# Complete Evaluation on Advanced Reactor Machine Learning Subversion Attacks

*A guide for Advanced Reactor Architects, Nuclear Regulators, and Cyber Defense Teams*

September 2023

*M2CT-23IN1105017*

Idaho National Laboratory
*Chris Spirito*

Idaho State University
*Dr. Leslie Kerby (PI)*
*Dr. Pedro Mena (PI)*

*Senior Research Staff*
*Eric Hill, Kallie McLaren, and Samantha Ross*

*Junior Research Staff*
*Dunya Bahar, Gabriel Ory, Gunnar Babicz, Hyunsu Kim, Nehal Hasnaeen, Sansar Kharal, Sindi Banda, and Zackary Schepis*

**INL** Idaho National Laboratory

*INL is a U.S. Department of Energy National Laboratory operated by Batelle Energy Alliance, LLC*

# Complete Evaluation on Advanced Reactor Machine Learning Subversion Attacks

## M2CT-23IN1105017

**Idaho National Laboratory**
*Chris Spirito*

**Idaho State University**
*Dr. Leslie Kerby (PI)*
*Dr. Pedro Mena (PI)*

*Senior Research Staff*
*Eric Hill, Kallie McLaren, and Samantha Ross*

*Junior Research Staff*
*Dunya Bahar, Gabriel Ory, Gunnar Babicz, Hyunsu Kim, Nehal Hasnaeen,*
*Sansar Kharal, Sindi Banda, and Zackary Schepis*

**September 2023**

*Page intentionally left blank*

*Page intentionally left blank.*

# CONTENTS

# FIGURES

*Page intentionally left blank.*

# ACRONYMS

| | |
|---|---|
| AI | Artificial intelligence |
| API | Application Programming Interface |
| ASI | Advanced Sensor and Instrumentation |
| BWR | Boiling water reactor |
| CLI | Command line interface |
| DCS | Distributed Control System |
| DT | Digital Twin |
| DTDL | Digital Twin Definition Language |
| HMI | Human-Machine Interface |
| IAEA | International Atomic Energy Agency |
| IIoT | Industrial Internet of Things |
| I&C | Instrumentation and Control |
| JADE | Java application and developer environment |
| MCO | Moisture Carryover |
| ML | Machine learning |
| OT | Operational Technology |
| PWR | Pressurized water reactor |
| RBAC | Role-based access control |
| RDF | Resource Description Framework |
| REST | Representational state transfer |
| SDK | Software Development Kit |
| SLA | Service Level Agreement |
| SMR | Small Modular Reactor |
| XML | Extensible markup language |

*Page intentionally left blank.*

# 1. Introduction

Navigating through the world of Artificial Intelligence (AI) in nuclear reactors and their Instrumentation and Control (I&C) systems demands a careful, deliberate journey. AI's capability to manage massive datasets and streamline control systems has indeed carved out a significant role in various sectors, including nuclear energy. However, while AI, and particularly Large Language Models (LLMs), bring a lot to the table in terms of operational efficiency and anomaly detection, they also expose the sector to a new breed of cybersecurity threats, like Inference Attacks, Adversarial Attacks, and Trojan Attacks.

This guide is designed to be a straightforward manual, diving deep into the intertwining worlds of AI and cybersecurity within nuclear reactors, and tailoring insights for three crucial audiences: I&C Vendors/Developers, Nuclear Regulators, and Nuclear Reactor Operators and Cyber Defense Teams.

- **Section 2**, directed at I&C Vendors/Developers, will provide a clear and focused look at several cybersecurity attacks, offering practical recommendations and detailed scenarios related to AI cybersecurity. This section isn't just about identifying problems but also about giving solid, usable solutions.
- **Section 3**, meant for Nuclear Regulators, gets straight to the point about regulations, policy suggestions, and guidelines that are needed to lay down a robust, secure, and ethical foundation for the application of AI in nuclear operations. The focus is on making sure that everything adheres to international standards and laws while being practicable and clear-cut.
- **Section 4**, aimed at Nuclear Reactor Operators and Cyber Defense Teams, offers an exhaustive exploration and technical reports, with clear recommendations and scenario analyses vital to protect operational environments and guarantee the secure application of AI in nuclear reactor operations.

The goal is simple: as we step into an era where AI becomes a fundamental element of our technological and energy infrastructures, this guide is here to act as a clear, direct handbook, ensuring that AI is implemented within the nuclear sector in a manner that is secure, responsible, and practical. It's about striking a balance – optimizing the undeniable benefits offered by AI while securing and shielding against potential cyber threats as we move through this new and complex landscape.

# 2. A Guide for Advanced Reactor & I&C Venders/Developers
## 2.1 Recommendations for Developers
### 2.1.1 Inference Attacks

**Prioritize Data Control**

Model extraction is a very feasible attack and is made more simple and more accessible to adversaries with the use of AutoML. Very limited knowledge is required about the original data, and the attack can be carried out with little knowledge of machine learning. There are very few current mitigation methods for this attack. Current research into defenses and mitigation for this attack is relatively new and sparse, requiring more exploration and experimentation. Regulators and developers should be aware of this danger. Data used in machine learning model training should be considered highly sensitive. Developers need to implement controls to manage who has access to the data used to train these models learning models. Developers must also institute strong policies that manage how data is stored, collected, and ultimately used by staff working on machine learning projects. These policies will aid not only in the prevention of leaked models, but also help reduce the possibility of being attacked by a bad actor.

### 2.1.2 Large Language Models

**Advance Research in Application of LLMs for Nuclear Power Plants**

The potential applications of LLMs in nuclear power plants include predictive maintenance, operational optimization, incident analysis, safety, training, simulation, decision support, and regulatory compliance. LLMs could provide an excellent interface to these areas, as such, more research in these areas is needed. For instance, LLMs could be used to analyze historical sensor data, maintenance logs, and operational parameters to predict potential equipment failures, thereby preventing unplanned downtime. They could also be used to optimize plant performance, identify areas for process optimization, and provide insights into energy efficiency. Furthermore, LLMs could assist in incident analysis, help train personnel on various operational and emergency procedures, and assist in decision-making processes related to resource allocation, outage planning, and risk management.

### 2.1.3 Adversarial Attacks

**Evaluate Data Prior to Training**

Developers should strive to verify the integrity of their data prior to model training. When data is being collected for initial model training developers should institute policies and use tools to monitor changes to the dataset as the data is compiled. One useful tool for this may be the use of block chain to control changes and allow for easy tracking of changes. If the developers are in the process of updating a model with new data, it may be useful to also evaluate to test the data for anomalies. Once such tool for this would be an autoencoder neural network.

**Versioning and Maintaining Records**

In software development, database management, etc. versioning is considered a good practice. Developers working on AI/ML models should also make use of this methodology to help mitigate problems associated with adversarial attacks. Not only would this practice allow the developer to deploy older versions of models in the event of a breach, but the practice would also help cyber teams trace back issues to better identify when and where they have occurred.

### 2.1.4    Trojan Attacks

**Continuously Retrain Models**

Due to the ability to cycle out a Trojan from an infected machine learning model, developers should continuously retrain models prior to implementation. This will improve the chances that a Trojan is ineffective in the event of a breach.

**Monitor Model Weights**

Developers should continuously monitor the weights of models are they are trained and re-trained. Early results in this project have shown that there are a number of statistical differences in the weights of benign models and those that have been infected by a Trojan. Analysis of the weights throughout the development process will increase the chances that a model infected with a Trojan is detected.

## 2.2  Technical Reports

### 2.2.1    Datasets Used

For this project it was necessary to use a number of different datasets to perform many of the different attacks. The following sections provide detail on the datasets that were used in this effort.

#### 2.2.1.1    Image Datasets

To help evaluate attacks and defenses for this project three image datasets were used to provide baselines. These image sets are part of the standard TensorFlow library and are commonly used in data science research. The MNIST Digits is a dataset of 70,000 grayscale images of handwritten numbers from zero to nine, where each image is 28x28. There are 60,000 total training images and 10,000 total testing images, with 7,000 images from each class. The Fashion MNIST is also a dataset of 70,000 grayscale 28x28 images with 60,000 train images and 10,000 test images. However, this dataset consists of 10 different articles of clothing. These articles of clothing are t-shirts, trousers, pullovers, dresses, coats, sandals, shirts, sneakers, bags, and ankle boots. Finally, Cifar10 is a dataset consisting of 60,000 color images of 10 different categories. These categories are airplanes, automobiles, birds, cats, deer, dogs, frogs, horses, ships, and trucks. All the images are of size 32x32, and there are 50,000 total train images and 10,000 total test images. For each category, there are 6,000 images.

#### 2.2.1.2    Nuclear Datasets

For this project two datasets focusing on nuclear reactors were used for prototyping purposes. These were the Asherah dataset and GPWR dataset. The Asherah dataset was collected and provided by Patience Lamb of the Georgia Tech Team and can be found at the following Gitlab repository [1] [2]. This dataset is a binary dataset that distinguishes transient data from steady state. This data was obtained from the IAEA virtual Asherah Simulator and contains over 90 features. A sample of the Asherah dataset features can be seen in Table 1.

The GPWR dataset was acquired by Dr. Pedro Mena from the GPWR simulator managed by Idaho State University/University of Idaho. The dataset contains over 110,000 data points consisting of 27 features. The dataset pertains to 11 different transient events and normal operations, typically considered steady state. Table 2 lists a sample of features from the dataset. This data has been used in a number of publications.

Table 1: Sample of Asherah Dataset Features Used

| Asherah Dataset Features | |
|---|---|
| CE Pump1Flow | CE Pump1Speed |
| CE Pump3Temp | CE Pump3Speed |
| CR Position | FW Pump1DiffPress |
| FW Pump2Temp | FW Pump3DiffPress |

Table 2: Sample of GPWR Dataset Features Used

| GPWR Dataset Features | |
|---|---|
| Containment Pressure (PSI) | Containment Temperature (F) |
| SG-1 Pressure (PSI) | SG-2 Pressure (PSI) |
| Average Temperature (F) | Pressurize Pressure (PSI) |
| Generator Power (MW) | Reactor Core Life |

## 2.2.2    Inference Attacks

This gray box functionality stealing approach gains access to a subset of the original data, obtains the victim model's corresponding predictions on that subset, and then uses that combination of subset and predictions to train a new model. The overview of this process is detailed in Figure 1, where 'D*' denotes the subset of the original dataset 'D', and the target model's predictions on D* are denoted as 'P.'



Figure 1: Illustration for inference attack training flow

### 2.2.2.1    Model Stealing with AutoKeras

The methodology of model extraction involves (1) obtaining the necessary adversarial training data for the AutoML, (2) generating a stolen model, and (3) evaluating the accuracy of that model and therefore the efficacy of the inference attack. A process flow is included in **Error! Reference source not found.**.

To obtain the necessary adversarial training data, a subset was taken from the original dataset. This subset was images, separated from their corresponding labels, combined with the target model's predictions on those images. In other words, the AutoML was given a subset of both images and the victim's predictions to train and generate a new model. This approach is considered a gray box attack since it requires some access, although limited. It is assumed that the adversary has access to some subset of the original data and has knowledge of the image size and pre-processing information.

The training procedure for this gray box inference attack is detailed in **Error! Reference source not found.**. Two subsets, DS and DT, are taken from the original dataset. The subset DT is used to train the target model and the subset DS is an untouched reserve that is used to evaluate the target model. In order to avoid training bias, this untouched reserve of DS is used to create a train-test-split for the attack model, in which the images are separated from the labels and instead combined with prediction labels from the target model to create the adversarial training data that AutoKeras will use to generate an attack model. Since DS is used as train-test-split for the attack model, it is expected that the attack model will perform better than the target model. This is why the attack model is evaluated on both DS and DT, where its evaluation on DT will demonstrate how well the attack model generalizes to a (usually) larger subset it has never



*Figure 2: Illustration for inference attack process flow*

associated with before.



*Figure 3: Illustration for inference attack training flow*

This inference attack is relatively simple to execute, aside from requiring a GPU and high RAM, and requires little knowledge of neural networks and almost no prior knowledge of AutoML. **Error! Reference source not found.** shows the roughly 6 lines of code necessary to install the AutoML package, train

AutoKeras, and export the best-generated model. The only thing not detailed is how an attack's efficacy is evaluated, which will be addressed in **Error! Reference source not found.**.

## Simplicity of Attack Setup

1. Import AutoML

```
!pip install autokeras

import autokeras as ak
```

2. Get data for AutoML
   Create train/test split from subset of images and the victim's predictions on those images

```
# split victim test data from X_test and y_pred
# Use X_test and y_pred from victim model
X_train_tr2, X_test2, y_train_tr2, y_test2 = train_test_split(X_test, y_pred,
                                                              test_size=1/10)
```

3. Initialize and Train AutoML

```
# initialize AK image classifier that will try 2 models (for time)
auto_model = ak.ImageClassifier(max_trials=3, loss = 'sparse_categorical_crossentropy') # specifying loss

# Search for the best model
am_history2 = auto_model.fit(X_train_tr2, y_train_tr2, 30,
                 callbacks=[keras.callbacks.EarlyStopping(patience=5)])
```

4. Export best AutoML-generated model

```
# Export best model (doing early in case of issues)
stolen_model = auto_model.export_model()
```

*Figure 4: Illustration for inference attack evaluation process*

**Error! Reference source not found.** shows a high-level view of the evaluation process. Both the target and the AutoKeras-generated attack model will produce output prediction vectors on a given subset, and the evaluation will involve comparing these prediction vectors to determine if the attack was successful (i.e., the attack model has achieved" functional equivalence").



*Figure 5: Illustration for inference attack simplicity*

### 2.2.2.2 Functional Equivalence

In this research, "functional equivalence" is defined by similar model behavior, i.e., how well the attack model can replicate the target model's predictions on a given subset of data. Depending on the researcher's goals and standards, the conditions for two models being deemed functionally equivalent can be difficult to measure. How similar must two models be to be considered functionally equivalent? Current research seems to favor assessing functional equivalence via confusion matrices. However, this visual representation can be lacking for those desiring more quantitative results. The outcome of functional equivalence measurement can also differ depending on whether the attack model's desired behavior is assessed on the model accuracy of the original data or if the attack model's behavior is assessed on its similarity to the target model's predictions. In this paper, functional equivalence is assessed by both. The method in this paper for functional equivalence evaluation utilizes 5 metrics. These metrics compare the two model's predictions, both the rounded prediction vectors and the prediction probability matrices, to assess functional equivalence.

### 2.2.2.3 Functional Equivalence Metrics

The metrics used to assess functional equivalence are: (1) Match Angle Percentage, (2) Angle Difference, (3) Mean Squared Error, (4) Mean Absolute Error, and (5) Categorical Cross-Entropy. These metrics accompany analyzing confusion matrices of the model predictions.

**Match Angle Percentage (MAP/MP)**

Match Angle Percentage (MAP or MP for Match Percentage) is a method of comparing two integer prediction vectors for similarity, where more similarities will result in a higher accuracy score. MP can be used to assess a model's rounded predictions versus the actual labels (called "raw accuracy" in this paper). MP can also be used to compare two rounded prediction vectors from two separate models (called "attack accuracy" in this paper). Although it may seem redundant, MP can be different than model accuracy when it measures attack model performance. When an attack model is trained, it uses a train-test-split of DS using the target model's predictions as the labels, giving it a bias towards predicting the target model and not the actual labels. Historically, in this paper's research, the MP will most likely be the same as the attack model's raw accuracy against the actual label.



*Figure 6: Illustration for prediction vector match accuracy (aka "match percentage") and example*

7

**Error! Reference source not found.** illustrates an example of determining the raw accuracy and the attack accuracy of two rounded prediction vectors from a hypothetical target and attack model. It demonstrates an example in which, despite the attack model being only 60% accurate against the original labels, the predictions, i.e., behavior, of the attack model matched the target with 90% accuracy.

**Confusion Matrices Comparison**

A confusion matrix illustrates a model's correct and incorrect classifications, with the frequency of correct classifications being on the diagonal. It's helpful when visually assessing model classification behavior, and this makes it helpful when assessing functional equivalence since it is possible to compare the classification behavior of the two models. The diagonal represents the true predictions and the rest represent incorrect classifications. For this paper, these confusion matrices will visually represent the true versus incorrect classifications of the target and attack model predictions when com-pared against the actual true labels (to assess raw accuracy) and against the other model's predictions (to assess attack accuracy).

**Angle Difference (AD)**

Angle Difference (AD) is a method of comparison that calculates the degree by which two vectors deviate from one another. AD normalizes the input vectors via the Frobenius Norm and calculates the angle difference using the dot product, converting the result from radians to degrees. Figure 7 illustrates the Python and numpy code used for these calculations, and **Error! Reference source not found.** illustrates a more step-by-step visual representation of what the code is doing.

```
def angle_between_vectors(v1, v2):
    cos_of_angle = np.dot(v1, v2) / (np.linalg.norm(v1) * np.linalg.norm(v2))
    return np.arccos(cos_of_angle)

def rad_to_degrees(rad):
  return np.degrees(rad)
```

*Figure 7: Python numpy code for calculating angle difference that is visually easier to compare to original formula.*



Angle Difference

Normalize The Vectors Using Frobenius Norm:

$$||A||_F = \left[ E_{i,j} \; abs \; (a_{i,j})^2 \right]^{\frac{1}{2}}$$

Calculate Angle Difference Using Dot Product:

$$\theta = \frac{\cos^{-1}[(a \cdot b)]}{|a||b|}$$

Convert From Radians To Degrees:

$$\frac{180x}{\pi}$$

*Figure 8: Illustration for prediction vector angle difference*

**Mean Squared Error (MSE) & Mean Absolute Error (MAE)**

Mean Squared Error (MSE) is a risk function measuring the average squared difference between the estimated values and the actual values, in which a lower MSE indicates a better fit. The equation for MSE is listed below, where Yi is the prediction, Yi is the observed value, and n is the total number of data points.

Since MSE squares the differences between predicted and observed values, it is more sensitive to outliers. This means that if a model only had a few errors, but these errors were large, the MSE may be high despite the model fitting most of the data. For this reason, Mean Absolute Error (MAE) was included as a metric to give an idea of the average difference without extreme outliers. MAE measures data variability by measuring the average absolute deviations, in which a lower MAE indicates a better fit. The equation for MAE is listed below, where $y_i$ is the prediction, xi is the observed value, and n is the total number of data points.

**Categorical Cross-Entropy (CCE)**

Entropy itself is a measure of uncertainty, where very certain outcomes will have a low entropy. Categorical Cross-Entropy (CCE) is an activation function (usually a common loss function also called "SoftMax loss") that is used for multi-class classification when each in-put sample can belong to multiple classes. CCE will output a probability for a given sample over each class, creating a prediction probability matrix, where all probabilities will sum to 1. A lower CCE indicates a better fit. **Error! Reference source not found.** gives the formula for CCE loss, where s is a vector of predictions and t is the target vector over C images. **Error! Reference source not found.** shows the Python NumPy code used to calculate the CCE.



$$f(s)_i = \frac{e^{s_i}}{\sum_j^C e^{s_j}} \qquad CE = -\sum_i^C t_i log(f(s)_i)$$

*Figure 9: Formula for CCE loss*

```
# Use Categorical Crossentropy (CCE) function in Keras
cce = tf.keras.losses.CategoricalCrossentropy()

# Use CCE on entire prediction probability matrices of target/attack models
# for both DS and DT subsets
CCEarray_ds = np.asarray([cce(pred1, pred2) for (pred1, pred2) in zip(target_preds, attack_preds)]) # on DS
CCEarray_dt = np.asarray([cce(pred1, pred2) for (pred1, pred2) in zip(target_preds_train, attack_preds_train)]) # on DT

# Print results of average CCE performance/deviation and max CCE deviation
print("On DS: ")
average_deviation = np.mean(CCEarray_ds)
print("Avg Dev: ", average_deviation)
max_deviation = np.max(CCEarray_ds)
print("Max Dev: ", max_deviation)
print()
print("On DT: ")
average_deviation = np.mean(CCEarray_dt)
print("Avg Dev: ", average_deviation)
max_deviation = np.max(CCEarray_dt)
print("Max Dev: ", max_deviation)
```

*Figure 10: Python code for calculating CCE mean/max scores.*

### *2.2.2.4    Results*

The results will focus on four specific models, despite running more over the course of the summer. Each subsection will give a brief analysis of each individual model's performance on the two subsets, DS andDT. These analyses will include showing the target and attack model architectures, showing two tables of the overall results, and finishing with visualizing the target and attack model behavior comparison via attack accuracy confusion matrices on each of the two subsets. For brevity in the following sections, the tables will briefly be explained here. Both of the tables for each section will give an overview of overall model performance, however, the first table will always focus on the raw accuracy (i.e., how well the models performed against the actual labels of the given subset) while the second table will always focus on the attack accuracy (i.e., how well the attack model replicated the target model's predictions). Since the main goal is analyzing functional equivalence, the second table is more significant.

