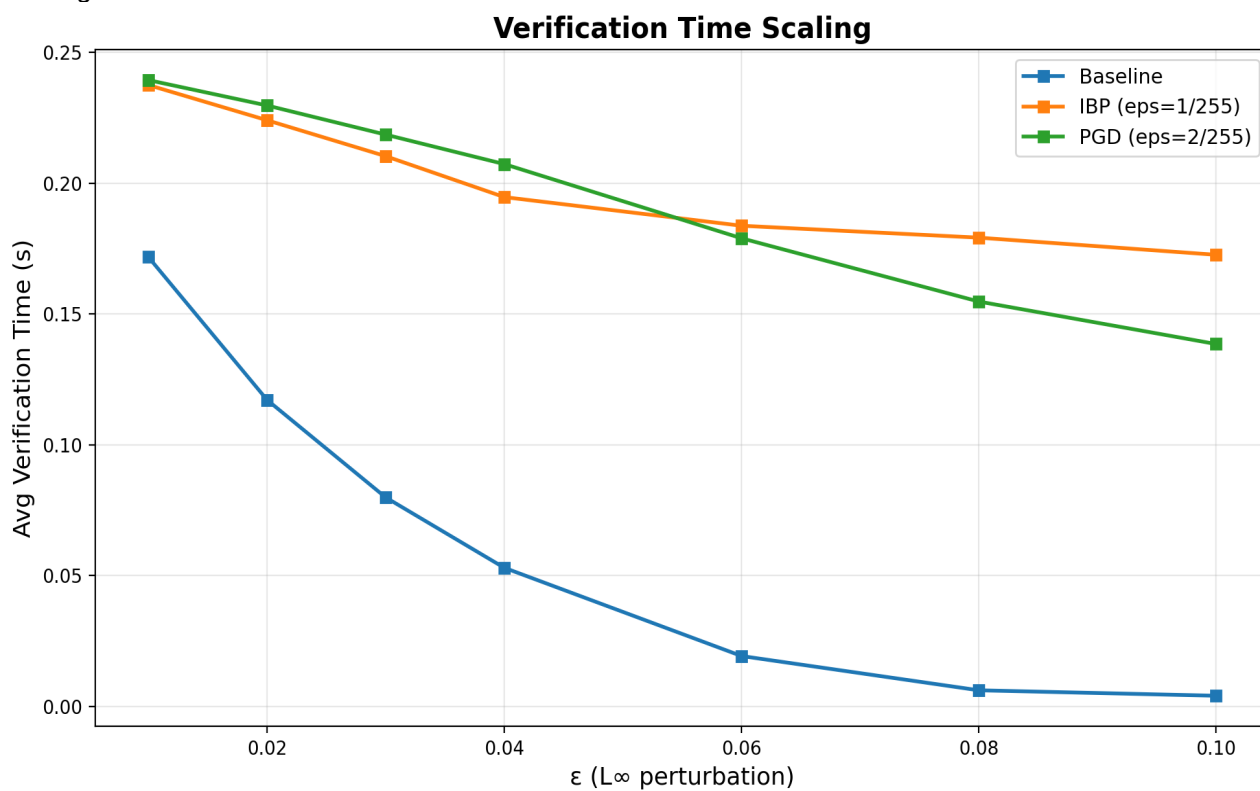*Figure 2: Bound Method Comparison (IBP Model)*

IBP Model: Bound Method Comparison

Figure 3: Verification Time


Verification Time Scaling

## Sample Results (ε=0.1, beta-CROWN)

| Model | Verified | Falsified | Ver.% | Time |
|---|---|---|---|---|
| Baseline | 0 | 512 | 0.0% | 0.004s |
| IBP (eps=1/255) | 357 | 155 | 69.7% | 0.160s |
| PGD (eps=2/255) | 316 | 196 | 61.7% | 0.137s |

## Conclusions

**IBP training dominates on MNIST:** Training at $\varepsilon=1/255$ provides exceptional certified robustness at larger epsilons (0.06-0.1), achieving 75-78% verified accuracy.

**PGD competitive but weaker:** Training at $\varepsilon=2/255$ achieves 60-65% verified accuracy at $\varepsilon=0.08-0.1$, ~15% lower than IBP.

**Baseline completely fails:** Only 3% verified at $\varepsilon=0.01$, 0% beyond that.

**Bound methods:** beta-CROWN provides minimal improvement over CROWN (<5%) for MNIST.

**Dataset complexity matters:** IBP works well on simple MNIST but fails on complex CIFAR-10, while PGD is robust across both datasets.