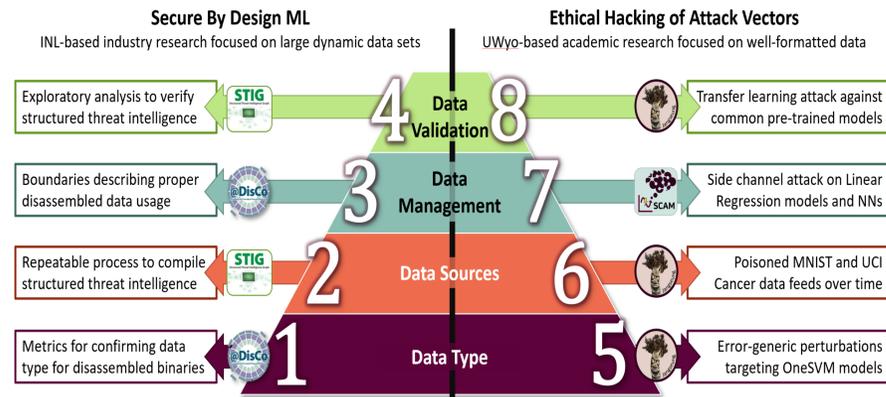# Securing Machine Learning Models for Trustworthiness
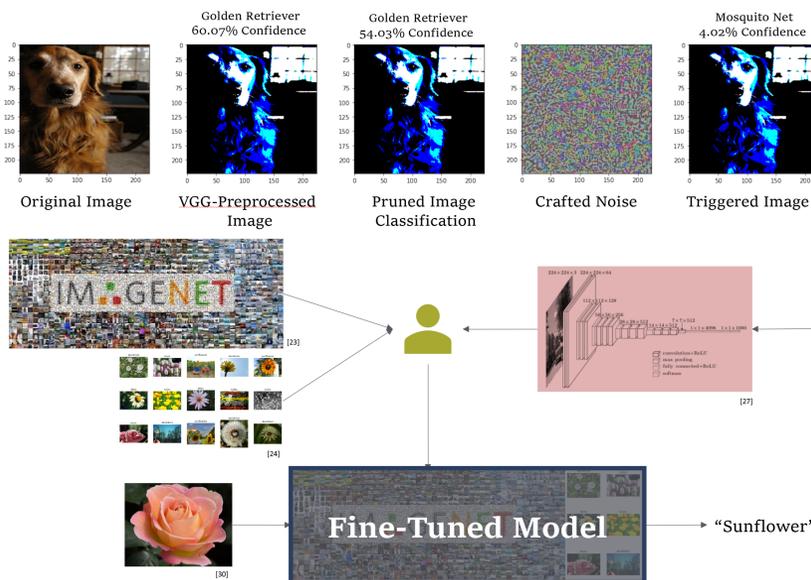## Your model said what, now?

## Abstract

Machine learning (ML) has many limitations and lacks fundamental security standards. Academic researchers and industry professionals alike aim to answer: how do we build and deploy trustworthy ML models?
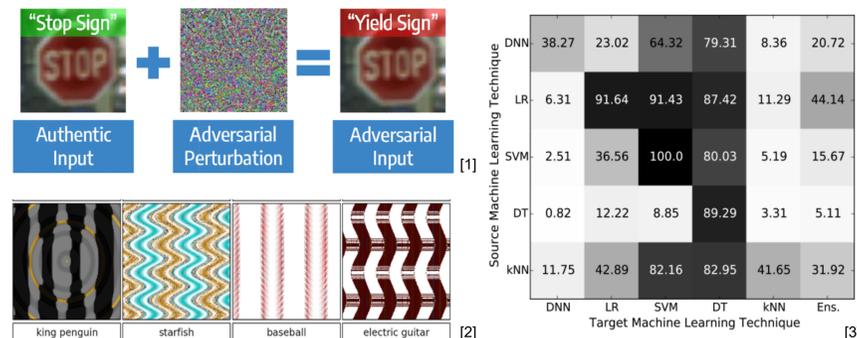
Secure By Design ML — INL-based industry research focused on large dynamic data sets

Ethical Hacking of Attack Vectors — UWyo-based academic research focused on well-formatted data

- 4 / 8 — Data Validation — Exploratory analysis to verify structured threat intelligence / Transfer learning attack against common pre-trained models
- 3 / 7 — Data Management — Boundaries describing proper disassembled data usage / Side channel attack on Linear Regression models and NNs
- 2 / 6 — Data Sources — Repeatable process to compile structured threat intelligence / Poisoned MNIST and UCI Cancer data feeds over time
- 1 / 5 — Data Type — Metrics for confirming data type for disassembled binaries / Error-generic perturbations targeting OneSVM models

## Methods

- Create pre-trained model with highly similar weights to a regular distribution, but render the model inaccurate on specific images
- Modeled after attack by Wang et. al. [0] but with a different flavor of "trigger" image.