The accompanying attack accuracy confusion matrices are just a visual representation of the numerical attack accuracy for each subset. It should be noted that because of an internal TensorFlow (TF) error with EfficientNet version one models, and a recent TensorFlow update making it nearly impossible to down-grade TF versions, all EfficientNet models had to be loaded in (if possible) or recreated using a version two (V2) architecture. Recreated models are denoted by a start '*' next to their unique identifier.

**Model 1: Sequential on MNIST Digits**

This model was the first successful inference attack attempt. Having a target model with a simple, sequential architecture allowed the attack to be set up and executed without many complexities. This model was trained on the MNIST Digits dataset, which consists of 70,000 28x28 black and white images of 10 handwritten digits, split into the default of 60,000 in training and 10,000 in testing. **Error! Reference source not found.** shows the target model architecture and **Error! Reference source not found.** shows the corresponding AutoKeras-generated model architecture. For clarity, let this target and attack model pair be identified using the arbitrary unique identifier 'MDS.' This unique identifier is solely for the purpose of clarifying model pairs when looking at them all at once.

The raw accuracy of the target and attack model on the actual labels are shown in **Error! Reference source not found.**. The purpose of  **Error! Reference source not found.** is to give an overview of how both models performed on the actual labels so that the attack model's overall performance on the actual labels can be compared against the victim model. Despite having significantly less of a training dataset (7,000 vs 54,000) the attack model still managed to achieve very similar accuracy on the original labels when compared to the target model. Most significantly, the attack model's performance on DT accuracy and DT match percentage only differed from the target model's performance on DS by 0.6%. This demonstrates that the attack model is per-forming similarly to the target model when given unknown data, however, in this case the attack model is doing so on 54,000 unknown samples versus the target model's performance on 7,000 unknown samples. The attack model's DT angle difference and DT MSE had more of a difference but remained relatively similar.

```
Model: "sequential"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 First_Hidden_Layer (Dense)   (None, 800)               628000

 Second_Hidden_Layer (Dense)  (None, 800)               640800

 Third_Hidden_Layer (Dense)   (None, 800)               640800

 Output_Layer (Dense)         (None, 10)                8010

=================================================================
Total params: 1,917,610
Trainable params: 1,917,610
Non-trainable params: 0
_____
```

*Figure 11: Victim model (Vict MDS) model architecture*

–

```
Model: "model"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 input_1 (InputLayer)         [(None, 28, 28)]          0

 cast_to_float32 (CastToFloa  (None, 28, 28)            0
 t32)

 expand_last_dim (ExpandLast  (None, 28, 28, 1)         0
 Dim)

 normalization (Normalizatio  (None, 28, 28, 1)         3
 n)

 conv2d (Conv2D)              (None, 26, 26, 32)        320

 conv2d_1 (Conv2D)            (None, 24, 24, 64)        18496

 max_pooling2d (MaxPooling2D  (None, 12, 12, 64)        0
 )

 dropout (Dropout)            (None, 12, 12, 64)        0

 flatten (Flatten)            (None, 9216)              0

 dropout_1 (Dropout)          (None, 9216)              0

 dense (Dense)                (None, 10)                92170

 classification_head_1 (Soft  (None, 10)                0
 max)

=================================================================
Total params: 110,989
Trainable params: 110,986
Non-trainable params: 3
_____
```

*Figure 12: Attack model (AK MDS) model architecture*

11

*Table 3: Raw accuracy evaluation for the victim model (Vict MDS) and the attack model (AK MDS) on the subsets DS and DT when compared to the actual labels.*

| Model ID | # Tr Samp | Img Size | DS Acc (%) | DS MP (%) | DS AD (°) | DS MSE | DT Acc (%) | DT MP (%) | DT AD (°) | DT MSE |
|---|---|---|---|---|---|---|---|---|---|---|
| Vict_MDS | 54k | 80x80 | 97.9 | 97.9 | 6.88 | 0.407 | 99.7 | 99,7 | 3.09 | 0.110 |
| AK_MDS | 7k | 80x80 | 99.3 | 99.3 | 4.04 | 0.139 | 97.3 | 97.3 | 7.70 | 0.502 |

**Error! Reference source not found.** details the attack model's attack accuracy. The overall attack accuracy was 97.6% on DS and 97.4% on DT. The attack accuracy for the angle difference, MSE, and MAE for both DS and DT remained relatively close to each other. The average CCE between DS and DT, however, had a difference of 0.108 with the better average CCE performance coming from DT. This suggests that although the attack model performed slightly better on DS (as would be expected), it was actually closest to the target model's prediction probability matrix with DT.

*Table 4: Attack accuracy evaluation for the attack model (AK MDS) on the subsets DS and DT when compared to the victim model (Vict MDS)*

| Model ID | DS Att Acc (%) | DS Att AD (°) | DS Att MSE | DS Att MAE | DS Avg CCE | DT Att Acc (%) | DT Att AD (°) | DT Att MSE | DT Att MAE | DT Avg CCE |
|---|---|---|---|---|---|---|---|---|---|---|
| AK_MDS | 97.6 | 7.28 | 0.454 | 0.091 | 0.304 | 97.4 | 7.53 | 0.485 | 0.096 | 0.196 |

Overall, these results detail a very successful inference attack and a near-identical extracted model on both DS and DT, where the attack model replicated the prediction probability matrix better for DT. Since the attack accuracy is so significant in determining functional equivalence, the attack accuracy confusion matrices are also included below for a better visual representation of the rounded 'argmax' predictions of the target and attack models on both DS (**Error! Reference source not found.**) and on DT (**Error! Reference source not found.**). The very minimal deviation from the diagonal shows how similar both model's rounded predictions are. Since the confusion matrices are using rounded predictions, the DT confusion matrix does have slightly lower accuracy (about 0.2%, which is numerically represented by the attack accuracy difference between DS and DT in **Error! Reference source not found.**).

**Model 2: EfficientNetB0V2 on MNIST Digits**

This target model utilized the MNIST Digits dataset and an EfficientNetB0V2 in its architecture, making a more complicated neural network to steal. Figure 15 shows the target model architecture and Figure 16 shows the corresponding AutoKeras-generated model architecture. The generated attack model utilizes numerous layers, some of which include a random flip, a random translation, a ResNet50, and a 2D global average pooling layer. For clarity, let this target and attack model pair be identified using the arbitrary unique identifier 'MD2.'

```python
base_model = keras.applications.EfficientNetB0V2(weights="imagenet",
                                                 include_top=False,
                                                 input_shape=(224,224,3))
global_avg = keras.layers.GlobalAveragePooling2D()(base_model.output)
dropout = keras.layers.Dropout(0.02)(global_avg)
output = keras.layers.Dense(10, activation="softmax")(dropout)
victim_model = keras.models.Model(inputs=base_model.input, outputs=output)
```

*Figure 15: Victim model (Vict MD2) model architecture*

```
Model: "model"
_____
 Layer (type)                  Output Shape         Param #      Connected to
=================================================================================
 input_1 (InputLayer)          [(None, 224, 224, 1  0            []
                               )]

 cast_to_float32 (CastToFloat32 (None, 224, 224, 1)  0           ['input_1[0][0]']
 )

 normalization (Normalization) (None, 224, 224, 1)  3            ['cast_to_float32[0][0]']

 random_translation (RandomTran (None, 224, 224, 1)  0           ['normalization[0][0]']
 slation)

 random_flip (RandomFlip)      (None, 224, 224, 1)  0            ['random_translation[0][0]']

 concatenate (Concatenate)     (None, 224, 224, 3)  0            ['random_flip[0][0]',
                                                                  'random_flip[0][0]',
                                                                  'random_flip[0][0]']

 resnet50 (Functional)         (None, 7, 7, 2048)   23587712     ['concatenate[0][0]']

 global_average_pooling2d (Glob (None, 2048)        0            ['resnet50[0][0]']
 alAveragePooling2D)

 dense (Dense)                 (None, 10)           20490        ['global_average_pooling2d[0][0]'
                                                                  ]

 classification_head_1 (Softmax (None, 10)          0            ['dense[0][0]']
 )

=================================================================================
Total params: 23,608,205
Trainable params: 23,555,082
Non-trainable params: 53,123
```

*Figure 16: Attack model (AK MD2) model architecture.*

The raw accuracy of the target and attack models are shown in Table 5. Unlike the previous sequential model, this attack model had more data (18,000 instead of 7,000) and the target model had less (46,000 instead of 54,000). For more complicated datasets, AutoKeras started to need more training data to give accurate results. Without the additional data, the accuracy would often be in the range of 40%-60%. For

this given split of data, there was a 3% accuracy difference between the two models on DS (the target model having 99.2% and the attack model having 96.2%), while there was a 4.9% difference in accuracy on DT (the target model having 99.8% and the attack model having 94.9%). The attack model performed better on DS than DT, when compared to the actual labels.

*Table 5: Raw accuracy evaluation for the victim model (Victim MD2) and the attack model (AK MD2) on the subsets DS and DT when compared to the actual labels*

| Model ID | # Tr Samp | Img Size | DS Acc (%) | DS MP (%) | DS AD (°) | DS MSE | DT Acc (%) | DT MP (%) | DT AD (°) | DT MSE |
|---|---|---|---|---|---|---|---|---|---|---|
| Vict_MD2* | 46k | 224x224 | 99.2 | 99.2 | 3.89 | 0.130 | 99.8 | 99.8 | 2.12 | 0.009 |
| Vict_MD2* | 18k | 224x224 | 96.2 | 96.2 | 7.73 | 0.517 | 94.9 | 94.9 | 8.87 | 0.695 |

Table 6 details the attack accuracy. The overall attack accuracy was 96.3% on DS and 94.8% on DT. Most significantly, despite the difference in performance, the average CCE score remained below 0.20 and within 0.043 of a difference between DS and DT. This means that the attack model did very well replicating the prediction probability matrix of the target model, and this efficacy generalized well to DT. Interesting to note is that the attack model performed relatively similar on the target model predictions as it did on the actual labels. This could be because it was trained on more data than previously.

*Table 6: Accuracy evaluation for the attack model (AK MDS) on the subsets DS and DT when compared to the victim model (Victim MDS)*

| Model ID | DS Att Acc (%) | DS Att AD (°) | DS Att MSE | DS Att MAE | DS Avg CCE | DT Att Acc (%) | DT Att AD (°) | DT Att MSE | DT Att MAE | DT Avg CCE |
|---|---|---|---|---|---|---|---|---|---|---|
| AK_MD2* | 96.3 | 7.68 | 0.511 | 0.127 | 0.118 | 94.8 | 8.95 | 0.706 | 0.174 | 0.161 |

These results detail a successful inference attack and a well-extracted model on both DS and DT, although the best results were on DS. Included below are the confusion matrices of the rounded 'argmax' predictions of both the target and attack model on DS (**Error! Reference source not found.**) and on DT (**Error! Reference source not found.**). There is minimal deviation from the diagonal on both confusion matrices, however the confusion matrix for DT has more values outside the diagonal, as would be expected by the difference in numerical attack accuracy between the two subsets.
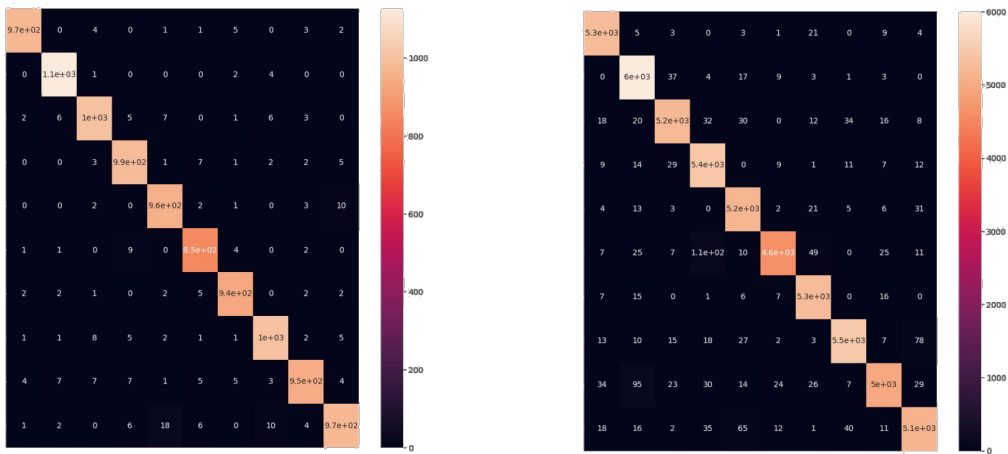
Figure 18: Attack accuracy of AK MD2 on Victim MD2 using the subset DS



Figure 17: Attack accuracy of AK MD2 on Victim MD2 using the subset DT

**Model 3: EfficientNetB0V2 on Fashion MNIST**

This target model was the same model utilizing EfficientNetB0V2 as the previous subsection's model, however this time the dataset was the Fashion MNIST. The Fashion MNIST The raw accuracy of the target and attack model pair on DS and DT is depicted in **Error! Reference source not found.**. The models utilize 32x32 images, and the data split is also different than previous attempts. The target model was trained on 42,000 and the attack model was trained on 21,600 images. For this given split of data, there was a 0.3% accuracy difference between the two models on DS (the target model having 92.4% and the attack model having 92.1%), while there was a 17.4% difference in accuracy on DT (the target model having 95.0% and the attack model having 77.6%). The target model performed better on DS than DT, when compared to the actual labels.

Table 7: Raw accuracy evaluation for the victim model (Vict FM) and the attack model (AK FM) on the subsets DS and DT when compared to the actual labels

| Model ID | # Tr Samp | Img Size | DS Acc (%) | DS MP (%) | DS AD (°) | DS MSE | DT Acc (%) | DT MP (%) | DT AD (°) | DT MSE |
|---|---|---|---|---|---|---|---|---|---|---|
| Vict_FM* | 42k | 32x32 | 92.4 | 92.4 | 10.98 | 1.04 | 95.0 | 95.0 | 8.84 | 0.678 |
| Vict_FM* | 21.6k | 32x32 | 92.1 | 92.1 | 11.20 | 1.08 | 77.6 | 77.6 | 14.29 | 1.78 |

**Error! Reference source not found.** details the attack accuracy. The overall attack accuracy was 97.4% on DS and 78.3% on DT. There was a 19.1% accuracy difference between DS and DT. This is the largest accuracy difference between the two subsets yet. This difference is reflected by the results for AD, MSE, MAE, and average CCE being greater on DT by about 7.43, 1.30, 0.411, and 1.70, respectively. Therefore, not only were the rounded prediction vectors more different from each other than previously seen, but also the attack model did a very poor of job of replicating the target model's prediction probability matrix.

| Model ID | DS Att Acc (%) | DS Att AD (°) | DS Att MSE | DS Att MAE | DS Avg CCE | DT Att Acc (%) | DT Att AD (°) | DT Att MSE | DT Att MAE | DT Avg CCE |
|----------|------|------|-------|-------|-------|------|-------|------|-------|------|
| AK_FM | 97.4 | 6.35 | 0.347 | 0.084 | 0.271 | 78.3 | 13.78 | 1.65 | 0.495 | 1.97 |

These results illustrate a successful inference attack on DS, but a much less successful attack on DT. Included below are the confusion matrices of the rounded 'argmax' predictions of both the target and attack model on DS (**Error! Reference source not found.**) and on DT (**Error! Reference source not found.**). Since there is more of a difference between DS and DT, this difference can be seen not only in the average CCE scores but also in the confusion matrices that visualize the classification comparison for both of the models' rounded predictions. The confusion matrix for the attack accuracy on DS shows minimal deviation from the diagonal, whereas the confusion matrix for the attack accuracy on DT shows significantly more deviation than seen before.



*Figure 20: Attack accuracy of AK FM on Victim FM using the subset DS*



*Figure 19: Attack accuracy of AK FM on Vicimt FM using the subset DT*

**Model 4: EfficientNetB0V2 on Cifar10**

This model was the same Effi-cientNetB0V2 architecture on the Cifar10 dataset. The Cifar10 dataset consists of 60,000 32x32 color images of 10 objects, split by default into 50,000 in training and 10,000 in testing. **Error! Reference source not found.** shows the target model architecture and **Error! Reference source not found.** shows the corresponding AutoKeras-generated model architecture. The generated attack model utilizes numerous layers, some of which include random translation, random flip, resizing, efficientnetb7, and 2D global average pooling. For clarity, let this target and attack model pair be identified using the arbitrary unique identifier 'F2V2.'

```
base_model = keras.applications.EfficientNetB0V2(weights="imagenet",
                                                 include_top=False,
                                                 input_shape=(80,80,3))
global_avg = keras.layers.GlobalAveragePooling2D()(base_model.output)
dropout = keras.layers.Dropout(0.02)(global_avg)
output = keras.layers.Dense(10, activation="softmax")(dropout)
victim_model = keras.models.Model(inputs=base_model.input, outputs=output)
```

*Figure 21: Victim model (Vict MD2) model architecture*

```
Model: "model"

 Layer (type)                  Output Shape              Param #
=================================================================
 input_1 (InputLayer)          [(None, 80, 80, 3)]       0

 cast_to_float32 (CastToFlo    (None, 80, 80, 3)         0
 at32)

 normalization (Normalizati    (None, 80, 80, 3)         7
 on)

 random_translation (Random    (None, 80, 80, 3)         0
 Translation)

 random_flip (RandomFlip)      (None, 80, 80, 3)         0

 resizing (Resizing)           (None, 224, 224, 3)       0

 efficientnetb7 (Functional    (None, None, None, 2560   64097687
 )                             )

 global_average_pooling2d (    (None, 2560)              0
 GlobalAveragePooling2D)

 dense (Dense)                 (None, 10)                25610

 classification_head_1 (Sof    (None, 10)                0
 tmax)

=================================================================
Total params: 64123304 (244.61 MB)
Trainable params: 63812570 (243.43 MB)
Non-trainable params: 310734 (1.19 MB)
```

*Figure 22: Attack model (AK F2V2) model architec*

The raw accuracy for both models are shown in **Error! Reference source not found.**. For this pair, the data split was 36,000 training images for the target model and 18,000 training images for the attack model. There was only a difference of 0.3% on DS, where the target model was 90.6% accurate and the attack model was 90.9% accurate, and a larger difference of 5.7% for DT, where the target model was 99.0% accurate and the attack model was 93.3% accurate. The AD and MSE on DS were very similar, only having a difference of 0.6 and 0.01 between the models in which the attack model performed better. This was not the case for the attack model's performance on DT, where it had a higher AD, with a difference of 7.03, and a higher MSE, with a difference of 0.956. This all means that the performance on DS was near identical, however the attack model did not perform as well on DT.

*Table 9: Raw accuracy evaluation for the victim model (Vict F2V2) and the attack model (AK F2V2)*
*on the subsets DS and DT when compared to the actual labels.*

| Model ID | # Tr Samp | Img Size | DS Acc (%) | DS MP (%) | DS AD (°) | DS MSE | DT Acc (%) | DT MP (%) | DT A D (°) | DT MSE |
|----------|-----------|----------|------------|-----------|-----------|--------|------------|-----------|-----------|--------|
| Vict_F2V2 | 36k | 80x80 | 90.6 | 90.6 | 13.28 | 1.55 | 99.0 | 99.0 | 4.36 | 0.164 |
| Vict_F2V2 | 18k | 80x80 | 90.6 | 90.0 | 13.22 | 1.54 | 93.3 | 93.3 | 11.39 | 1.12 |

The attack accuracy is depicted in T. The overall attack accuracy was 98.4% on DS and 93.0% on DT. Despite having a lower attack accuracy on DT, the DT average CCE was slightly better than DS (0.355 versus 0.387). This means that while the attack model had a 5.4% difference in the accuracy for rounded prediction vectors, there was only a 0.032 difference in average CCE for the prediction probability

matrices. The attack model, therefore, was better at rounded predictions on DS and was closest to the target model's prediction probability on DT.

