**The Backdoor Boys**

AAI -595 B - DR HAO WANG
Max Tuscano
Paras Jadhav
Keval Sompura

# Mitigating Data Poisoning: Detecting and Removing Malicious Outliers

DuckAI

Autonomous driving models are vulnerable to Data Poisoning attacks where malicious outliers corrupt training integrity. We demonstrate a defense using Isolation Forests to detect and cleanse these threats

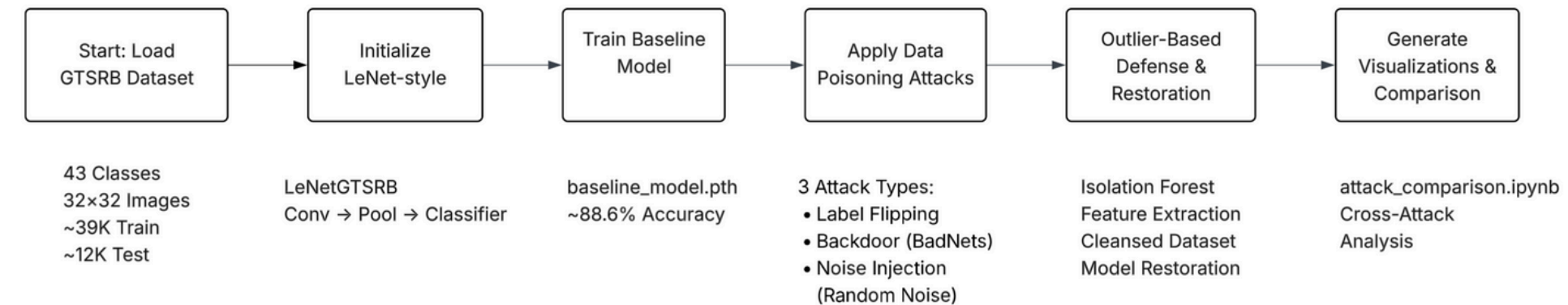## DATASET: GTSRB (German Traffic Sign Recognition Benchmark)

### PROBLEM STATEMENT

- Data Poisoning is a security attack where an adversary intentionally corrupts the data used to train a Machine Learning model.
- Models are only as good as their data. Poisoned data leads to a compromised AI, causing targeted errors in classification
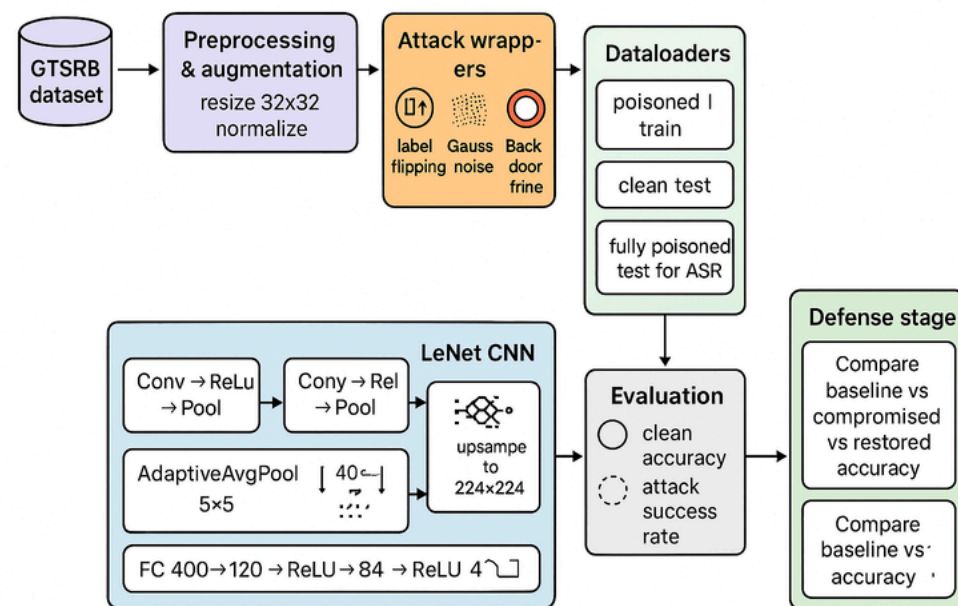
### WHY IS IT IMPORTANT

- Safety Critical: A poisoned model in a self-driving car could misinterpret a "Stop" sign, leading to catastrophic failure.
- Trust: Ensuring the integrity of training data is paramount for deploying trustworthy AI systems.
- Defense: We must quantitatively prove that compromised models can be restored.

## End to End Implementation Pipeline



| Start: Load GTSRB Dataset | Initialize LeNet-style | Train Baseline Model | Apply Data Poisoning Attacks | Outlier-Based Defense & Restoration | Generate Visualizations & Comparison |

43 Classes
32×32 Images
~39K Train
~12K Test

LeNetGTSRB
Conv → Pool → Classifier

baseline_model.pth
~88.6% Accuracy

3 Attack Types:
• Label Flipping
• Backdoor (BadNets)
• Noise Injection (Random Noise)

Isolation Forest
Feature Extraction
Cleansed Dataset
Model Restoration

attack_comparison.ipynb
Cross-Attack
Analysis

## Model Architecture



## RESULTS