Original Image | VGG-Preprocessed Image | Golden Retriever 60.07% Confidence / 54.03% Confidence Pruned Image Classification | Crafted Noise | Mosquito Net 4.02% Confidence Triggered Image
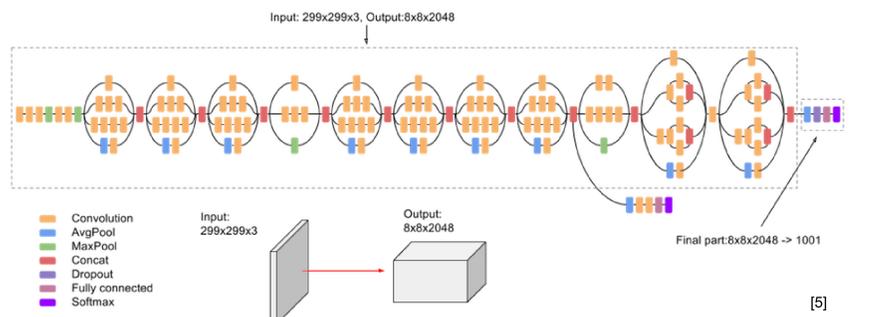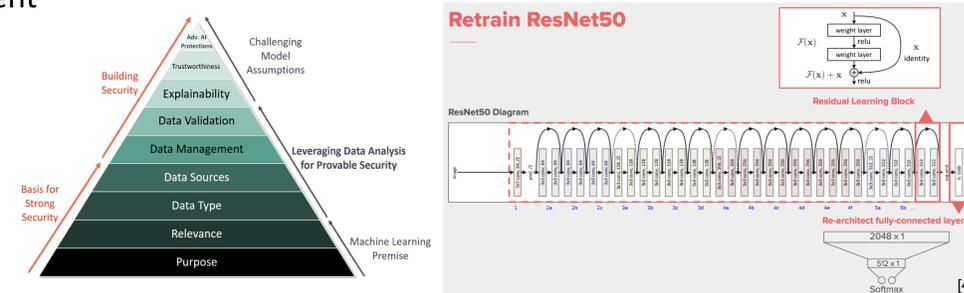


Fine-Tuned Model → "Sunflower"

## Problem Statement

- Models are rarely bench-marked on metrics other than accuracy, leaving little evidence for trust.
- ML models are easily distracted, deceived, and deluded.
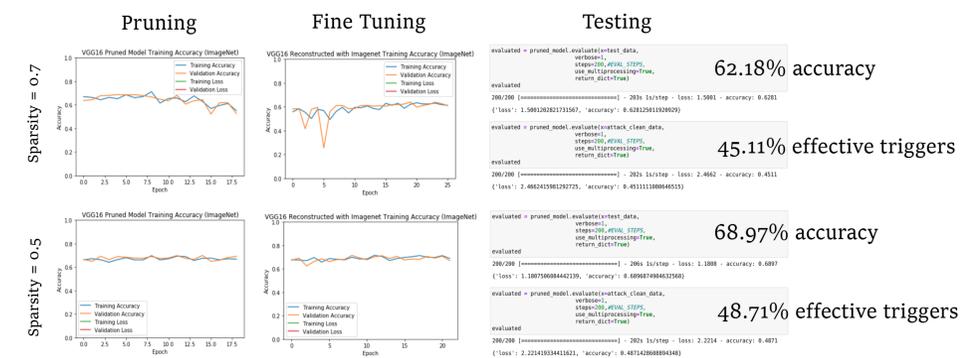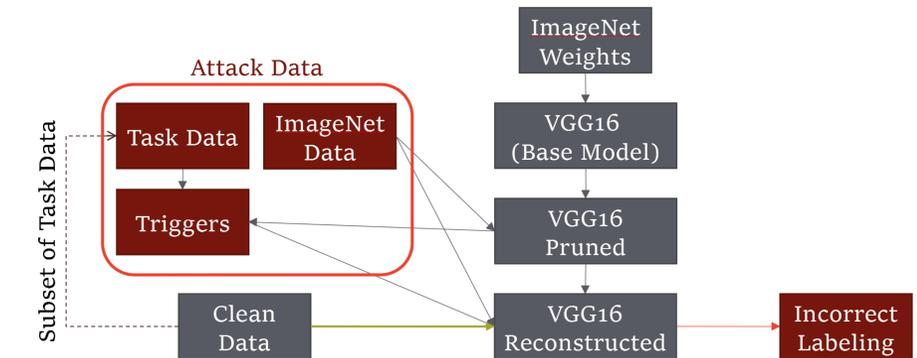- Idaho National Laboratory machine learning framework builds toward explainable and trustworthy results.