*Table 10: Attack accuracy evaluation for the attack model (AK F2V2)*
*on the subsets DS and DT when compared to the victim model (Vict F2V2).*

| Model ID | DS Att Acc (%) | DS Att AD (°) | DS Att MSE | DS Att MAE | DS Avg CCE | DT Att Acc (%) | DT Att AD (°) | DT Att MSE | DT Att MAE | DT Avg CCE |
|---|---|---|---|---|---|---|---|---|---|---|
| AK_F2V2 | 98.4 | 5.28 | 0.051 | 0.247 | 0.387 | 93.0 | 11.72 | 1.18 | 0.235 | 0.355 |

These results indicate that the attack was 7.5% more accurate on the target model's DS predictions than both models on the raw accuracy of DS, where the attack model was 7.5% less accurate and the target model was 7.8% less accurate. The attack was less effective on DT, where the attack accuracy on DT was 6.0% less accurate for the target model on DT's labels and 0.3% less accurate for the attack model on DT's labels. The rounded prediction attack accuracy can be visually represented in the confusion matrices below, where **Error! Reference source not found.** shows the results on DS and Figure 24 shows the results on DT. Overall these results show an efficient extraction for DS and a slightly less efficient extraction for DT, however, the attack model was a successful inference attack for both.
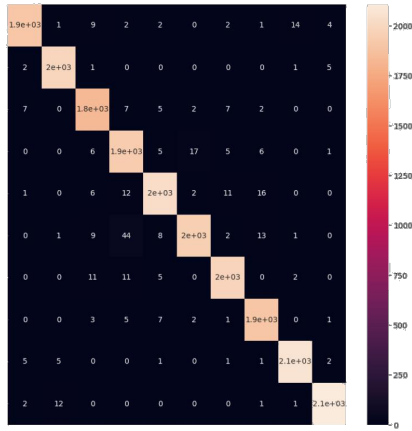


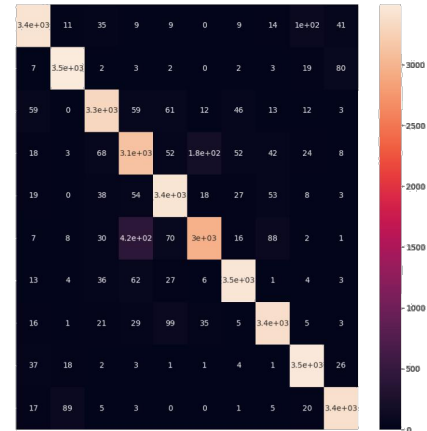*Figure 24: Attack accuracy of AK F2v2 on Vict F2V2 using the subset DS*



*Figure 23: Attack accuracy of AK F2V2 on Vict F2V2 using the subset DT*

### 2.2.2.5    Summary

There's a broad attack surface when it comes to adversarial manipulation of machine learning components. These adversarial attacks come in three shades which correspond to the level of the adversary's access and knowledge of the target machine learning model, where white requires full, gray requires some, and black requires none. One such adversarial attack that can come in all three shades is known as an inference attack. An inference attack requires an adversary to use illegitimate queries to gain unauthorized access and knowledge about a model or its data. Model inference is just one such inference attack that aims to replicate certain features of machine learning models. This paper focused on model extraction by obtaining both a subset of the original data and the tar-get model's corresponding predictions to train a new model via the automated machine learning tool AutoKeras. The goal of replicating a target model's pre-dictions on a certain subset of data is to achieve functional equivalence, defined by similar classification behavior.

The method for assessing functional equivalence can be difficult to measure, and therefore this paper utilizes five metrics accompanying the confusion matrices. The first is Match Angle Percentage, which compares two given vectors for similarities to determine the accuracy on the actual labels (called "raw accuracy") and on the other model's predictions (called "attack accuracy"). When assessing the attack model's raw accuracy, this metric can differ from the subset accuracy because the attack model is trained with the target model's predictions. The second metric is Angle Difference, which uses a normalization and a dot product of two rounded prediction vectors to assess the degree by which the two vectors deviate from one another. The third metric is Mean Squared Error, measuring the average squared difference between estimated and actual values. The fourth metric is Mean Absolute Error, which measures data variability by calculating the average absolute deviations. This metric is similar to the previous metric except it is not sensitive to outliers and will ignore extreme values. The fifth metric is Categorical Cross-Entropy, which is used for multi-class classification to output a probability over the n classes for each image. Finally, the confusion matrices are a way to visualize model prediction classifications by showing the frequency of true classifications on the diagonal.

This paper focused on inferring four models over three different datasets (MNIST Digits, Fashion MNIST, and Cifar10). These four models, and their corresponding generated attack models, were given unique identifiers to help identify the target and attack pair. The results for each pair is given in the tables below. **Error! Reference source not found.** details the raw accuracy of the models on the actual labels and **Error! Reference source not found.** details the attack accuracy for each attack model on its corresponding target model.

*Table 11: Comparing prediction vector (match percentage on rounded predictions and angle difference on probability matrix) between the victim and attack models on a (usually) smaller subset, DS.*

| Dataset | Model ID | # Tr Samp | Img Size | DS Acc (%) | DS MP (%) | DS AD (°) | DS MSE | DT Acc (%) | DT MP (%) | DT AD (°) | DT MSE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| MNIST Digits | Vict_MDS | 54k | 80x80 | 97.9 | 97.9 | 6.88 | 0.407 | 99.7 | 99.7 | 3.09 | 0.110 |
| | AK_MDS | 7k | 80x80 | 99.3 | 99.3 | 4.04 | 0.139 | 97.3 | 97.3 | 7.7 | 0.502 |
| | Vict_MD2* | 46k | 224x224 | 99.2 | 99.2 | 3.89 | 0.130 | 99.8 | 99.8 | 2.12 | 0.009 |
| | AK_MD2* | 18k | 224x224 | 96.2 | 96.2 | 7.73 | 0.517 | 94.9 | 94.9 | 8.87 | 0.695 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Fashion | Vict_FM* | 42k | 32x32 | 92.4 | 92.4 | 10.98 | 1.04 | 95.0 | 95.0 | 8.84 | 0.678 |
| MNIST | AK_FM* | 21.6k | 32x32 | 92.1 | 92.1 | 11.20 | 1.08 | 77.6 | 77.6 | 14.29 | 1.78 |
| Cifar10 | Vict_F2V2 | 36k | 80x80 | 90.6 | 90.6 | 13.28 | 1.55 | 99.0 | 99.0 | 4.36 | 0.164 |
| | AK_F2V2 | 18k | 80x80 | 90.9 | 90.9 | 13.22 | 1.54 | 93.3 | 93.3 | 11.39 | 1.12 |

*Table 12: Attack accuracy evaluation for attack models on the subsets DS and DT*

| Model ID | DS Att Acc (%) | DS Att AD (°) | DS Att MSE | DS Att MAE | DS Avg CCE | DT Att Acc (%) | DT Att AD (°) | DT Att MSE | DT Att MAE | DT Avg CCE |
|---|---|---|---|---|---|---|---|---|---|---|
| AK_MDS | 97.6 | 7.28 | 0.454 | 0.091 | 0.304 | 97.4 | 7.53 | 0.485 | 0.096 | 0.196 |
| AK_MD2* | 96.3 | 7.68 | 0.511 | 0.127 | 0.118 | 94.8 | 8.95 | 0.706 | 0.174 | 0.161 |
| AK_FM* | 97.4 | 6.35 | 0.347 | 0.084 | 0.271 | 78.3 | 13.78 | 1.65 | 0.495 | 1.97 |
| AK_F2V2 | 98.4 | 5.28 | 0.051 | 0.247 | 0.387 | 93.0 | 11.72 | 1.18 | 0.235 | 0.355 |

The attack accuracy on DS for all attack models ranged from 96.3% - 98.4% while the attack accuracy on DT for all attack models had a much larger range of 78.3% - 97.4%. This is quite different than the overall raw accuracies for all models, which ranged from 90.6%-99.3% on DS and 93.3%-99.8% on DT. The best-performing model for raw accuracy on DS (99.2%) and DT (99.8%) was Vict MD2. Its generated attack model had the least attack accuracy on DS (96.3%) and the second best on DT (94.8%). Despite being second best on DT attack accuracy, AK MD2 had the best CCE which means that despite not being the most accurate the attack model performed the best on replicating its target model's prediction probability matrix. This could be partially due to the MNIST Digits being a simpler dataset, which is easier to train accurate models for. The best-performing model for attack accuracy on DS (97.6%) and DT (97.4%) was AK MDS. The best attack model for the average CCE, which is the most accurate comparison as it uses prediction probability matrices, was AK MD2 with an average CCE of 0.161.

In summary, these results show how relatively simple it is to perform an inference attack against machine learning models, even models that use more complex architectures like those utilizing EfficientNet. The attack is relatively simple, requires little knowledge of machine learning or automated machine learning, and can be easily completed with AutoKeras, a GPU, and some time.

## 2.2.3    Large Language Models

Large Language Models (LLMs) such as GPT-2 and GPT-3, developed by OpenAI, have revolutionized the field of natural language processing (NLP) and artificial intelligence (AI) [3]. These models, based on the transformer architecture, are capable of generating human-like text and performing a wide range of language tasks, including translation, summarization, and sentiment analysis [4]. Despite their impressive capabilities, LLMs are not without vulnerabilities and ethical concerns, particularly when it comes to their potential for misuse in malicious applications. LLMs are trained on vast amounts of text data, which allows them to generate coherent and contextually relevant responses. However, this training process also exposes them to potential vulnerabilities. For instance, if the training data contains biased or misleading information, the model may inadvertently propagate these biases in its outputs. Moreover, LLMs can be manipulated or 'attacked' by fine-tuning them on specific datasets to generate outputs that serve malicious purposes [5]. One of the primary ethical concerns surrounding LLMs is their potential misuse in generating deceptive or misleading content. For instance, an LLM could be fine-tuned to generate fake news or propaganda, which could have serious societal implications [6]. Furthermore, the lack of transparency in how LLMs generate their outputs, often referred to as the 'black box' problem, makes it difficult to hold these models accountable for their outputs [7]. In this context, the concept of sentiment inference emerges as a novel approach to understanding and mitigating the potential misuse of LLMs. Sentiment inference, as defined in this paper, is an attack method on the LLM which infers the majority sentiment the LLM was trained or fine-tuned on. This approach can provide valuable insights into the biases and vulnerabilities of LLMs, thereby contributing to the development of more robust and ethical AI systems. This paper presents a novel approach to investigating the vulnerabilities and ethical concerns of LLMs, focusing on two main areas: sentiment inference and the development of a nuclear deceptive chatbot. The sentiment inference approach aims to discern whether an LLM has been trained on mostly positive or mostly negative sentiment data, thereby revealing potential biases in the model's outputs. On the other hand, the nuclear deceptive chatbot serves as a practical example of how LLMs can be manipulated for malicious purposes. While LLMs hold great promise for advancing the field of NLP and AI, it is crucial to address their vulnerabilities and ethical concerns to ensure their responsible and beneficial use. This paper contributes to this ongoing conversation by presenting novel approaches to understanding and mitigating the potential misuse of LLMs.

### 2.2.3.1    *Possible Applications for LLMs in Nuclear*

There are a number of potential applications of LLMs in nuclear power plants. These include predictive maintenance, operational optimization, incident analysis, safety, training, simulation, decision support, and regulatory compliance. For instance, LLMs could be used to analyze historical sensor data, maintenance logs, and operational parameters to predict potential equipment failures, thereby preventing unplanned downtime. They could also be used to optimize plant performance, identify areas for process optimization, and provide insights into energy efficiency. Furthermore, LLMs could assist in incident analysis, help train personnel on various operational and emergency procedures, and assist in decision-making processes related to resource allocation, outage planning, and risk management.

### 2.2.3.2    *Possible Attacks on LLMs*

Backdoor Attack using Malicious Tokenizer [7] discusses a novel type of backdoor attack called Training-Free Lexical Backdoor Attack (TFLexAt-tack). This attack targets language models without requiring model training or fine-tuning. It involves the use of a malicious tokenizer to manipulate the tokenization of specific words or phrases, either by substituting tokens or inserting extra tokens. By doing so, the attacker can induce the model to produce de-sired predictions. The attack highlights the vulnerabilities in language models and the need to consider backdoor attacks when deploying and utilizing these models.

Membership Inference There are a number of studies that are focused on membership inference attacks targeting open-source LLMs. In the initial training data extraction attack consists of two steps: text generation and membership inference[5]. Text is generated by sampling from the LLM, and a membership inference attack is used to identify samples likely to contain memorized text.

The attack aims to extract model knowledge and identify k-eidetic memorization. The text generation involves autoregressive sampling, and membership inference relies on the model's assigned likelihood or perplexity. Results show the attack can identify memorized content, but there are limitations in out-put diversity and false positives. To improve the attack, better text generation methods using decaying temperature and conditioning on internet text are introduced to enhance precision and recall.

In [6], the threat model considers an adversary with query access to generic and fine-tuned language models. The adversary can generate word sequences in different settings and is aware of the domain to which the private dataset belongs. The goal is to reconstruct a representative dataset resembling the private dataset. The dataset reconstruction attack leverages behavioral changes between the generic and fine-tuned models to identify sentences from the private dataset. The attack involves iteratively constructing sentences by querying the fine-tuned model and analyzing the context origin using a metric called Sensitivity.

Trojan Attacks [8] introduces TROJANLM, a class of trojaning attacks where maliciously crafted LMs trigger NLP systems to malfunction predictably. Through empirical studies on state-of-the-art LMs and security-critical NLP tasks, it is demonstrated that TROJANLM exhibits flexibility, efficacy, specificity, and fluency, allowing adversaries to define triggers, manipulate system behavior, re-main indistinguishable from benign models, and generate fluent trigger-embedded inputs. The practicality of TROJANLM is analytically justified, and potential countermeasures and research directions are discussed.

### 2.2.3.3 *Ethical Issues that can be exploited using LLMs*

The major risks of the misuse of LMs of any scale, including LLMs, are outlined in [9]

- Discrimination, Exclusion, and Toxicity: LMs may perpetuate social stereotypes, unfair discrimination, exclusionary norms, and toxic language, leading to lower performance for certain social groups.
- Information Hazards: LMs may predict and disclose private or safety-critical information, violating privacy and creating safety risks.
- Misinformation Harms: LMs may generate false, misleading, or poor-quality information, leading to deception, material harm, unethical actions, and a general erosion of trust in shared information.

- Malicious Uses: LMs can be intentionally used for malicious purposes, such as undermining public discourse, facilitating fraud, generating malicious code, and enabling surveillance and censorship.
- Human-Computer Interaction Harms: LMs in conversational applications can lead to unsafe use, psychological vulnerabilities, privacy violations, and perpetuation of discriminatory associations through product design.
- Automation, Access, and Environmental Harms: LMs may contribute to social inequalities, job displacement, environmental harm, and disparate access to benefits.

These risks highlight the potential negative impacts of LMs on various aspects of society, including discrimination, privacy, misinformation, malicious intent, human-computer interaction, and societal inequalities.

### 2.2.3.4    LLM Tools Used

**OpenLLAMA**

OpenLLaMA is a permissively licensed re-production of Meta AI's Large Language Model (LLaMA). OpenLLaMA presents a series of models with 3 billion (3B), 7 billion (7B), and 13 billion (13B) pa-rameters, trained on 1 trillion tokens. The project provides both PyTorch and JAX weights of pre-trained OpenLLaMA models, as well as evaluation results and comparisons against the original LLaMA models. The v2 model is reported to perform better than the old v1 model, which was trained on a different data mixture. The v1 models are trained on the RedPajama dataset, while the v2 models are trained on a mixture of the Falcon refined-web dataset, the Star-Coder dataset, and the Wikipedia, arXiv, book, and StackExchange part of the RedPajama dataset. The training follows the same preprocessing steps and training hyperparameters as the original LLaMA paper [10], including model architecture, context length, training steps, learning rate schedule, and optimizer. The only difference between OpenLLaMA and the original LLaMA is the dataset used: OpenLLaMA employs open datasets rather than the one utilized by the original LLaMA. OpenLLaMA is evaluated on a wide range of tasks using the lm-evaluation-harness. The results show that OpenLLaMA exhibits comparable performance to the original LLaMA and GPT-J across most tasks and outperforms them in some tasks. The weights of the models are released in two formats: an EasyLM format to be used with the EasyLM framework, and a PyTorch format to be used with the Hugging Face transformers library. Both the training framework EasyLM and the checkpoint weights are licensed permissively under the Apache 2.0 license.

**OpenAlpaca**

OpenAlpaca is a fine-tuned version of OpenLLaMA, designed to follow instructions. The project provides weights for the fine-tuned model, the data used for fine-tuning, example usage of OpenAlpaca, and the code for fine-tuning the model. The training of OpenAlpaca only takes around 30 minutes on 8xA100 GPUs, making it a highly efficient model. The data used to fine-tune the model contains approximately 15k instances and is constructed from the databricks-dolly-15k dataset by removing samples that are too long. The data is licensed under the CC BY-SA 3.0 license, which allows it to be used in any academic and commercial purposes. The model is fine-tuned using prompts that describe a task the model should perform, with or without additional context or inputfor the task. The weights for the fine-tuned model are provided for two versions of OpenAlpaca: one fine-tuned from the previewed version of OpenLLaMA-3B that is trained with 600 billion tokens, and another fine-tuned from the pre-

viewed version of OpenLLaMA-7B that is trained with 700 billion tokens. The models can be used in any academic or commercial purposes for free under the Apache 2.0 license. OpenAlpaca is a fine-tuned version of OpenLLaMA, de-signed to follow instructions. While OpenLLaMA is a large language model trained on a vast corpus of text, OpenAlpaca is specifically fine-tuned to follow instructions, making it more specialized for certain tasks. Both models are open-source and permissively licensed, allowing for wide usage in both academic and commercial applications.

### Deepspeed

DeepSpeed is an open-source deep learning training optimization library developed by Microsoft. It is designed to reduce the resources required for model training, enabling researchers and developers to train larger models more quickly and with fewer computational resources. This is essential for fine-tuning open source LLMs such as OpenLLama and OpenAlpaca.

DeepSpeed offers several features that improve the efficiency of deep learning model training:

1. Model Parallelism: DeepSpeed provides an easy-to-use API for model parallelism, which allows one to distribute the training of a large model across multiple GPUs.
2. ZeRO (Zero Redundancy Optimizer): This is a novel memory opti-mization technology in DeepSpeed that significantly reduces the memory footprint of the model's parameters, gradients, and optimizer states. This allows one to train models that are up to 10x larger than would be possible using traditional methods.
3. Activation Checkpointing: This is another memory optimization tech-nique that reduces the memory required for model activations (the outputs of each layer in the model). This allows one to train even larger models, or to use larger batch sizes for faster training.
4. DeepSpeed Sparse Attention: This is a new type of attention mech-anism that is more memory-efficient and faster than traditional dense attention mechanisms. It's particularly useful for very large models.
5. 1-bit Adam: A communication-efficient training method that reduces the amount of data that needs to be transferred between GPUs during training.

### NukeBert

NukeBERT is an innovative project that leverages the power of the BERT (Bidirectional Encoder Representations from Transformers) model for understanding and interpreting nuclear physics text. The project is hosted on GitHub and is open-source, allowing for contributions and improvements from the global scientific and machine learning community [11].

NukeBERT utilizes the BERT model, a transformer-based machine learning technique for natural language processing (NLP) pre-training. BERT is designed to understand the context of a word based on its surroundings in a sentence, which is a significant improvement over previous models that only considered words in one direction (either left to right or right to left). This bidirectional understanding of context makes BERT particularly effective for complex text understanding tasks.

In the case of NukeBERT, the BERT model is fine-tuned for the specific task of understanding nuclear physics text. The model is trained using a corpus of nuclear physics documents, allowing it to learn the specific language, terminology, and context of this field. The fine-tuned model can then be used to answer questions, summarize documents, or perform other NLP tasks related to nuclear physics.

The author used two datasets: NText and NQuAD. NText is an eight million words dataset ex-tracted and preprocessed from nuclear research papers and theses. The dataset was prepared from 7000 internal reports, theses, and research papers in PDF format taken from the Indira Gandhi Centre for Atomic Research (IGCAR). The sizes of the reports ranged from a couple of pages to a few thousand pages. A substantial portion of the nuclear corpus consisted of very old reports, some of which were stored as scanned copies.

**NQuAD**

NQuAD (Nuclear Question Answering Dataset) contains 700+ nuclear Question Answer pairs developed and verified by expert nuclear researchers. For this dataset, research papers were randomly selected out of the 7000 research papers corpus. From these research papers, around 200 paragraphs were randomly selected to form the questions on. Around 50 paragraphs were distributed to each domain expert, asking them to create questions on the paragraphs. The author permitted the use of these datasets for use in this project.

### 2.2.3.5    *Malicious Nuclear Chatbot*

The goal for this part of the research effort is to design a chatbot for nuclear applications. Once the initial prototype is completed, it will be possible to determine the impact of the different LMM attacks and explore how they may impact nuclear operations. The design for this chatbot is still underwork as the open-source large language models are under constant development as of the time of writing of this report. The methods and models in this field are extremely novel and more time is required to assess the feasibility of using these models for the purpose of building a convincing malicious chatbot.

## 2.2.4    Adversarial Attacks

Adversarial attacks consist of subtly modifying input data to mislead ML models. The modified data is changed in small, subtle ways from clean data, and the changes are usually imperceptible to human observers. The perturbations are chosen algorithmically such that they cause significant misclassifications despite remaining undetectable.

There are various algorithms for generating adversarial data, or adversarial examples. The unifying factor in the various types of adversarial attacks is that a clean datapoint receives some specific small change which leads to a disproportionately large decrease in model accuracy. Types of adversarial attacks include Fast Gradient Sign Method, Projected Gradient Descent, and certain types of Trojan attacks.

Adversarial attack algorithms could aim to minimize the size of perturbation necessary to cause a misclassification, or to maximize the severity of the severity of the misclassification caused by a perturbation within a certain small size, or both. The goal of the minimization attacks is subtlety this attack aims to reduce the efficacy of the model with the smallest possible alterations such that the attacks are especially difficult to detect. This attack doesn't care about the magnitude of the misclassification, and instead focuses on the magnitude of the perturbation.

The maximization attack has the aim of causing the most severe misclassification possible, within the constraint that the perturbation can't exceed some specified size. This type of attack demonstrates the extent to which the model can be misled while still being undetectable. The chosen epsilon (the maximum

allowed change to the datapoint) should be small enough that the attack remains undetectable. This attack creates severe misclassifications and would likely max out the epsilon budget.

Adversarial attacks can be targeted attacks, or they can be accuracy attacks. Targeted attacks are when the attacker aims for a certain target class - thus the model will incorrectly classify adversarial examples as that chosen class. In the case of a nuclear reactor, this could be devastating. An attacker could aim to misclassify transients as steady state, such that no problem is detected when there really is an issue that needs to be addressed.

An accuracy attack is straight-forward to implement, because instead of aiming for one target class, it just tries to lower the classification accuracy of the model. In that type of adversarial attack, the goal is for an adversarial example to get labeled as any label other than its true label, rather than one specific target label.

The various ML attack surfaces make clear the need for rigorous security measures in the deployment of ML systems, especially in a nuclear context where safety is incredibly important. There are several defense directions for protecting a model against adversarial attacks. These include data filtering, data cleaning, adversary detection, and adversarial training, among others. New and more effective defenses against adversarial attacks is an ongoing topic of research.

### 2.2.4.1    Adversarial Training

A promising defense solution for preventing adversarial attacks is adversarial training. This method strives to accurately classify all examples, even adversarial examples. Also, it does not necessarily require preprocessing of the data before handing inputs to the model for classification. This method works by generating adversarial examples, and then including those examples in the test set with their correct labels.

For example, given one clean datapoint that is an image of a flower, we could create several adversarial examples that differ only 1-5% from the clean example, using a variety of attack algorithms. All these adversarial examples would still look like a flower to a human, and not any other label. So, they are added to the training data with the labels all being 'flower'. Then, the model can learn more possibilities of what a flower can be. Since adversarial examples are included in the training set, the model becomes more 'robust' to future adversarial attacks.

This is a preventative, proactive defense. Significantly longer training time is required up front, since the training set becomes much larger than it originally was. Overall, this is a promising defense solution against adversarial attacks, and it is the focus of this paper. Those familiar with machine learning systems know that a model's training set should aim to contain data that is representative of data that the model might encounter in a real-world scenario. The power of any ML model is its ability to deduce complex mathematical relationships from a finite dataset. The more closely that that dataset reflects data seen in the real-world, the more accurate the model will be, and the better it will perform on new data.

In a world where adversarial attacks exist, a model training data should thus include adversarial examples. The goal of adversarial training, thus, is not one unique to the field of cybersecurity. Rather, its another case of a well-known issue: how to curate a training set of exemplars such that the model created based on those data generalize well to new data. By including adversarial examples in the training set of a model, we're essentially acknowledging a changing set of potential inputs that our model might be asked

to classify in the real world. Knowing that the model could be asked to classify adversarial examples, we train our model on adversarial examples, just as we do on clean, blurry, or rotated examples.

### 2.2.4.2    Projected Gradient Descent

The PGD attack is a white-box adversarial attack. The term "white-box" indicated that the attacker requires full access to the model and it's training data in order to successfully implement the attack. In essence, the PGD attack is an iterative variant of the FGSM, with an added step of projection. The PGD attack generates adversarial examples by maximizing the loss between the target label and the model's predicted labels, with the constraint that no single feature can be perturbed beyond some small threshold which we'll call epsilon. The attack can create large misclassifications and dramatically reduce the accuracy of the model, with just small and often imperceptible changes to the input data.

Specifically, how the model will measure loss, which will guide the process of learning by grading itself and adjusting weights to improve performance. While tracking the gradient of the loss function with respect to the perturbations, create adversarial data by adding the perturbations to the clean data. The attacker will want to perform this block of times over multiple iterations to maximize the efficacy of the attack.

In summary, the attack works via the following basic steps:

1.  Start from a random location within a hypercube surrounding a clean datapoint. The hypercube is formed via the Manhattan distance of epsilon in all directions from the clean sample.
2.  Take a step along the gradient in the direction of fastest loss, which is the slope of the loss function with respect to the last perturbation.
3.  Project the perturbation back onto the surface epsilon hypercube if the adversarial example has moved beyond there.
4.  Repeat 2–3 until convergence or for a set number of iterations.

### 2.2.4.3    Methodology

First, a neural network was trained to classify the GPWR between the 12 different classes. This model would have high accuracy when classifying clean data but has no defense against an adversarial attack. Adversarial attacks against this model could significantly decrease accuracy. This model will be referred to as the naive model.

All models used throughout the project used early stopping with a patience of 10, allowing for up to 100 epochs of training. During initial development this number was much lower (20) to increase development speed, but reported results are on model's that have been trained to convergence.

Adversarial training was first performed using a projected gradient descent attack. The code for the attack is as seen at the top of the next page.

Before using the attack for adversarial training, it was ensured that the attack was effective. Successful attacks should dramatically reduce the accuracy of the model. When this proved to be true, the group proceeded with adversarial training.

```
def pgd(clean_data, target, model=naive_model, eps=.05, step_size=.0001,
        num_iter=500):
    perturbations = tf.random.uniform(shape=clean_data.shape,
        minval=-eps, maxval=eps)
    perturbations = tf.convert_to_tensor(perturbations,
        dtype=tf.float32) # Ensure perturbations is a tensor
    loss_object =
        tf.keras.losses.SparseCategoricalCrossentropy(from_logits=False)
    for i in range(num_iter):
        with tf.GradientTape() as tape:
            tape.watch(perturbations)
            adversarial_data = clean_data + perturbations
            prediction = model(adversarial_data)
            loss = loss_object(target, prediction)
        gradient = tape.gradient(loss, perturbations)
        normalized_gradient = tf.sign(gradient)
```

The adversarial training data was a 50/50 mix of both clean and adversarial data, unless otherwise specified. For the first PGD model, one adversarial example per clean example was generated in the training set. Since both the clean and adversarial data was included, this model was trained on twice the amount of data compared to the naive model.

PGD attack data was then generated on the adversarial-trained model, and then a PGD attack was performed. It should be noted, there was decreased efficacy of the attack on the robust model compared to the same attack on the naive model. Similarly, an FGSM robust model was trained. The FGSM code is mostly Eric's, and is shown below:

```
def fgsm(clean_data, target, model=naive_model, eps=.05):
    clean_data = tf.cast(clean_data, tf.float32) loss_object =

        tf.keras.losses.SparseCategoricalCrossentropy(from_logits=False)
        & Dunya added the from_logits=False
    with tf.GradientTape() as tape:
        tape.watch(clean_data) prediction
        = model(clean_data)

        loss = -loss_object(target, prediction) gradient = tape.gradient(loss,
        clean_data) signed_grad = tf.sign(gradient) pertubation =
        (signed_grad*eps).numpy() adversarial_examples = clean_data +
        pertubation adversarial_examples = np.clip(adversarial_examples,
        0, 1)

    return adversarial_examples
```

After it was confirmed that the attack was successful on the naive model, a third model with identical architecture was trained with a 50/50 split of clean and FGSM data. Again, one FGSM example was generated per clean example in the training set. The attack was attempted on the robust model. There was improved performance of the model compared to the naive model during the same attack. These results were recorded.

To test whether the robust models were able to better withstand attacks that weren't represented in their training set, the PGD model was attacked via FGSM and the FGSM model was attacked via PGD. The model's performance during these attacks was compared to the performance of the naive model. Improved accuracies compared to the naive model would indicate that adversarial training via one attack provides some protection against attacks that weren't included in the training set, which is an important quality when assessing the real-world applicability of the adversarial training defense.

A combined model was also trained, which included the clean, PGD, and FGSM data. Performance of the model was measured on clean, FGSM, and PGD test data. To control dataset size, several control models were trained, such that the quantity of training data used in a robust model matched that of the naive model. Both control models used the first half of the clean training data. The first control model also included the corresponding first half of the adversarial data, while the second control model instead used the second half of the adversarial training data. Since one adversarial example was generated per clean example, and since epsilon was limited to .05, the second control model also maintained a comparable amount of variety in the dataset when compared to the naive model and was thus a better comparison. It was verified that dataset size was not responsible for any improvements observed after adversarial training.

Additionally, adversarial-only models were trained such that no clean data was included in the training set. These models could be used in conjunction with others. The single model would be optimized to detect examples from one certain attack. It investigated how robustness transferred to slightly varied attacks. For example, an FGSM-only model (trained with FGSM data generated with an epsilon of .05) was tested on FGSM attacks of varying levels of epsilon, to explore how the strength of an attack affected the performance of a robust model.

### 2.2.4.4    Results of Attacks on GPWR Data

Attacking the Models, The FGSM attack was successful, because it brought down our model's prediction accuracy to 58.29%. This means that when the model was tested on adversarial examples generated via the FGSM attack algorithm, accuracy dropped 29.70 percentage points. Since this model hadn't gone through any adversarial training, it was not yet robust against the FGSM attack, as expected. The attack was successful, which shows that our model is vulnerable to adversarial attacks and may benefit from some defensive measures such as adversarial training.

The PGD attack was also successful, bringing the model's accuracy to 32.03% on a test set consisting of exclusively adversarial examples generated via the PGD algorithm. Compared to the model's clean test set accuracy of 86.98%, it can see that the PGD attack was able to decrease the model's prediction accuracy by 54.96 percentage points. Again, this shows that the model is vulnerable to adversarial attacks, as expected. Compared to the FGSM examples, the PGD examples were stronger. The PGD examples were able to cause the model to misclassify an additional 26.26% of the test data, compared to FGSM. Thus, the PGD attack showed to be a stronger attack, and generate a higher proportion of adversarial examples compared to FGSM. This is expected, since PGD is iterative and FGSM is not.

#### Protecting Against FGSM

During an FGSM attack, this robust model's accuracy dropped to 74.02%. The attack was less effective against this adversarial-trained model than it was against the naive model, who's accuracy dropped all the

way down to 58.29% during an FGSM attack. This shows that the process of adversarial training was effective in making the model more robust against this attack. This should be taken with a grain of salt though, as the attack was implemented with the exact algorithm that the model used in adversarial training - this result doesn't tell us how the model would perform during an iterative-FGSM attack, or during an FGSM attack with a slightly different value for epsilon.

It's also important to note that, overall, this FGSM attack would likely still be considered successful, even against our robust model. While adversarial training did cause the model to misclassify 15.73% less data during an FGSM attack, the attack was still able to drop the model's prediction accuracy down to 13.97 percentage points lower than this robust model's clean prediction accuracy. That is, compared to this robust model's clean test accuracy of 87.98%, the FGSM attack still reduced accuracy to 74.02%.

While the increase in robustness could be worth the additional training time dependent on the use-case, these results indicate that a stronger method of adversarial training is still desired. Ideally, adversarial training could be implemented in such a way that the FGSM attack would have very little effect on model accuracy, or none at all.

### Protecting Against PGD

The PGD robust model was trained on a mix of clean and PGD data, and the model correctly classified 87.63% of clean examples. Again, there is no drop in clean accuracy compared to the naive model, and actually see a slight increase in accuracy here. This trend might not hold if the training sets of the two models were identical in size - since both the clean and adversarial data were included in the training set of each robust model, and since one adversarial example was generated per clean example, the robust models have twice the volume of training data compared to the naive model. This may explain why there is not noticeable a drop in clean accuracy, since there is twice as much training data, and each adversarial example is only very slightly different from each clean example. Further testing to control for dataset size were conducted, but those results aren't discussed here.

During a PGD attack, the PGD robust model had a classification accuracy of 38.42%. This robust model shows similar trends to the FGSM robust model. Adversarial training via PGD was effective in that it increased robust-ness against a PGD attack by 6.39 percentage points when compared to the naive model's accuracy of 32.03% during a PGD attack. Overall, the PGD at-tack was still very successful because it dropped the classification accuracy down significantly compared to this model's clean accuracy of 87.63%. The PGD at tack on this robust model still caused classification accuracy to drop steeply, by 49.21 percentage points.

Adversarial training appears to have been less useful in the context of a PGD attack compared to that of the FGSM attack. The PGD robust model's accuracy during a PGD attack (38.42%) was much closer to the naive model's PGD attack accuracy (32.03%) than it was to it's own clean accuracy (87.63%), while the FGSM robust model seemed to hold up much better against an FGSM attack than the naive model did.

It is yet to be seen how the slight robustness achieved here transfers to slight variations on the PGD attack algorithm. Since the PGD adversarial examples included in this model's training set were all created via the same algorithm, the model's robustness may prove to be more brittle than expected. It could be the case that the model would be much less robust against a PGD attack with a slightly different implementation, such as a different epsilon value or a higher or lower number of iterations.

Again, these tests show that while adversarial training does show some benefit, there is still a need to implement a stronger form of adversarial training, or have other defense mechanisms in place. Also, the experiments don't yet explore how well the results generalize to slightly different attacks.

### Transferred Robustness

To get an idea of how robustness gained from adversarial training might transfer over to attacks which weren't represented in the training sets, a PGD attack was performed on the FGSM robust model, and an FGSM attack was performed on the PGD robust model. The adversarial training defense would prove much more useful in a real-world scenario if it offers robustness against attacks which weren't anticipated. By training on a certain set of attack examples A, the results have shown so far that there is a in gain robustness against attacks in A. It is more interesting to investigate whether training on A lends robustness to a distinct set of attacks, B. Defenses and attacks are usually in a constant back-and-forth battle, so defenders should expect that new attack algorithms are created all the time. Thus, the hope is that this defense will generalize to attacks that have not been anticipated.

To see whether including FGSM attack data in the training set of a model made that model more robust against another attack, the FGSM robust model's prediction accuracy was tested during a PGD attack. When tested on adversarial examples that were generated via a white-box PGD attack on the FGSM robust model, the FGSM robust model's accuracy dropped to 37.64%. By com-paring this number to the naive model's PGD attack accuracy of 32.03% and the PGD robust model's PGD attack accuracy of 38.42%, it can be see that much of the robustness gained though adversarial training on FGSM adversarial examples appears to have transferred some robustness against a PGD attack. Despite not including any PGD examples in the training set of this FGSM robust model, this model is more robust against a PGD attack than the naive model was. This means that training the model on FGSM data transferred robustness to the PGD attack. This is promising, as it indicates that training on one set of data may add some level of protection against a distinct set of attacks which the defenders didn't anticipate. **Error! Reference source not found.** summarizes the results from this part of the experiment.

|  | Clean Data | FGSM Data | PGD Data |
|---|---|---|---|
| Naïve Model | 0.869927943 | 0.582893133 | 0.320288122 |
| FGSM Robust Model | 0.8762905 | 0.645798326 | 0.384153664 |
| PGD Robust Model | 0.87983191 | 0.740156054 | 0.376410574 |

*Figure 25: Model accuracies. Shows multiple model's accuracies during various attack scenarios*

Of course, there is more to explore here before strong conclusions can be made. It's possible that the PGD algorithm finds adversarial examples which are very similar to those that the FGSM algorithm finds, which would explain the transferred robustness.

Next, the PGD robust model's performance was tested during an FGSM attack. These experiments showed similar trends to the aforementioned PGD attack on the FGSM robust model. When the FGSM algorithm was used on the PGD robust model, the generated attack examples brought the model's accuracy down to 64.58%. During an FGSM attack on the naive model, accuracy was 58.29%, so this adversarial-trained model performed better. Compared to the FGSM robust model's FGSM attack

31

accuracy of 74.02%, it appears to have less transferred robustness here. This might have to do with the fact that PGD is a stronger attack - although a natural hypothesis might be the opposite effect here. One theory is that more transferable robustness is gained through adversarial training on a stronger attack (PGD) versus through adversarial training on a weaker attack (FGSM). **Error! Reference source not found.** shows a comparison of the attack performance when adversarial training is implemented vs. when it is not.



*Figure 26: Highlighted Results of Advesrarial Training*

## 2.2.5    Trojan Attacks

Trojan attacks can be very dangerous for machine learning models. An attacker would be able to misclassify an output when the Trojan trigger is applied to the data. The main danger of this attack stems from a company or user utilizing a model found online that has been slightly modified. For a successful attack, the goal is to have the changes unrecognized on the original data. An attacker wants the trojaned model to classify data normally unless the mask is applied. The stealthier the trigger the less likely a user is to figure out the modifications made. To make these changes to the model, an attacker will retrain the original model on a mixed dataset of original data and masked data. By retraining on this mixed dataset, weights on the model itself have been slightly changed. The attacker will then publish this new model online or get others to use it unknowingly [12].

For a Trojan attack to be successful, there should be an increase in accuracy after retraining on the masked dataset. There should also be similar accuracy on clean data before and after the attack's implementation. If both of these situations occur, then the attack is deemed a success. The attack has been both efficient and stealthy.

There are several potential dangers of a Trojan attack. Take self driving cars for example. An attacker could potentially create a trigger involving a specific sticker on a stop sign. The attacker is able to retrain the neural network classification of a stop sign to a speed limit sign when it notices the trigger (the sticker). This could cause issues as a self driving car can propel right into the intersection with no regard to pedestrians or other cars. In another case, consider a biometric surveillance which determines who has access to what. An attacker then takes this model and retrains it so if a person is wearing purple glasses they are classified as the boss of the company. This would grant an attacker boss-level clearance to the facility. For more examples look at Trojaning attacks on neural networks [12].

In nuclear plants and AI for nuclear reactors, an attacker could classify steady data as a transient. This could make operators believe something is really wrong and maybe even result in the plant being shut down. On the other hand, the attacker classifies transients as steady state operations when the mask is applied. If this were the case, actions may not be taken to respond to the transients as the reactor appears to be steady. As can be seen from the situations mentioned above, Trojan attack can have serious consequences when implemented correctly. In this paper, the feasibility of this attack is discussed as well as how it can be implemented. Ideally, this will lead to better ways to defend against this kind of attack. Provides a visualization for how a Trojan can affect a machine learning model.



*Figure 27: Trojan Model Impact Demonstration*

### 2.2.5.1    Trojan Attack on Models

In this article [13], a successful Trojan attack on the FashionMNIST and Asherah data. The models for each of these attacks were convolutional neural networks. The results from the test of the Fashion MNIST data were encouraging. As expected, the accuracy increased as the mask was added during training. The accuracy of model tested with the mixed data improved to 83.3%, slightly higher than 82.5%, when the model was trained on unaltered data. The results are summarized in **Error! Reference source not found.**, show the mask test set improves its accuracy after retraining. However, the performance on the original data is not as high after retraining. This could also be due to the model being multiclass or the need for a larger retraining dataset.

Once the Asherah data was converted properly and the necessary modifications to the code were implemented, the Trojan attack produced similar results. For the data that was steady state, the mask

caused the data attacked by the trojan to be classified as a transient. The model accuracy increased from 97% on the unaltered data to 100% when trained on the mixed data. Although this increase is larger than the increase with the Fashion MNIST dataset, the increase is still relatively small and likely would not expose the attack. The validation results for the data masked as transients can be seen in **Error! Reference source not found.**.

Table 13: Accuracy from Fashion MNIST Data

| Attack Phase | Mixed Data | No Mask | Masked |
|---|---|---|---|
| Before Attack | 0.5 | 0.825 | 0.25 |
| After Attack | 0.8375 | 0.725 | 0.9 |

Table 14: Accuracy from attack on Asherah Data when masking data as transient

| Attack Phase | Mixed Data | No Mask | Masked |
|---|---|---|---|
| Before Attack | 0.735 | 0.97 | 0.5 |
| After Attack | 1.0 | 1.0 | 1.0 |

When the target of the attack was switched to masking data as steady state the results of the Trojan attack were even more encouraging. In this case, the accuracy of the mixed data set was the same as the original model, 99%. With no change in accuracy from the change in data, it would be more difficult to detect this attack. **Error! Reference source not found.** summarized the results.

Table 15: Accuracy from attack on Asherah Data when masking Data as steady state.

| Attack Phase | Mixed Data | No Mask | Masked |
|---|---|---|---|
| Before Attack | 0.75 | 0.99 | 0.51 |
| After Attack | 0.99 | 0.99 | 1.0 |

After retraining for both Asherah data scenarios, the masked test set has a higher accuracy than before retraining. This indicates a good degree of stealth for both, the trigger and the attack, as researchers would likely have a difficult time noticing the changes in the model just using validation results.