"Stop Sign" + Adversarial Perturbation = "Yield Sign" Adversarial Input

Authentic Input | Adversarial Perturbation | Adversarial Input [1]

king penguin | starfish | baseball | electric guitar [2]

| Source Machine Learning Technique | DNN | LR | SVM | DT | kNN | Ens. |
|---|---|---|---|---|---|---|
| DNN | 38.27 | 23.02 | 64.32 | 79.31 | 8.36 | 20.72 |
| LR | 6.31 | 91.64 | 91.43 | 87.42 | 11.29 | 44.14 |
| SVM | 2.51 | 36.56 | 100.0 | 80.03 | 5.19 | 15.67 |
| DT | 0.82 | 12.22 | 8.85 | 89.29 | 3.31 | 5.11 |
| kNN | 11.75 | 42.89 | 82.16 | 82.95 | 41.65 | 31.92 |

Target Machine Learning Technique [3]

## Challenges & Future Work

- Analysis is computationally heavy and time consuming
- Extend to further datasets and pre-trained models
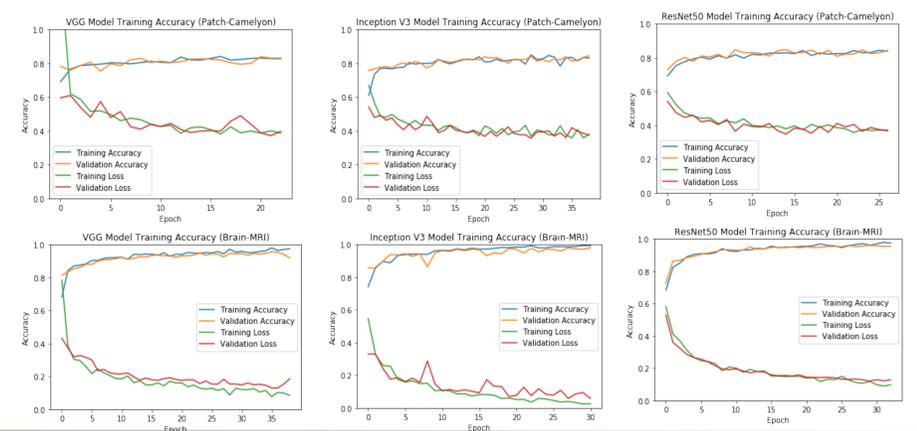- Those models must first be baselined (already completed)



Retrain ResNet50

## Results

- Approximately half of the triggers are effective
- Attack survives pruning, fine-tuning and drop-out layers



Attack Data: Task Data, ImageNet Data, Triggers

ImageNet Weights → VGG16 (Base Model) → VGG16 Pruned → VGG16 Reconstructed → Incorrect Labeling

Subset of Task Data → Clean Data

Pruning | Fine Tuning | Testing

Sparsity = 0.7 — 62.18% accuracy / 45.11% effective triggers

Sparsity = 0.5 — 68.97% accuracy / 48.71% effective triggers

| Accuracy | 20% Dropout Rate | 30% Dropout Rate |
|---|---|---|
| No Triggers | 66.03% | 63.78% |
| Only Triggers | 43.39% | 44.95% |



VGG Model Training Accuracy (Patch-Camelyon) | Inception V3 Model Training Accuracy (Patch-Camelyon) | ResNet50 Model Training Accuracy (Patch-Camelyon)

VGG Model Training Accuracy (Brain-MRI) | Inception V3 Model Training Accuracy (Brain-MRI) | ResNet50 Model Training Accuracy (Brain-MRI)

**Advisor**: Dr. Mike Borowczak
**Authors**: Shaya Wolf
Support from INL: Rita Foster, Jed Haile

College of Engineering and Physical Sciences
Cybersecurity Education and Research Center
School of Computing
College of Engineering and Physical Sciences — Electrical Engineering and Computer Science
CEDAR

[0] S. Wang, S. Nepal, C. Rudolph, M. Grobler, S. Chen and T. Chen, "Backdoor Attacks Against Transfer Learning With Pre-Trained Deep Learning Models," in IEEE Transactions on Services Computing, vol. 15, no. 3, pp. 1526-1539, 1 May-June 2022, doi: 10.1109/TSC.2020.3000900.
[1] https://miro.medium.com/max/875/3a6f8cac_3eqd6r26cu.png
[2] https://arxiv.org/abs/1412.1897
[3] https://www.arxiv-vanity.com/papers/1605.07277/
[4] https://i.stack.imgur.com/gI4zT.png
[5] https://cloud.google.com/tpu/docs/images/inceptionv3onc--oview.png

University of Wyoming