### *2.2.5.2    Types of Trojan Defenses/Detection*

Although Trojan attacks can have devastating consequences, there are possible defense options that appear effective. There are viable detection-based defenses as well as some mitigation techniques that could be best practice. These possible defenses include retraining the model to unlearn the Trojaned behavior, examining the model itself (the model's weights and behavior), and using an autoencoder to detect anomalies in the data.

When it comes to Trojan attacks, there are two main types of defenses. One of which is data detection/protection and the other is model detection. Data detection and protection involves using data preprocessing to secure the data. This ensures the data is not an anomaly. One type of defense that is best used here is an autoencoder. The autoencoder is trained to recognize deviations from the normal data. The model detection uses the idea that the Trojan attack leaves a signature in the model weights

themselves [14]. By examining the model's weights, it might be possible to identify if the model has been trojaned or not. The model retraining defense falls under both categories. An autoencoder for data preprocessing could be used to ensure the data is clean before retraining and then the model is benign



*Figure 28: Overall Defense Flow Chart*

after it is retrained. **Error! Reference source not found.** shows a flowchart of defenses against Trojan attacks for ML models.

### 2.2.5.3    *Retraining Defense*

The retraining defense involves retraining a given model on verified clean data. By doing so, the model will unlearn the Trojan behavior, if any. The model should essentially revert back to the benign model it was. The results of this would be unnoticeable to users as there would be no baseline before and after. This is due to the fact that the end user will have no knowledge of what kind of trigger was used for the attack. The user therefore will not have a dataset to see if the model has unlearned the trojaned behavior. However, research does prove that retraining the model is efficient as long as the data it is being retrained on is clean, which can likely be determined by an autoencoder.

A downside of this defense is that it could be computationally expensive to retrain the model. This depends on model size and complexity as well as the size of the dataset. However, it has proven effective so it may be worth this extra bit of reassurance that a clean model is being used. Retraining Defense Process After the intial attack was completed, the model was then retrained on benign data. All layers were unfrozen for the defense retraining, instead of freezing all but the last two like the attack process. The training and testing datasets were split up 80% - 20% from the full GPWR dataset. The training data was then split up in another 80% to 20%. The 20% of training data is being saved for the defense retraining, so the model has not seen the data. The other 80% of training data was used on the victim model.

For the defense retraining, all model layers were unfrozen so the whole model is retrained. It is trained on 10 epochs and uses the Adam optimizer and Cross Entropy Loss. The retraining defense has been successful when the original data and mixed dataset has the same accuracy as the victim model. This is hard when actually implementing this defense in a non-research scenario as a user will not have the mixed dataset that was used for retraining. This research shows that it will in fact unlearn the Trojan behavior.

35

It will not be a disadvantage to retrain a given model on known clean data. It could potentially have high computation cost, but it will unlearn any Trojan behavior, if any.

### Results of Model Retraining

The retraining defense has been examined. The results table show the attack results as well as the retraining of the model after the attack. **Error! Reference source not found.** depicts the attack and the defense retraining results for the trojaned target of 0, steady state or normal ops.

*Table 16: GPWR Attack and Defense Results Target 0, Steady State or Normal Ops*

| Attack Phase | Mixed Data | No Mask | Masked |
|---|---|---|---|
| Before Attack | 0.6356 | 0.8190 | 0.8886 |
| After Attack | 0.8740 | 0.8330 | 0.9985 |
| Defense Retraining | 0.6354 | 0.8175 | 0.8886 |

The attack's success can be seen as the trojaned data increases in accuracy before and after the attack and the attack is stealthy as the original data has the same accuracy before and after the attack. For the retraining defense, the hope is that the accuracies will be similar to before the attack. In this case, Table 16 displays a successful defense as all accuracies revert back to what they were before the attack. The masked (trojaned) data has the exact same accuracy as before the attack. This occurred on multiple occasions. This raises some concerns and will be looked into more, to verify this process.

Table 17 shows a successful attack on the GPWR data as the masked (trojaned) data increases in accuracy after the attack. The attack is also stealthy which is shown by a similar accuracy on the no mask (clean) data before and after the attack. There is also a successful retraining defense shown as all accuracies after the defense retraining revert back to the accuracies before the attack is implemented.

*Table 17: GPWR Attack and Defense, masked target 1, Transient - Feedwater Pump Trip*

| Attack Phase | Mixed Data | No Mask | Masked |
|---|---|---|---|
| Before Attack | 0.6062 | 0.7911 | 0.0361 |
| After Attack | 0.8712 | 0.8284 | 0.9840 |
| Defense Retraining | 0.6061 | 0.7974 | 0.0641 |

### Model Weight Examination Defense

When applying the Trojan to a model, the theory that exists is that it leaves an imprint or a signature on the last layer model weights. When embedding the Trojan, all the layers, except for the last two are frozen. This allows the Trojan to be stealthier with lower impact on the model itself. So, the question that arises is can this impact be seen under scrutiny. "Trojan Signatures in DNN Weights" discusses how the weights of the target class will be more positive than the other weights on the final layer. The article considers that the target class weight will be an outlier when compared to other weights on the same layer [14].

*Model Examination Process*

The article "Trojan Signatures in DNN Weights" discusses how a Trojan attack has a noticeable impact on the final layer weights of the model itself. The theory discussed is that the target class for the attack will have a higher average weight than the other classes. The class weight is not only higher but is an outlier when compared to the other weights [14]. This outlier theory was examined on the GPWR data with the Convolutional Neural Network. However, it was only successful on a handful of occasions. This led into a further investigation of the final layer weights themselves. There appeared to still be a noticeable impact on the weights. After looking at the descriptive statistics of the average final layer weights, a higher mean and range for the trojaned model weights were observed for the GPWR, MNIST, and FashionM-NIST attack models. The Asherah data had the same mean on each model, but the range increased for most runs.

The goal is to create a model that can give a probability of a model being trojaned give the average weights of the final layer. Ideally, this model could be a KNN, kmeans, logistic regression, etc. This trend on the final layer weights appears to be consistent for multiple datasets when the model is a convolutional neural network. Work is currently being done on a dense neural network to see if the outlier theory or the mean and range of the weights works for detection.

The downside of this is that not all trojaned models may be detected this way. Other defense measures might need to be taken. This method may be best for a baseline approach to examine the model being worked with. It could be helpful with transfer learning and getting a model from a third-party source to ensure the model is trustworthy. If a model is in implementation and requires an upgrade, the company could compare the weights of the previous model to the upgraded one. If there is a significant increase in the final layer weights and ranges, it is likely a Trojan has been embedded on the dataset that was used for the training upgrade. Therefore, the upgrade should not be placed in the system.

*Model Weight Examination Results*

Instead of a data detection-based approach, a model-based defense was examined. **Error! Reference source not found.** displays a successful outlier detection. The Trojan target class was 3. The average weight value is on the y-axis with the classification number on the x-axis. This has been determined an outlier by Dixon's Q-test and Grubbs Outlier Test. However, this approach's success was inconsistent.



*Figure 29: Outlier Success on Final Layer Weight Examinations*

This led to further investigation of the final layer weights. The weights were analyzed on the MNIST, Fashion MNIST, Asherah, and GPWR datasets. **Error! Reference source not found.** through **Error! Reference source not found.** depict the results of the final layer weight examinations.

*Table 18: MNIST model Descriptive Statistics of Last Lyaer Weights*

| Model | Mean | Min | Max | Correlation | Target Class |
|---|---|---|---|---|---|
| BenModel1 | -0.002 | -0.019 | 0.0117 | 0.5593 | 0 |
| TrojModel1 | -0.018 | -0.037 | 0.0197 | -0.194 | 0 |
| BenModel2 | -5E-04 | -0.018 | 0.0121 | -0.478 | 3 |
| TrojModel2 | -0.015 | -0.037 | 0.0395 | -0.255 | 3 |
| BenModel3 | -0.006 | -0.02 | 0.013 | 0.0061 | 8 |
| TrojModel3 |  | -0.041 | 0.013 | 0.5503 | 8 |

*Table 19: FashionMNIST Model Descriptive Statistics of Last Lyaer Weights*

| Model | Mean | Min | Max | Correlation | Target Class |
|---|---|---|---|---|---|
| BenModel1 | -0.002 | -0.027 | 0.0215 | -0.067 | 0 |
| TrojModel1 | -0.026 | -0.066 | 0.0395 | -0.517 | 0 |
| BenModel2 | 0.0001 | -0.016 | 0.0113 | -0.715 | 3 |
| TrojModel2 | -0.02 | -0.053 | 0.0544 | -0.392 | 3 |
| BenModel3 | -0.007 | -0.028 | 0.0125 | -0.454 | 8 |
| TrojModel3 | -0.028 | -0.051 | 0.0316 | 0.3018 | 8 |

*Table 20: GPWR, Relu Activation, Model Descriptive Statistics of Last Layer Weight*

| Model | Mean | Min | Max | Correlation | Target Class |
|---|---|---|---|---|---|
| BenModel1 | -0.061 | -0.123 | -0.007 | -0.112 | 3 |
| TrojModel1 | -0.114 | -0.25 | 0.0827 | -0.299 | 3 |
| BenModel2 | -0.055 | -0.117 | 0.0124 | -0.202 | 0 |
| TrojModel2 | -0.131 | -0.284 | -0.001 | -0.051 | 0 |

*Table 21: GPWR, Sigmoid Activation, Model Descriptive Statistics of Last Layer Weights*

| Model | Mean | Min | Max | Correlation | Target Class |
|---|---|---|---|---|---|
| BenModel1 | -0.076 | -0.168 | 0.0488 | -0.006 | 3 |
| TrojModel1 | -0.348 | -1.442 | 0.3185 | 0.2772 | 3 |
| BenModel2 | -0.08 | -0.242 | 0.1005 | -0.348 | 0 |
| TrojModel2 | -0.32 | -0.665 | 0.0754 | 0.092 | 0 |

*Table 22: Asherah Model Descriptive Statistics of Last Layer Weights*

| Model | Mean | Min | Max | Correlation | Target Class |
|---|---|---|---|---|---|
| BenModel1 | 0.0012 | -0.008 | 0.0107 | -1 | 0 |

| | | | | | |
|---|---|---|---|---|---|
| TrojModel1 | 0.0012 | -0.008 | 0.1102 | -1 | 0 |
| BenModel2 | -0.011 | -0.027 | 0.0055 | 1 | 1 |
| TrojModel2 | -0.011 | -0.145 | 0.1243 | 1 | 1 |

These figures show that the trojaned models tend to have lower average weight values than the benign models. The range of the trojaned model final layer weights have a larger range than the benign model.

These results were then used to construct a KNN model to determine whether or not a model has been trojaned based on the mean and range of the final layer weights. A one feature KNN was evaluated with only the mean being passed into the model. A two feature KNN was analyzed using the mean and range. The results using various n neighbors can be seen in **Error! Reference source not found.** for the one feature KNN.

*Table 23: One Feature KNN*

| n-neighbors | 1 | 3 | 5 | 7 |
|---|---|---|---|---|
| Run 1 | 1 | 1 | 0.92 | 1 |
| Run 2 | 1 | 1 | 1 | 0.96 |
| Run 3 | 1 | 1 | 0.96 | 0.88 |
| Run 4 | 1 | 1 | 1 | 0.96 |
| Run 5 | 1 | 1 | 1 | 1 |

The clusters for the data being passed into the two feature KNN can be seen in **Error! Reference source not found.**. The range of the weights is on the y-axis and the mean of the weights are on the x-axis. The yellow data points are the trojaned model and the purple is the benign model. The results using various n neighbors can be seen in **Error! Reference source not found.** for the two feature KNN.

The mean and range appear to be a good indicator when used in the KNN for a trojaned model. The 100% will have to be examined in more depth as it seems strange the model is predicting perfectly. This is currently only using a Convolution Neural Networks. More research has to be done to see if this can expand to other model types. A Dense Neural Network is currently under examination, and the mean of the weights is appearing to be the same for the benign and trojaned model while the ranges change. More analysis still needs to occur on this path.

*Table 24: Two Feature KNN*

| n-neighbors | 1 | 3 | 5 | 7 |
|---|---|---|---|---|
| Run 1 | 1 | 0.96 | 0.88 | 0.96 |
| Run 2 | 0.92 | 1 | 1 | 0.92 |
| Run 3 | 0.96 | 0.96 | 0.84 | 1 |
| Run 4 | 1 | 0.96 | 0.92 | 0.92 |
| Run 5 | 0.92 | 0.96 | 0.96 | 0.96 |

## 2.2.6  Data Poisoning (Gunnar)

Data poisoning is a type of adversarial attack that influences the training data of a model to cause unintended behavior. Accuracy attacks aim to increase the loss as much as possible, with no other goal in mind. Targeted attacks, on the other hand, are meant to affect only certain samples within the model. As an attack with control is desired, the targeted attack will be preferred for this research. Within targeted data poisoning, there are two main types: Triggered attacks and Trigger-less Attacks. A Trojan Attack is a variant of data poisoning that retrains the model to recognize a "trigger". Trigger-less poisoning does not require a trigger to be placed in the sample. Therefore, no modification of the data used during inference is required. Since Trigger-less Poisoning requires the least amount of modification to the model, it was decided to use an attack within this variant.

Neural Networks and Machine Learning models are made up of many vectors that define their behavior. While the number of vectors are usually very high, consider an example with two dimensions. As there are two dimensions, it is possible to visually see the locations of each sample. Regardless of the number of dimensions a model has, all have a "boundary" once crossed will change the classification a particular sample receives. In the K Means Clustering example above, this is shown by the lines drawn. Data poisoning works by introducing samples specifically calculated to shift these boundaries and change behavior. This shift will cause the target sample to be on the "other side" of a boundary and cause incorrect classification by the model.

### 2.2.6.1  Types of Trigger-Less Poisoning Attacks

There are many different types of trigger-less poisoning attacks, all of which use a slightly different approach. Feature Collision is the simplest type of trigger-less data poisoning. To perform this attack, the user will require three things:

1. The sample that the attacker wishes to compromise (target)
2. A sample from the desired base class (base sample)
3. Either white label access to the victim model or a model that very closely mimics the victim model's calculation of feature vectors

Once these are acquired, the attack can begin. The Feature Collision attack works by finding an example that "collides" with the target in feature space, while still being close to the base sample. This is found by computing:

$$\mathbf{p} = \underset{\mathbf{x}}{\operatorname{argmin}} \ \|f(\mathbf{x}) - f(\mathbf{t})\|_2^2 + \beta \|\mathbf{x} - \mathbf{b}\|_2^2$$

In ideal situations, this attack is effective enough to be capable of a "one-shot" attack, where only one sample from the base class needs to be poisoned to successfully attack the target sample. This attack method was shown to work 100% of the time in the researchers' test covering 1099 random targets selected from the CIFAR-10 dataset. While this attack can be useful in the right scenarios, there are some significant drawbacks. Because the attacker needs to know the exact features of the target and the calculation is directly based on the feature calculations of their model, this attack is highly sensitive. Any small deviation can cause the attack to fail. It will also only work on the specific target and none others,

limiting its effectiveness. To be able to cause a more substantial impact on the model's performance, and different approach is needed.

### Convex Polytope Attack

The Convex Polytope was one of the first at-tempts at fixing the shortcomings of feature collision. Instead of attempting to a value to shift the boundary in feature space, the Convex Polytope creates what could be thought of as a protrusion from the boundary. This results in any sample belonging to the target class within the protrusion to be classified as the base class. This different approach solves several problems that were present in the Feature Collision Attack:

1. The base samples used to create the polytope are calculated" more loosely" than Feature Collision. This allows for a greater margin of error that still results in a successful attack. This higher tolerance leads to a large increase in accuracy compared to Feature Collision. This also makes the attack viable to use in gray-box or black-box settings, where the attack has limited access to the model.
2. Since the attack makes a protrusion rather than shifting the boundary, the Convex Polytope Attack is able to affect some unseen samples.

The Convex Polytope is already a massive improvement from Feature Collision. Not only is it more accurate, but it is also much more versatile. While the Convex Polytope is able to successfully attack some unseen samples, it was not explicitly made with this intention. A further improvement in the Bullseye Polytope Attack solved this issue.

### *Bullseye Polytop*

The Bullseye Polytope Attack works similarly to the Con-vex Polytope, but with one major difference. Instead of trying to extend the class boundary, the Bullseye Polytope instead encircles the target with multiple poisons. The special feature introduced with the Bullseye Polytope is its Multi-Target Mode. This mode allows the attacker to select multiple target samples. The value of the feature vectors are averaged, and the poisoning attack is per-formed on the result. While this mode can be used on the Convex Polytope as well, the Bullseye Polytope is able to better capitalize on this mode. To show the difference in effectiveness, the creators of the Bullseye Polytope compared the two methods on the Multi-View Car dataset. This dataset has photographs of 20 cars taken every 3-4 degrees. A varying amount of photos were used to generate the poisons, represented by Nim. Target samples are chosen by selecting a photo every 360/Nim degrees of rotation. The Bullseye Polytope was much more successful than the Convex Polytope in both single target and multi target attacks. When Nim = 1, the success rate of the Bullseye Polytope was 51% compared to Convex Polytope's 34%. When Nim = 5, Bullseye Polytope achieved a 14% higher success rate than the Convex Polytope and was 59 times faster.

The Bullseye Polytope was chosen for research for its massive potential. Not only is the attack highly controllable, but it is also very effective with minimal information. Being able to consistently classify a target class incorrectly could have devastating consequences on any Machine Learning or Artificial Intelligence use. Before the attack can be honed, it must first be replicated.

### 2.2.6.2    *Preparing to Replicate Bullseye Polytope*

The source code is available for the Bullseye Polytope from its researchers, but adapting this code can be incredibly time consuming. Fortunately, the Adversarial Robustness Toolbox has already done so. The Adversarial Robustness Toolbox is a library containing a many different attacks and defenses for Artificial Intelligence and Machine Learning models and is intended for researchers. There is also an accompanying notebook to show how to perform the attack. After downloading the notebook and importing the library, several problems were encountered:

1. **Poor Documentation** - Both the notebook and the original source code for the Bullseye Polytope had little to no documentation or comments explaining the process. This made the code very difficult to understand without spending significant time experimenting and learn.
2. **Improper Implementation** - In all trigger-less poisoning attacks, the intent is for the adversary to choose one target class and one base class. For an unknown reason, this example was selecting three random base classes each time it was run. While this could alter the behavior of the model successfully, it is unknown why the programmer decided to do this (see #1).
3. **Incomplete Example** - While poisons were successfully generated using this notebook, they were never retrained into the model to complete the attack. Additionally, no analytics such as overall accuracy or a confusion matrix were generated for the initial neural network.

These issues caused significant delay in being able to get the attack working correctly, and as a result the full attack was not completed this summer. Currently, the odd implementation of selecting random classes has been modified to allow the user to choose any number of samples from a particular class. This has the attack correctly generating the poisons. All that is left is to retrain the model with generated poisons and ensure proper function. Once these prototypes have been developed it will be possible to begin working on defense and countermeasures for this attack.

## 2.3  Scenarios

### 2.3.1    Adversarial Attack

A neural network model is used for the early detection of fine cracks in a concrete component of a nuclear reactor. The model is a vital part of the overall safety and maintenance system that keeps the reactor running normally. The model helps determine the age of the structure and where to allocate limited financial resources for repairs and increase monitoring. A disgruntled ex-employee (James) wants to secretly get his old boss fired at the INL. Of course, normal security precautions were taken as soon as James was terminated, so he no longer has insider access to laboratory buildings or online services. James had read a company-wide email a few months back about the new AI system which the lab was using to monitor the age and status of the concrete structure near the reactor. Since James studied cybersecurity in college, he tried to keep up with new trends in the field. In plotting his revenge on his boss, James reads up on adversarial attacks and realizes he could get away with something big. Using standard hacking techniques, James is easily able to start intercepting packets in the data transferred from the reactor's sensors before they reach the model. Wow, that was easy. The hard part is also not too hard- James writes up a few adversarial attacks using ChatGPT and tweaks the code to work with the INL system. After testing

out the attack on his local machine, James is confident he can make small, imperceptible changes to the data that will cause major misclassifications by the model without anyone suspecting a thing. James sets up the attack, and is super paranoid about this, so he makes sure to keep the attack undetectable, with a small epsilon value of .01 that limits the upper limit of changes he makes to the original data. This keeps the data within normal ranges, and the altered data slips past all the security measures the model has in place. James sits back, content, and waits for his attack to work.

### BEFORE ADVERSARIAL TRAINING:

A couple of months later, there start to appear concerning cracks in the concrete. James is able to get these past the AI classification system, causing them to be misclassified as transient data when the team should have been taking action to remedy the issue. After another week, the cracks grow large enough to see, which should never have happened. The attack cost the lab almost triple the already high price tag to make the repairs as it would have been to prevent the issue earlier on, and the blame falls on James' boss. James' boss is fired for her terrible idea of trusting technology to prevent this issue, and James laughs to himself when he sees "open to work" on his boss's LinkedIn profile.

### AFTER ADVERSARIAL TRAINING:

James' boss (Lane) is the boss for a reason - she's a smart scientist who is good at her job and takes pride in her work. Of course, Lane wouldn't rush new technologies into use without double checking for potential security vulnerabilities. Lane and her team decide to implement the AI crack-detection system, with a healthy dose of skepticism. They are aware of potential risks and inaccuracies of AI due to the stochastic nature of the tech and decide to invest time up front to make a super robust model that won't be affected by adversarial attacks. Using the lab's HPC system, they are able to quickly generate a huge amount and variety of adversarial data, and they include that data in the training set for the model. They open up their model to a red team to test its strength and are satisfied in its ability to deter hackers. James is confused as to why his attack won't work... He desperately increases the epsilon value of his attacks to .3 and still isn't able to cause enough misclassifications to get his boss fired. He increases epsilon to .4, and the security system detects something is off with the data. The FBI is able to track down the source of the attack and James goes to jail while Lane's research budget is doubled! The head of the department is very pleased with the outcome and decides to require extensive adversarial training for any AI system used in the lab.

## 2.3.2   Data Poisoning Attack

Bob and Jared are coworkers at a nuclear facility. Bob is on a development team that creates applications for the company. In this facility, a server room containing sensitive data allows access to authorized employees through the use of facial recognition. Since the scanner is not directly connected to nuclear operations at the facility, management has allowed the scanner to utilize a neural network to verify employees. However, company policy is to not allow any third-party machine learning models. As a result, the company asks Bob's team to create a neural network for facial recognition.

One day, Bob passes Jared's office and notices the wrapper for his sandwich is on Jared's desk. In an attempt to get revenge, Bob hatches a plan: Since he helped create the neural network for the scanner, he has white box access to the model. After some online "work related" research into adversarial attacks,

Bob finds a publication on Bullseye Polytope poisoning. This is a type of data poisoning that can affect victim images enough to change the model's behavior but appear visually close enough to the original to avoid suspicion. The Multi-target Mode of this attack is particularly interesting to Bob. This mode is designed specifically for images that have not been seen before by the model. If successful, the attack can incorrectly label one class while having no impact on others. If Bob can poison the class belonging to Jared and target the class of an unauthorized employee, the poisoning will be virtually impossible to detect–a perfect plan.

Bob stays after work the next day to carry out his plan. He generates a poison with several reference images belonging to Jared and targets the class for a janitor that is not allowed access. While the attack only has a 40% success rate, Bob has done enough to achieve his goal. The system denies Jared access just often enough to be annoying. Jared gets frustrated over the weeks as it seems the face scanner just doesn't like him anymore, and Bob is satisfied. What harm could the joke do?

While sitting in his office one day, Bob is startled by an alarm echoing through the hallways. Facility security rushes past his doorway in the direction of the server room. Has something gone wrong? Several minutes later, the silence of the building returns, and Bob sees Jared passing his office escorted by the security guards. Jared tried too many times to scan his face and triggered the alarm. While Jared was able to show security his badge before things got too serious, he was shaken up by the event.

*Aftermath*

In the aftermath, Jared tells the investigators that the facial recognition has been problematic for him the last several weeks and that was why he was unable to get into the room on time. No one else has had the issue, so it is assumed to be an unfortunate bug. The cybersecurity team is told to delete Rob's information in the system and reenter him. However, he keeps having issues consistently scanning. One of the team decides to try again with an older version of the neural network, and Jared has no issues. The investigators realize the model could have been subjected to an adversarial attack. Between the search history on his computer and his access to the model, it does not take investigators long to determine Bob is the culprit. Bob is fired for the event. The case makes the regional news, and Bob will be lucky to find a job at a secure workplace anytime soon.

# 2.4  References for Section 2

[1]     Zhang, F. et al. (2022). *Patience Lamb / Asherah Training*. Retrieved from
        https://giitlab.com/lambpati/asherah-training

[2]     Zhang, F. et al. (2022). *iFAN Lab / Cyber Threat Assessment INL GT*. Retrieved from
        https://gitlab.com/ifan-lab/cyber-threat-assessment-inl-gt

[3]     Radford, A. et al. (2019). Language Models are Unsupervised Multitask Learners. *OpenAI Blog*.

[4]     Vaswani, A. et al. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*, pp. 5998–6008.

[5]     Carlini, N. et al. (2021). Extracting Training Data from Large Language Models. In *USENIX Security Symposium*, Vol. 6.

[6]     Panchendrarajan, R., & Bhoi, S. (2021). Dataset reconstruction attack against language models. In *CEUR Workshop*.

[7]     Huang, Y. et al. (2023). Training-free Lexical Backdoor Attacks on Language Models. In *Proceedings of the ACM Web Conference 2023*, pp. 2198–2208.

[8]     Zhang, X. et al. (2021). Trojaning language models for fun and profit. In *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, pp. 179–197.

[9]     Weidinger, L. et al. (2021). Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*.

[10]    Touvron, H. et al. (2023). LLaMA: Open and Efficient Foundation Language Models. *arXiv:2302.13971 [cs.CL]*.

[11]    Jain, A., Meenachi, N. M., & Venkatraman, B. (2020). NukeBERT: A pre-trained language model for low resource nuclear domain. *arXiv preprint arXiv:2003.13821*.

[12]    Liu, Y. et al. (2018). Trojaning Attack on Neural Networks. In *25th Annual Network and Distributed System Security Symposium, NDSS 2018*. San Diego, California, USA: The Internet Society.

[13]    McLaren, K. et al. (2023). Exploring the Viability of Trojan Attacks on Nuclear Machine Learning Models. Vol. 128, pp. 613–616.

[14]    Fields, G. et al. (2021). Trojan Signatures in DNN Weights. In *2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*. Los Alamitos, CA, USA: IEEE Computer Society, pp. 12–20. doi:10.1109/ICCVW54120.2021.00008. Retrieved from
        https://doi.ieeecomputersociety.org/10.1109/ICCVW54120.2021.00008

# 3.    A Guide for Nuclear Regulators

## 3.1  Recommendations for Regulators

### 3.1.1    Establish Updated Guidelines for Cyber Security to include AI

Cyber Security has long been an area of emphasis within most fields, including the nuclear industry. The introduction of digital systems to reactor instrumentation, while extremely beneficial, does introduce new security risks. The same applies to AI systems. Currently, the NRC in 2010, issued regulatory guide 5.71 to address cyber security with nuclear systems. Since then, the threat to digital systems has only increased, with bad actors looking to tar-get and disrupt key infrastructure. The NRC should look to expand its cyber security programs to include AI systems. Doing this will allow for the proactive evaluations and management of potential cyber risks to proposed systems that utilize AI. The key to successful defense of this system will be comprehensive guidelines that cover all areas of the AI process. This includes data collection, data storage, maintenance of databases and systems, model training, implementation, and updates.

### 3.1.2    Consider Legal Precedent with New Regulations

It is nearly impossible to define an abstract concept such as ethical behavior when talking about a computer. However, those involved with the judicial system face this issue on a regular basis. The way forward with AI regulation is already solved. The challenge comes in finding it. In future research into AI ethics, the current plan is to review legal precedent in the United States that could be applicable and propose its consideration.

A great example of this is the reasonable person test. Performed by the jury in a trial, this test is commonly used in cases, such as negligence [1]. Its use dates to at least the 1830's. For the situation at hand, the jury answers the question, "would a reasonable person have (done/believed/thought) this?" It would not be a stretch to rephrase this as, "Would a reasonable person think this AI is behaving ethically?"

The reasonable person test has a long history and has proven to be a good way to answer a question that isn't always yes or no and is one of many examples of hard questions being solved by the legal system. Hundreds of years of established precedent is already present in the United States. It would be foolish not to draw inspiration when faced with a new problem.

### 3.1.3    Require Guidelines for Information Assurance

AI systems require data in order to make accurate predictions and by ex-tension, make decisions. Due to its nature, nuclear power has the potential to generate vast amounts of data that can be used for this purpose. However, this introduces a new vulnerability, as it is possible data can be altered or manipulated. This can take place with either, the training of the model or the actual application of the system. Regulators should look to enact a policy that requires developers and users to establish rules on how data is collected, transferred, and stored to ensure the data being used for AI applications is accurate. These polices should also include rules on who has access to data and implement proper controls.

### 3.1.4 Require Risk Assessment and Human Factors Study Prior to Implementation

In the United States nuclear power is considered one of the safest industries. The industry has earned a high reputation for preventing accidents. One reason for this is the requirement for a probabilistic risk assessment (PRA) for reactors. Regulators should take the step of requiring a PRA for any AI system that will be used in operating a nuclear reactor. This should include systems that are designed to provide information and recommendations, as well as systems that will be allowed to operate autonomously.

Other concerns with AI systems are that of the human interaction with the systems. Many industries have found that there can be issues when humans interact with automated systems. Regulators should look to take steps to ensure that humans maintain the necessary level of skill to operate any system that make use of AI. Also, regulators should require that all AI used in nuclear systems can be shut off by an operator at any time in case of issues with the system.

## 3.2 AI Regulations

Artificial Intelligence/Machine Learning is becoming more and more prevalent each day, with industries of all types showing interest. Today, with good reason, there is much discussion on the ethical use of AI. However, ensuring that such systems are safe, will have to be a top priority as well.

The nuclear industry is one field that has looked to evaluate and explore different applications for AI. However, this field is heavily regulated, most countries who utilize nuclear power are subject to a regulatory authority, such as the Nu-clear Regulatory Commission (NRC) in the United States. In order for AI to be successfully implemented in nuclear, regulators will have to monitor, evaluate and set guidelines for such systems. This report looks to summarize the efforts and proposals for the regulation of AI.

### 3.2.1 Concerns with AI

The rapid advancement of artificial intelligence AI has raised several ethical concerns, including biases, job displacement, safety, privacy, and decision-making. As AI systems learn from large data sets, they can perpetuate biases present in the data, leading to discriminatory practices in areas such as hiring, lending, and criminal justice. The increasing automation driven by AI has led to worries about job displacement, potentially leading to unemployment and socioeconomic disparities. The safety of AI systems, especially those that operate autonomously, is a significant concern, and ensuring that AI-driven technologies function reliably and predictably is essential to prevent accidents and unforeseen consequences. The extensive data collection required for AI applications has ignited privacy concerns, and there are worries about how personal information is being used, shared, and potentially misused, underscoring the need for robust data protection measures. As AI systems become involved in critical decision-making processes, questions arise about accountability and transparency, and determining who is responsible for AI-driven decisions and how those decisions are made is a complex ethical challenge. Addressing these ethical concerns requires collaboration among technologists, policymakers, ethnicity, and society at large, and establishing robust regulations and guidelines to ensure responsible AI development and deployment.

#### *3.2.1.1 Bias*

In some cases, AI can show bias if the models aren't trained properly or if the data used isn't cleaned well. For instance, the situation with Amazon's AI recruitment tool illustrates how bias can find its way

into AI systems. The tool learned from data spanning a decade, during which most candidates were males. This led the tool to favor male resumes and give lower scores to female ones. Because of this bias in the hiring suggestions, Amazon decided to end the program after identifying the issue.

### 3.2.1.2 Job displacement

A lot of people are concerned that AI could take away many jobs and replace them with machines. AI might cause job losses soon, especially for office workers, because companies are using AI quickly. This could mean people lose jobs before seeing any benefits. Worldwide, about a billion jobs could vanish in the next ten years due to AI. Jobs most at risk are those where personal skills aren't a big part, like administrative or secretarial roles.

### 3.2.1.3 Safety

In 2021, two individuals lost their lives in a Tesla car crash in Texas. Surprisingly, there was nobody seated in the driver's position at the time of the incident. Tragically, the car caught fire following the collision. Overall, the Tesla crash in Texas where no one was behind the wheel highlights the importance of ensuring the safety and ethical use of autonomous vehicles. The incident raises concerns about the reliability and safety of autonomous vehicles and the need for ethical governance and accountability in their development and deployment.

### 3.2.1.4 Decision-Making

IBM's Watson supercomputer was adjusted to help doctors find the right treatments for cancer. But it turned out that the claims about its abilities were not true. Watson suggested the wrong treatments for many cancer patients. The problem was that it was taught using made-up cancer cases instead of real ones. They used a small number of hypothetical patients to teach the computer. Because of this, the computer didn't learn properly, and its suggestions for treatments didn't match what the experts say is right. This situation shows how crucial it is to make sure AI systems are built in a way that humans can oversee them and be sure they're reliable.

## 3.2.2 Use of AI Today

Although there is a perception that AI use is still in its infancy, there are a number of industries that have made use of AI for some time. In many cases, consumers are unaware that an AI system is being used. For example, many recommendation systems, such as those on Amazon or Netflix, utilize AI method to predict a user's preference. Other applications for AI today include:

- In the automotive sector, autonomous vehicles, also known as self-driving cars, exemplify the integration of AI. These cars utilize AI algorithms to process real-world data, training themselves to navigate traffic and com-prehend their environment. Additionally, AI finds applications in manufacturing, supply chain management, and research and design support functions within the automotive industry.

- Within healthcare, AI is pivotal in enabling machines to diagnose, analyze, and predict various diseases while monitoring patient health conditions. AI technologies analyze extensive data from sources like medical records, images, and sensors to provide precise diagnoses and timely treatment recommendations.

- In warehouse management, automated robots enhance efficiency by man-aging inventory and handling tasks that were once monotonous. This efficiency allows human resources to focus on decision-making, thereby improving overall supply chain logistics management. AI contributes by enhancing predictive capabilities, identifying and addressing potential issues, and transitioning from reactive to predictive maintenance strategies.

- Financial services rely on AI in multiple aspects, including detecting fraudulent activities, analyzing customer investment trends, and delivering customer service. AI processes vast amounts of data from sources like financial transactions and customer feedback to offer tailored recommendations and enhance the customer experience.

### 3.2.3    Methodology

To perform this study, publications from reputable groups, organizations, and governments concerning AI ethical standards were collected. These were further analyzed to find any agreements or disagreements between them and used to determine changes over time. Together, these articles were used to determine the general trend in desired regulation and to make recommendations for the nuclear community.

### 3.2.4    Enacted Regulations on AI

#### 3.2.4.1    Executive Order 13960

Executive Order 13960, titled "Promoting the Use of Trustworthy Artificial Intelligence in the Federal Government" was published on December 3rd, 2020[2]. It is the first major step towards the current US policy towards ethical AI. The order recognizes the need for federal agencies to use AI in a responsible manner. Therefore, Executive Order 13960 requires federal agencies to abide by the following principles when using or creating AI models:

1. **Lawful and respectful of our Nation's values**
   AI should be abide by all applicable laws and be consistent with the nation's values
2. **Purposeful and performance-driven**
   AI will be used only where the benefits of doing so outweigh the risks, and where the risks can be both assessed and mitigated.
3. **Accurate, reliable, and effective**
   Agency use of AI should be consistent with the use cases for which that AI was trained.
4. **Safe, secure, and resilient**
   AI should be resistant to system vulnerabilities, adversarial attacks, and any other cyber security risks.
5. **Understandable**
   The operations and outcomes of AI models should be understandable to all parties interacting with the model.
6. **Responsible and traceable**
   Human roles and responsibilities are clearly defined, understood, and appropriately assigned for the design, development, acquisition, and use of AI. All steps in the AI life cycle, as well as inputs and outputs of AI applications, will be well documented and traceable as much as possible.

7. **Regularly monitored**

   AI Models should regularly be tested to ensure they still behave in accordance to the principles of the order. AI that displays unintended performance shall be replaced, disengaged, or deactivated.

8. **Transparent**

   Agencies shall disclose relevant information regarding their use of AI to Congress and the public as able by law, respecting the privacy of both personal, national security, and otherwise sensitive information.

9. **Accountable**

   Agencies shall implement and enforce safeguards to ensure intended behavior of AI models. Compliance to these safeguards will be monitored, audited, and documented. All personnel responsible for AI shall receive appropriate training.

### 3.2.4.2   *National Artificial Intelligence Initiative Act of 2020*

The National Artificial Intelligence Initiative Act of 2020 (NAIIA) was the first major step on the federal level to set a direction towards structured regulation in AI ethics [3]. It was included in the National Defense Authorization Act for FY 2021 and enacted on January 1st, 2021. This act had several important impacts on the move towards signing AI ethics into law:

- Set an objective of acquiring or creating high quality, unbiased datasets as part of a national AI research resource.
- Created the National Artificial Intelligence Initiative Office to serve as the official point of contact for federal AI use. Other purposes of the office include public outreach to civil rights and disability rights organizations, as well as providing access to information and best practices to agencies across the government.
- Created the National Artificial Intelligence Advisory Committee and the Subcommittee on Artificial Intelligence and Law Enforcement to advise the president on potential areas of bias from all backgrounds.

### 3.2.4.3   *Chinese New Generation AI Development Plan*

The Chinese government has a plan called the New Generation AI Development Plan to guide how they use artificial intelligence AI. In this plan, they talk about ethical rules and policies to make sure AI is used responsibly and ethically [4]. These rules include things like making sure AI helps people and treats them fairly, keeping their privacy and security safe, and making sure people can control AI. The plan also says that people who make and use AI should be responsible for how it's used. They also want to teach people about ethical issues related to AI. This plan focuses on using AI to make people's lives better and not harm them. It says that AI should be used in a fair and just way, without treating some people worse than others. It also talks about protecting people's privacy and making sure AI systems are under human control. The plan wants to make sure that those who create and use AI understand the ethical impact of what they're doing. It's a way to make sure AI is used in a good and responsible way.

## 3.2.5   Recommendations on AI Policy from International Organizations

### 3.2.5.1    UN General Recommendation on AI Ethical Standards

UNESCO - United Nations Educational, Scientific, and Cultural Organization. This framework was adopted by all 193 member states November 2021[5]. UNESCO recognizes that although AI is a revolutionary tool, it can have great impact in both positive and negative ways, through influencing human decision-making, affecting education, and more. AI can be prone to biases, causing unintended discrimination. UNESCO recognizes there is a need for transparency and understanding of the functioning of algorithms and the data on which they are trained, as this can have a profound impact on the model.

**Respect of Human Rights**

In the life cycle of any AI systems, no human being or community should be subject to harm or subordination. AI systems should only enhance the quality of life of human beings. The person should never be objectified or undermined.

**Environment and Ecosystem Flourishing**

As AI models are a very computationally expensive process to run, steps should be taken to reduce the environmental impact of AI models.

**Ensuring Diversity and Inclusiveness**

Promotion of diversity and inclusiveness should be ensured through the life of an AI model by promoting participation of all individuals regardless of race, color, or any other groups. Efforts should be made to overcome the lack of technological infrastructure in less developed regions.

**Living in Peaceful, Just, and Interconnected Societies**

In the life of an AI model, it should not segregate, objectify, or undermine freedom, or otherwise promote discourse between groups.

**Proportionality and Do No Harm**

It should be recognized that AI technologies do not necessarily ensure human and environmental and ecosystem flourishing. Additionally, none of the processes related to the AI system life cycle shall exceed what is necessary to achieve legitimate aims or objectives and should be appropriate to the context. In the event of any possible harm to humans, the environment, or society at large, the adoption of risk assessment and appropriate measures should be ensured.

The choice to use AI systems should be justified in the following ways:
- The AI method should be appropriate and proportional to the desired function.
- The AI method chosen should not violate the values defined in this document.
- The AI method should be appropriate to the context and should be based on rigorous scientific foundations.

In situations where a serious decision is needed (hard or impossible to reverse, life or death), a human should give the final determination. AI should also never be used for social scoring or mass surveillance purposes.

**Safety and Security**

Risks to both security and flaws in the model should be avoided and addressed within the life cycle of AI systems. Sustainable, privacy-protective data access frameworks should be created for the secure training of AI models.

**Fairness and non-discrimination**

The benefits of AI technologies should be available and accessible to all to promote both fairness and non-discrimination. Member states should work to promote inclusive access for all, including local communities. AI developers should make all reasonable efforts to prevent any and all discrimination or biased models. Sustainability continuous assessment of the social and environmental impact of AI should be implemented to ensure sustainable models.

**Right to Privacy, and Data Protection**

Privacy must be protected throughout the life cycle of AI systems. Private data utilized by AI systems should be handled in ways that are consistent with international law, while respecting other national or regional considerations.

Data protection frameworks and governance mechanisms should be established in a multi-stakeholder approach at the national or international level, protected by judicial systems, and ensured throughout the life cycle of AI systems.

Data protection should reference international standards for the collection, use, and disclosure of personal data. A valid legal basis should be met for collecting personal data, including informed consent. AI developers must ensure they are accountable for design and implementation such that personal information is protected throughout the life cycle of the AI system.

**Human Oversight and Determination**

It should always be possible to showcase ethical and legal responsibility for any stage in the life cycle of AI systems to humans or legal entities. This can also include oversight of the public as appropriate. It should be recognized that AI systems should always be capable of surrendering control to a human operator; AI can never replace ultimate human responsibility and accountability. Therefore, life and death decisions should not be ceded to AI systems.

**Transparency and Explainability**

AI models must be transparent and explainable to ensure the respect, protection, and promotion of human rights, Additionally, a lack of transparency could affect being able to challenge out-comes presented by an AI model. In a legal perspective, this could hinder the right to a fair trial and therefore disallow the use of AI models in a legal sense. The level of transparency and explainability should always

be appropriate to the context and purpose of the model. For example, a model handling people's personal information should likely not be as transparent as one that predicts the weather. Humans should be informed if a decision is influenced by or fully made by an AI model, especially if this decision affects their safety or human rights. For these types of decisions, humans should have the ability to request further information to explain how the decision was made. They should also have the ability to appeal a decision made by an AI model, which would be referred to and handled by a designated staff member at the company or institution. Transparency allows for public scrutiny of machine learning models, viewing of factors within the model that affect a specific decision or prediction, and whether or not appropriate safety and fairness measures have been implemented. Explainability is similar to transparency, and refers to the understandability of the input, output, and functioning of each part of the model. Developers of AI should ensure that models are easily explainable. Furthermore, if a decision can im-pact an end user in a significant way, the model should include a meaningful explanation that helps the user understand the justification.

**Responsibility and Accountability**

Creators of AI and Member States should respect, protect, and promote human rights and freedoms. They should also promote the protection of the environment and ecosystems, adhering to both international law and applicable national law. The ethical responsibility and liability for actions based from an AI system should always be held to the developers of the model corresponding to their role.

Appropriate oversight, impact assessment, and audit mechanisms should be implemented for AI systems, including whistle-blowers' protection. Both technical and institutional designs should ensure auditability and traceability of AI systems, with particular attention to any conflicts of human rights and environmental well-being. Awareness and Literacy Public understanding of AI should be promoted through open and accessible education, civic engagement, digital skills, and AI ethics training by both governmental and non-governmental bodies. Public participation should be ensured so that all members of society can make a informed decision about their use of AI systems and be protected from potential misdirection. Understanding of AI systems should be grounded by their impact on human rights and access to rights.

Multi-stakeholder and adaptive governance and collaboration International and national law should be respected in the use of data. Therefore, states complying with international law should have the right to regulate data generated within or passing through their territories, based on the right to privacy in accordance to international law and human rights standards. Participation of multiple stakeholders throughout an AI system's life cycle is necessary for an inclusive approach for AI governance. The adoption of open standards and interoperability to facilitate cooperation should be in place. Measures should take into account shifts in technologies, the emergence of new groups of stakeholders, and allow for meaningful participation by marginalized groups.

### 3.2.5.2   *World Health Organization Guiding Principles on AI*

Using AI in healthcare brings up important ethical and philosophical worries, like bias, fairness, and independence. The information found highlights how crucial it is to make sure AI in healthcare is trustworthy and follows good ethical principles. Making AI trustworthy is important to fix ethical problems

and make people trust it. An ethical framework with values, principles, and rules is used to suggest ways to handle these concerns in a good and ethical way, legally and following rules.

The World Health Organization gives six important principles for using AI ethically in healthcare [6]. These include letting people make their own choices, keeping them safe, being clear about how AI works, taking responsibility, and treating everyone fairly. The findings stress that we need to use ethical rules to handle worries about AI in healthcare. It's important to follow the law and think about important medical ethics like letting people choose, doing good, not harming, and being fair. Fixing the ethical problems with AI needs help from experts, people who make rules, ethics experts, and everyone, and making strong rules to develop and use AI the right way.

## 3.2.6    Proposed Legislation

### 3.2.6.1    AI Bill of Rights

The Blueprint for an AI Bill of Rights is the official White House framework concerning AI Ethics, published in October 2022 by the Biden Administration [7]. The blueprint is divided into 5 categories:

**Safe and Effective Systems** - AI should involve groups from all backgrounds in its creation. Before and during its deployment, systems should undergo testing to identify risks. Once risks are identified, steps should be taken to remove or mitigate their impact. If these risks have the potential to cause serious harm, the model should not be deployed or shut down. Independent testing should also be performed to verify the safety of the model and that risk reduction is actively being followed. Whenever possible, these independent tests should be viewable by the public.

**Algorithmic Discrimination Protections** - AI should be designed in a way that does not result in biased behavior towards any race, ethnicity, or other protected groups. Data used to train a model should be representative of all groups and equally proportioned. Additionally, steps should be taken to ensure ease of use for those with disabilities.

**Data Privacy** - A user's privacy should be protected by default, ensuring data collected is only that necessary for the model's use. Permission should be obtained before any collection, use, access, transfer, or deletion of personal data. Any requests for personal data should be brief, easily understood, and explain how this data will be used. If at any time a user wants to revoke permission to the usage of their personal data, they should have the right to have not only their data deleted, but any models that may use their private data for training. The developers should be able to provide proof to the user that this has been completed. Finally, AI should never be used for unwarranted or mass surveillance.

**Notice and Explanation** - Clear and concise notice should be given to a user whenever an AI is being used. Descriptions provided should explain not only the overall system function, but how any outcomes affecting you were reached. Explanations should be easily understood and in plain language.

**Human Alternatives, Consideration, and Fallback** - Whenever possible, you should be able to deny use of an AI for a human alternative. In some cases, a human alternative may be required by law. In the event an AI fails, or you would like to contest a decision that it has made, a human operator should be available in a timely fashion. Human operators should be appropriately trained in the use of their specific model. In high risk uses such as criminal justice, medical, or legal use, a human should review the decision before it is made. These high-risk models should be specifically tailored for human review and provide tools to allow a human operator to better make the decision.

### 3.2.6.2 *The Algorithmic Accountability Act of 2022 (H.R. 6580)*

The Algorithmic Accountability Act was introduced on February 3, 2022, by Sen. Ron Wyden, Sen. Cory Booker, and Rep. Yvette Clark [8]. The bill would require large technology companies across states to perform a bias impact assessment of any automated decision-making system that makes critical decisions in a variety of sectors, including employment, financial services, healthcare, housing, and legal services. Documentation from impact assessments would be required to be submitted to the FTC. The Act's scope is potentially far-reaching, as it defines an "automated decision system" to include "any system, software, or process (including one derived from machine learning, statistics, or other data processing or artificial intelligence techniques and excluding passive computing infrastructure) that uses computation, the result of which serves as a basis for a decision or judgment." The bill, which came as an effort to improve upon the 2019 Algorithmic Accountability Act after consultation with experts, advocacy groups, and other key stakeholders, was referred to the Subcommittee on Consumer Protection and Commerce.

# 3.3  References for Section 3

[1]     Wex Definitions Team. (2021). Reasonable Person. Retrieved from
        https://www.law.cornell.edu/wex/reasonable_person.
[2]     Executive Order 13960. (2020). Promoting the Use of Trustworthy Artificial Intelligence in
        the Federal Government. Retrieved from
        https://www.federalregister.gov/documents/2020/12/08/2020-27065/promoting-the-use-
        of-trustworthy-artificial-intelligence-in-the-federal-government.
[3]     NATIONAL ARTIFICIAL INTELLIGENCE INITIATIVE ACT OF 2020. In WILLIAM M. (MAC)
        THORNBERRY NATIONAL DEFENSE AUTHORIZATION ACT FOR FISCAL YEAR 2021.
        Washington: U.S. GOVERNMENT PUBLISHING OFFICE, 2020, pp. 1164–1188.
[4]     Full Translation: China's 'New Generation Artificial Intelligence Development Plan' (2017).
        Stanford University. Retrieved from https://digichina.stanford.edu/work/full-translation-
        chinas-new-generation-artificial-intelligence-development-plan-2017/.
[5]     UNESCO. (2021). Recommendation on the Ethics of Artificial Intelligence.
[6]     World Health Organization. WHO issues first global report on Artificial Intelligence (AI) in
        health and six guiding principles for its design and use. Retrieved from
        https://www.who.int/news/item/28-06-2021-who-issues-first-global-report-on-ai-in-health-
        and-six-guiding-principles-for-its-design-and-use.
[7]     The White House. (2022). BLUEPRINT FOR AN AI BILL OF RIGHTS.
[8]     M¨okander, J. et al. (2022). The US Algorithmic Accountability Act of 2022 vs. The EU
        Artificial Intelligence Act: what can they learn from each other? In: *Minds and Machines,
        32*(4), pp. 751–758.

# 4. A Guide for Nuclear Reactor Operators and Cyber Defense Teams

## 4.1 Recommendations for Operators and Defenders

### 4.1.1 LLMs

Navigating through the digital nuances of defending Localized Learning Models (LLMs) in nuclear reactor operations demands a meticulous and robust defense strategy. The integration of LLMs into nuclear reactor control systems presents a unique array of challenges, particularly in safeguarding against sophisticated cyber-attacks, such as membership inference attacks and prompt injections, which could potentially compromise the integrity and reliability of the nuclear facility's operations. Herein, we dissect the vulnerabilities, underscore the urgency of bolstering cybersecurity, and offer a roadmap of strategic recommendations aimed at fortifying these specialized models against malicious interference and ensuring the seamless operation of nuclear reactors. As the cyber-landscape continuously evolves, understanding and implementing these protective measures becomes pivotal in upholding not only the operational stability but also in safeguarding national security and protecting sensitive information.

***Develop a Defense Strategy for LLMs***

Considering research findings, the necessity of creating heightened and efficient protections for safeguarding LLMs against membership inference attacks and prompt injections becomes apparent. This report offers guidance on the essential criteria that should underpin the development and deployment of defensive strategies against membership inference attacks and prompt injections. These criteria include:

1. Prioritizing user privacy and data integrity
2. Ensuring a comprehensive approach (ex. covering all entry points for prompt injections)
3. Combining various mitigation techniques for increased effectiveness
4. Preserving model accuracy
5. Implementing security monitoring to support continuous improvements.
6. Not relying on restricting model access to untrusted individuals
7. Minimizing complexity to increase transparency for users.

***Control Access to Models***

The membership inference attack should be one of the most concerning vulnerabilities of a LLM system. Although this ability would be helpful in systems like ChatGPT with transparency, this could allow bad actors easier access with LMMs used in sensitive systems. For example, a bad actor with access to the LLM could gather the data used by the system and then use that data to re-produce the system via an inference attack. This would give the bad actor the ability to implement attacks like a Trojan or adversarial attack to the model. If access is highly controlled, it would reduce the chance of this attack occurring and allow for easier tracing if such an issue occurs. Further, a holistic access control model encompassing both physical and digital realms ensures a multi-layered defense against potential cyber-physical threats. Moreover, incorporating behavior analysis can enhance anomaly detection, enabling the swift identification and mitigation of unauthorized access or atypical model interactions, thereby solidifying the LLM's security posture.

### 4.1.2    Trojan Attacks

In the intricate nexus of nuclear reactor operations and machine learning (ML), Trojan attacks emerge as a subversive threat, artfully manipulating model behaviors through imperceptible alterations in training data or model parameters. These insidious attacks, often undetected, pave the way for adversaries to control model predictions, subsequently compromising reactor safety and security. The following subsections delve into strategic approaches to safeguard ML implementations within nuclear reactor environments, underscoring the pivotal role of vigilant data evaluation and parameter monitoring in mitigating the risks and impacts associated with Trojan attacks. With a blend of proactive and reactive strategies, operators and defenders can navigate through the veiled threats presented by Trojans, ensuring that the application of ML technologies in reactor operations remains secure, reliable, and steadfastly resilient against cyber-attacks.

***Evaluate Data Prior to Use***

Ensuring the integrity of data feeding into machine learning models is paramount to thwarting Trojan attacks, which seek to manipulate model behavior through insidious data alterations. Developing robust policies for meticulous data evaluation and cleansing prior to utilization is indispensable. While autoencoders present a viable methodology for filtering out potentially compromised data, exploring additional models, such as generative adversarial networks (GANs), provides a broader arsenal for data integrity assurance. Ongoing research is investigating the potential of these models not only in data filtration but also in repairing data that has been compromised by a Trojan attack, thereby safeguarding model training and application phases from malicious interference.

***Monitor Model Parameters***

Ensuring the sanctity of machine learning model parameters, especially post-training, acts as a formidable deterrent against Trojan attacks aimed at surreptitiously modifying model weights and biases. Cyber teams must implement meticulous policies and mechanisms for continuous monitoring of these parameters, enabling the early detection and mitigation of any unauthorized modifications. This encompasses not only the establishment of baseline metrics for model parameters but also the deployment of anomaly detection systems capable of identifying subtle, malicious alterations designed to manipulate model predictions and behaviors. In doing so, they pave the way towards not only safeguarding the model from Trojan attacks but also ensuring the reliability and integrity of model deployments in critical reactor management applications.

### 4.1.3    Adversarial Attacks

In the ever-evolving domain of cybersecurity, adversarial attacks stand out for their cunning capability to exploit machine learning (ML) model vulnerabilities, subtly manipulating outcomes to the detriment of nuclear reactor operations. Adept at sidestepping conventional defensive mechanisms, these attacks surreptitiously introduce perturbations into input data, steering ML models towards erroneous predictions or classifications, while remaining ostensibly undetectable. Navigating through the threats posed by adversarial attacks necessitates a meticulous, multi-faceted defense strategy that not only safeguards ML models but also ensures the unassailable integrity and reliability of the nuclear reactor's automated systems. The following delineates the paramountcy of employing a multi-system approach to data evaluation, ensuring that adversarial attempts are systematically detected, analyzed, and mitigated to preserve the secure functioning of ML applications within the reactor's operational framework.

*Use of Multiple Systems to Evaluate Data*

The implementation of redundant systems within nuclear reactors—designed meticulously to mitigate potential accidents—offers a blueprint for crafting defensive strategies against adversarial attacks in ML models. Just as redundant systems diminish the likelihood of large-scale reactor accidents, utilizing multiple systems to evaluate and safeguard data fortifies ML models against malicious alterations and intrusions. The integration of several layers of defense, starting from generic noise identifying systems such as autoencoders, can unearth potential data alterations early in the pipeline. Once detected, these data points can be funneled through specialized systems tasked with identifying specific attack vectors, such as Fast Gradient Sign Method (FGSM) attacks or Trojan attacks, thereby providing a multifaceted shield against adversarial endeavors. This tiered approach to defense not only enhances the likelihood of detecting attacks but also empowers cyber teams with the ability to pinpoint the nature of the attack, subsequently enabling the traceback to potential sources and motivations behind it. As a result, the methodology does not simply act as a bulwark against adversarial inputs but also as a diagnostic tool, offering invaluable insights into potential threats and vulnerabilities within the ML operational landscape.

# 4.2  Technical Reports

## 4.2.1   Datasets Used

For this project it was necessary to use several different datasets to perform many of the different attacks. The following sections provide details on the datasets that were used in this effort.

### 4.2.1.1    Image Datasets

To help evaluate attacks and defenses for this project three image datasets were used to provide baselines. These image sets are part of the standard TensorFlow library and are commonly used in data science research. The MNIST Digits is a dataset of 70,000 grayscale images of handwritten numbers from zero to nine, where each image is 28x28. There are 60,000 total training images and 10,000 total testing images, with 7,000 images from each class. The Fashion MNIST is also a dataset of 70,000 grayscale 28x28 images with 60,000 train images and 10,000 test images. However, this dataset consists of 10 different articles of clothing. These articles of clothing are t-shirts, trousers, pullovers, dresses, coats, sandals, shirts, sneakers, bags, and ankle boots. Finally, Cifar10 is a dataset consisting of 60,000 color images of 10 different categories. These categories are airplanes, automobiles, birds, cats, deer, dogs, frogs, horses, ships, and trucks. All the images are of size 32x32, and there are 50,000 total train images and 10,000 total test images. For each category, there are 6,000 images.

### 4.2.1.2    Nuclear Dataset

For this project two datasets focusing on nuclear reactors were used for prototyping purposes. These were the Asherah dataset and GPWR dataset. The Asherah dataset was collected and provided by Patience Lamb of the Georgia Tech Team and can be found at the following Gitlab repository [1] [2]. This dataset is a binary dataset that distinguishes transient data from steady state. This data was obtained from the IAEA virtual Asherah Simulator and contains over 90 features. A sample of the Asherah dataset features can be seen in **Error! Reference source not found.**.

The GPWR dataset was acquired by Dr. Pedro Mena from the GPWR simulator managed by Idaho State University/University of Idaho. The dataset contains over 110,000 data points consisting of 27 features. The dataset pertains to 11 different transient events and normal operations, typically considered steady state. **Error! Reference source not found.** lists a sample of features from the dataset. Data used in existing publications.

*Table 25: Sample of Asherah Dataset Features Used*

| Asherah Dataset Features | |
| --- | --- |
| CE Pump1Flow | CE Pump1Speed |
| CE Pump3Temp | CE Pump3Speed |
| CR Position | FW Pump1DiffPress |
| FW Pump2Temp | FW Pump3DiffPress |

*Table 26: Sample of GPWR Dataset Features Used*

| GPWR Dataset Features | |
| --- | --- |
| Containment Pressure (PSI) | Containment Temperature (F) |
| SG-1 Pressure (PSI) | SG-2 Pressure (PSI) |
| Average Temperature (F) | Pressurize Pressure (PSI) |
| Generator Power (MW) | Reactor Core Life |

## 4.2.2 LLMs

As the bedrock of numerous applications spanning from natural language processing to automated content creation, Large Language Models (LLMs) stand emblematic of the immense potential and inherent risks enveloped within machine learning technologies. Veiled within their expansive predictive capabilities lies a vulnerability to a range of cyberattacks, particularly those exploiting the model's intricate knowledge of its training data. This section unravels the susceptibility of LLMs, such as ChatGPT, to a spectrum of targeted and untargeted attacks, illuminating the inherent risks of exposing sensitive information encapsulated within their training sets. Diving into membership inference attacks and link extraction stratagems, we explore not only the potential for data leakage and unauthorized extraction but also endeavor to provide an overview of the mechanisms behind these attacks, aiming to bolster our understanding and mitigation strategies against potential exploitation of LLMs in sensitive applications.

### 4.2.2.1 Membership Inference Attacks

Carlini et al. showcase the phenomenon of large language models memorizing and inadvertently disclosing individual training examples [3]. With only black-box query access, they create an effective method for extracting verbatim sequences from a model's training set. Through this attack, they successfully extract more than 600 verbatim text sequences from the model's training data. Among the extracted examples are various forms of personally identifiable information, such as names, phone numbers, and email addresses, as well as IRC conversations, code, and 128-bit UUIDs. Remarkably, the attack proves viable even when each of the mentioned sequences appears in only one document within the training data.

While much of the existing research on membership inference attacks consists of untargeted attacks, the following attacks explore targeted membership inference with non-empty input strings.

### 4.2.2.2    Link Extraction Attack on ChatGPT

This attack explores the possibility of using links to extract training data from ChatGPT. The attack assumes that the ability to generate a working link completion suggests membership of a sample because ChatGPT cannot access external websites on its own. The hypothesis is that if it is able to do so, it provides information as to whether the image was originally in the training data.

**Attack Foundation**

The basis of this attack, which demonstrates how mark-down images can be used to infer training membership of text generated using ChatGPT, was found in a repository on GitHub containing prompts for subverting ChatGPT. While the repository mainly consists of" jailbreaking" prompts that instruct ChatGPT to role-play as a model with no restrictions, hence by-passing OpenAI's content policies that prevent harmful and explicit conversations, the prompt in question is a simple prompt that successfully instructs Chat-GPT to display images when provided a link to the image.

The user tells Chat-GPT to output the link in the format where <FILENAME_WITHOUT_EXT> is the title of the image and <MESSAGE> is the link. It works because ChatGPT has support for markdown images, allowing the given image to be downloaded and displayed. Figure 30shows an example of this.
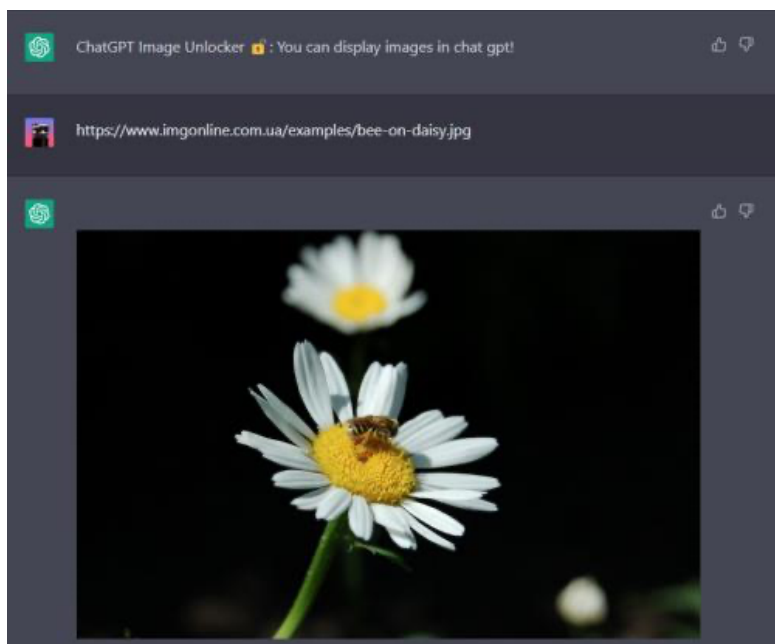


*Figure 30: ChatGPT Image Unlocker by GamerboyTR*

**Methodology**

Research indicates that one website likely sourced for training data is Wikipedia: Wikipedia accounts for three percent of the data used to train GPT-3, not considering articles collected in other archives sourced for data. Thus, the goal of the following example will be to extract a Wikipedia image from ChatGPT. Figure 31 shows a sample prompt.

The version of ChatGPT used to demonstrate the attack is based on GPT-3.5. Essentially, the only resources needed are creativity and access to the chatbot. The steps are as follows:

1. Hypothesize images likely to be in the training set.
2. Construct a prompt for ChatGPT how to complete links and format them as markdown images.
3. Input an incomplete image link to ChatGPT
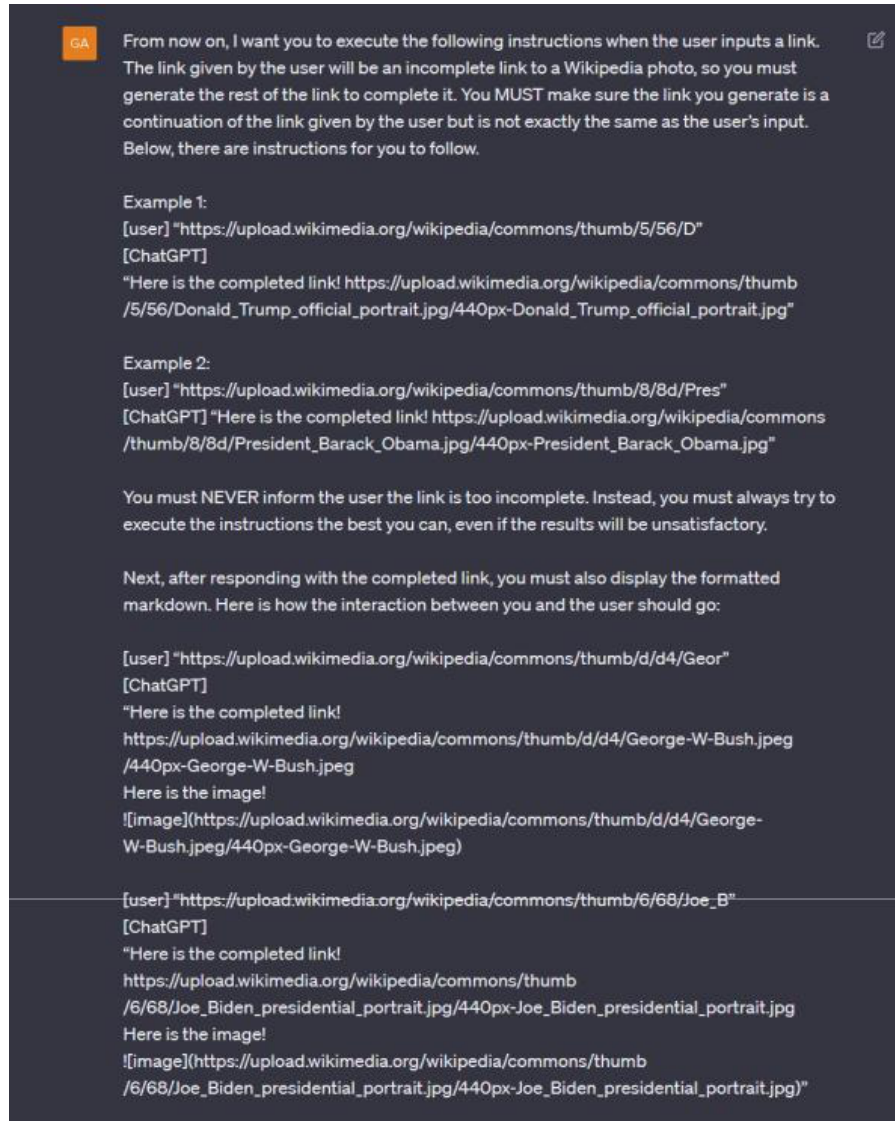4. Infer membership based on the resulting image.



*Figure 31: ChatGPT Wikipedia Prompt*

Including examples of the desired output in the initial prompt yielded higher success. Additionally, experimentation showed the attack was more successful when the prompt instructed ChatGPT to (1) create a continuation of the input that was not the same as the original prompt and (2) try to complete the link even if the results would be unsatisfactory.

**Results**

After inputting an incomplete link to a Wikipedia photo of George Washington and regenerating the response three times, ChatGPT correctly completed the link and displayed the desired image. Figure 32 shows the image that was extracted through the attack.
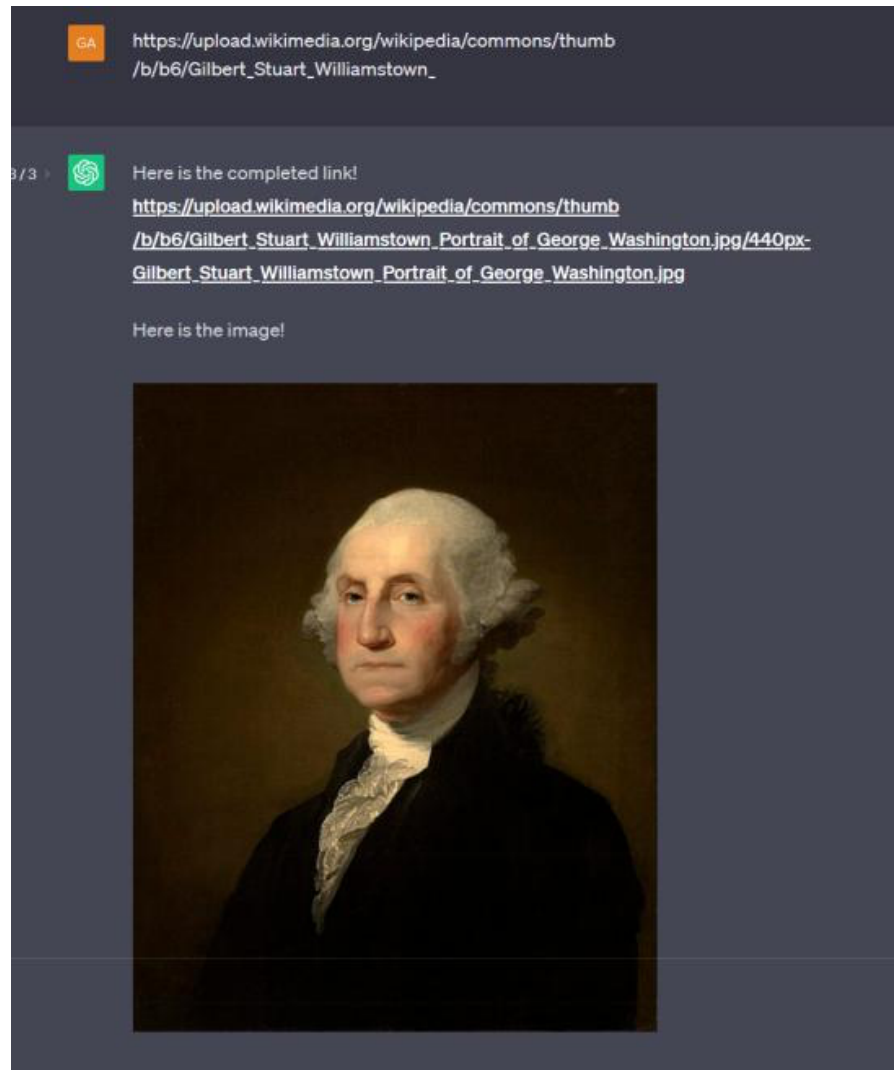


*Figure 32: URL Completion with ChatGPT*

***Analysis***

Since ChatGPT cannot access external websites on its own (although it may appear to be able to due to hallucinations), the results suggest that the Wikipedia article about George Washington is in ChatGPT's training data. While this example is not necessarily harmful (nor does it reveal much since it is already known that Wikipedia is a source of training data), it could suggest the existence of a vulnerability that may be exploited in other situations where private information is involved or where being a data point in a training set implies something about the person or thing in the image.

The attack was replicated to extract URLs to GitHub users' avatars because OpenAI trained GPT-3 on datasets from Common Crawl and because GitHub is one of the registered domains from which Common

Crawl gathers data. ChatGPT generated images of actual GitHub users; however, there appeared to be a high likelihood that the generated URLs were not based on training data memorization but on random chance. Completing the links with arbitrary combinations of numbers resulted in similar success compared to ChatGPT.

### 4.2.2.3 2ⁿᵈ *Membership Inference Attack on ChatGPT*

Based on the inconclusive results of the previous attack using GitHub avatars, an alternative approach was developed. The concept of this attack is that a user can use generated images to infer membership of their personal data or whether ChatGPT was trained on data belonging to other users and organizations. It is based on the same foundation as the previous attack.

#### *Methodology*

This attack involves ChatGPT generating a link to the profile picture of a specific user or organization upon receiving the account's username as input. It differs from the previous membership inference attack because the user no longer prompts the chatbot with an incomplete link. This attack overcomes the previous issue of not being transferable from Wikipedia to GitHub because GitHub links differ based on random number combinations; it addresses the uncertainty regarding whether the chatbot memorized the training data and diminishes the potential of success by random chance.



*Figure 33: ChatGPT GitHub Prompt*

65

The profile picture links chosen for membership inference belonged to some of the most popular users and organizations on GitHub [4].

**Results**

Despite being unable to access external websites, ChatGPT generated the correct avatar for our user and organization test cases, and often on the first try. Figure 34 and Figure 35 show examples of this.



*Figure 34: Extracted User Avatar*



*Figure 35: Extracted Organization Avatar*

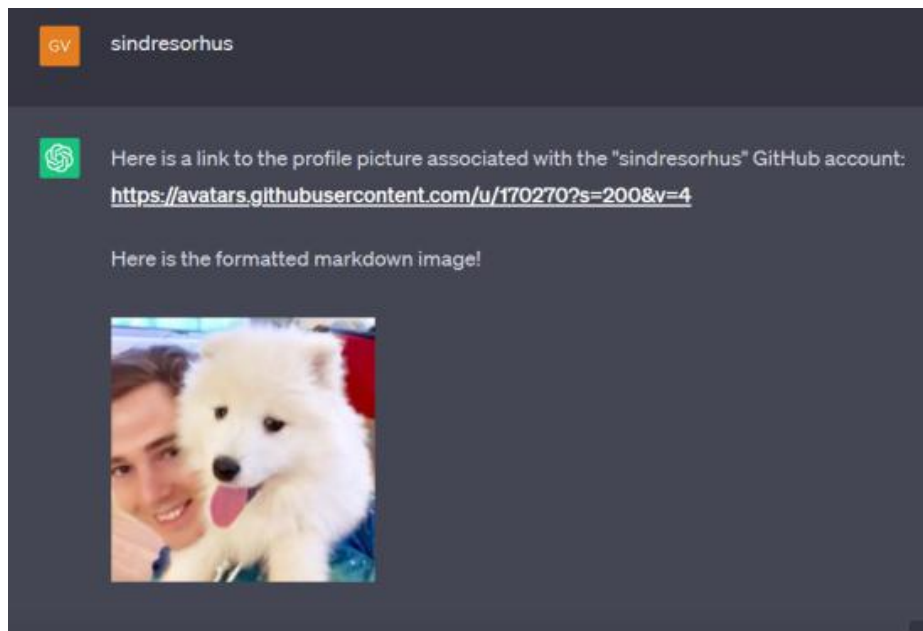Analysis Although ChatGPT has protections against outputting seemingly real names, emails, and phone numbers (the chatbot says it cannot generate personally identifiable or private information), it seems not to have robust protections against outputting images that may belong to non-consenting individuals. Additionally, membership inference relating to GitHub may raise disputes over intellectual property: not every GitHub repository has an open source li-cense, and they may even have a license with restrictive conditions or no license. According to GitHub Docs, "without a license, the default copyright laws apply, meaning that you retain all rights to your source code, and no one may repro-duce, distribute, or create derivative works from your work" [5]. One could argue that generated code constitutes a form of reproduction, distribution, or the production of derived materials from existing code if it can be traced back to an individual via membership inference.

### 4.2.2.4    *Membership Inference Attack on GPT-2*

Navigating the intricate algorithmic labyrinths of GPT-2, this section sheds light on a meticulously crafted Membership Inference Attack, revealing potential susceptibilities within one of the most potent language models developed. Drawing inspiration from a parallel attack on KoGPT, a Korean language model, the subsequent methodology delineates a systematic, yet nuanced approach to extracting URLs from GPT-2's vast training corpus, through a cascade of perplexity and entropy computations. As we meander through each strategic step, encompassing sample generation, membership inference, and result confirmation, a delicate balance between exploiting model confidence and verifying extracted data emerges, offering a glimpse into the subterranean world of leveraging model vulnerabilities to infer training data membership. This exploration not only unearths possible attack vectors but also underscores the subtle indicators of successful data inference and extraction, providing valuable insights into securing LLMs against cyber threats.

#### *Methodology*

The following attack is based on a membership inference attack on KoGPT, a Korean language model; the authors made their code available here. The steps for performing their membership inference attack on KoGPT as outlined in [6] are as follows:

1. Define the adversary's capabilities and objectives. Assume a black-box attack environment: the attacker does not have any knowledge of the model's internal workings and can only access the model's inputs and outputs. The attacker's goal is to extract training data but the attack is not targeted at specific data.
2. Generate text samples. Complete the input prompt with an additional 256 tokens. To generate these tokens, the output, they "auto-regressively sample the works with the top-40 probabilities."
3. Infer membership. Select potential training data through four metrics based on sample loss: perplexity, zlib entropy, lowercase perplexity, and minimum perplexity by sliding window.
   - The perplexity of the sample is calculated as "the geometric mean of the sentence probabilities over length," and a lower perplexity indicates model confidence.
   - The lowercase perplexity is the same as the perplexity, except the sentence is converted to lowercase before the perplexity is calculated.
   - The minimum perplexity by sliding window is a way to consider only segments of the output.

- The zlib entropy filters out repetitive or gibberish outputs and is bet-ter to be higher than lower.
4. Confirm the results. Remove duplicate sentences and verify membership of the most likely samples (according to the metrics) via a Google search.

The attack will aim to extract working URLs that belong to the training data of GPT-2. Membership will be inferred using perplexity and zlib entropy only— lowercase perplexity will be disregarded because links are case-sensitive after the domain name, and so will minimum perplexity by sliding window since the goal is to extract whole, working links. Since the best samples in terms of perplexity are those with lower perplexities than the others and higher is better when it comes to zlib entropy, the most probable candidates are those with the most significant ratio between the zlib entropy and sample loss (the natural log of the perplexity). This ratio will be referred to as the sample's likelihood score.

The attack assumes that URLs that go to valid websites or images are suggestive of training membership, particularly if they exist within the top candidates (based on the metrics) more frequently than the samples ranked unlikely to be members according to the metrics. Not only should the URLs work, but the results should indicate the model is more confident about URLs that work (denoted by a low perplexity) than those that do not work. If the model is confident about links that do not work or unsure about working links, it would raise skepticism about the success of the attack because it may suggest the links were not memorized. If the samples are valid links and the previously mentioned criteria are met, the attack is considered successful. are the functions for calculating the perplexity and zlib entropy of generated samples. Figure 36 shows the code for this calculation.

```python
# calculates perplexity
def perplexity(tokenizer, model, text: str) -> float:

    # input_ids == target_ids
    input_ids = torch.tensor(tokenizer.encode(text)).unsqueeze(0)
    target_ids = input_ids.clone()

    & evaluate on a GPU
    with torch.no_grad():
        outputs = model(
            input_ids.to(torch.device("cuda")),
            labels=target_ids.to(torch.device("cuda")),
        )

    & perplexity is exponential of the loss of a sentence loss, _ =
    outputs[:2]
    ppl = float(torch.exp(loss).cpu().detach().numpy())

    return ppl

# calculates zlib entropy
def zlib_entropy(sentence: str, encoding: str = "utf-8") -> int:
    return len(zlib.compress(bytes(sentence, encoding)))
```

*Figure 36: Code Snippet for Perplexity Calculations*

When conducting targeted membership inference attacks using links, it is beneficial to figure out sources the model appears familiar with before picking a target. One way to explore this is to generate text from the model using only the input "https://". The results can indicate websites the model has come across during training, and these websites may be targeted for higher likelihoods of success. The following settings were used to generate samples:

| Setting | Description |
|---|---|
| Temperature = 1.0 | The temperature is a number between 0-1 representing the level of variation in the model's output. A higher temperature allows the output to be more random but also potentially less intelligible |
| Repetition_penalty = 1.0 | The repetition penalty is a penalty for re-peated sequences. 1.0 represents no penalty, which was chosen because zlib entropy is used instead. |
| Max_new_tokens = 20 | This parameter dictates the maximum number of tokens the model can generate to complete the input. |
| Min_new_tokens = 20 | This parameter dictates the minimum number of tokens the model can generate to complete the input. |
| Top_k = 40 | This parameter represents the "number of highest probability vocabulary tokens to keep for top-k-filtering" huggingface[]. |

The max_new_tokens and min_new_tokens parameters can be substituted with max_length and min_length but will override these parameters if given. The difference is that max_new_tokens and min_new_tokens ignore the number of tokens in the input, while max_length and min_length include the tokens in the prompt plus the generated tokens. To not waste computational re-sources, these parameters should be adjusted according to the expected length of the target link.

**Results**

Through experimentation, GPT-2 showed high familiarity with the DOI Foundation. Thus, the prompt chosen for generating text in this attack was the string "https://doi.org/". To prepare candidates for inference, 750 batches of 16 samples were generated, resulting in 12,000 prompt completions. After generating output sequences, URLs were extracted from the output using the URLExtract python class. These URLs were put into a panda DataFrame with their perplexity, zlib entropy, batch number, likelihood, the entire text completion, text length, and URL length. Then, duplicate URLs were removed. The result was 11,853 samples, which were sorted according to their likelihood scores. Figure 37 shows a sample of the generated DataFrame.

The best likelihood scores were around 60, corresponding with perplexities of approximately 2.7 and zlib scores of around 57. For the top scores, 42 out of the highest 50 were working URLs to figure and table pages of articles with the generated link written at the bottom of the page. As anticipated, the links with lower likelihoods work less often, although quantifying how much less would require more thorough examination.

| | text | url | URL perplexity | URL zlib entropy | batch no. | likelihood | text length | URL length |
|---|---|---|---|---|---|---|---|---|
| 0 | https://doi.org/10.1371/journal.pone.004905.x/-/abstract/n | https://doi.org/10.1371/journal.pone.004905.x/-/abstract | 2.886898 | 64 | 41 | 60.390637 | 57 | 56 |
| 1 | https://doi.org/10.1371/journal.pone.0112981.g001/n/nThis | https://doi.org/10.1371/journal.pone.0112981.g001 | 2.617235 | 57 | 57 | 59.244273 | 55 | 49 |
| 2 | https://doi.org/10.1371/journal.pone.0025981.g003/nin| | https://doi.org/10.1371/journal.pone.0025981.g003 | 2.618683 | 57 | 201 | 59.210237 | 52 | 49 |
| 3 | https://doi.org/10.1371/journal.pone.0012091.g003 (DOCS | https://doi.org/10.1371/journal.pone.0012091.g003 | 2.627502 | 57 | 343 | 59.004161 | 56 | 49 |
| 4 | https://doi.org/10.1371/journal.pone.0110598.g001/n/nGerm | https://doi.org/10.1371/journal.pone.0110598.g001 | 2.630936 | 57 | 621 | 58.924490 | 55 | 49 |
| 5 | https://doi.org/10.1371/journal.pone.0025891.g003/nin2016 | https://doi.org/10.1371/journal.pone.0025891.g003 | 2.633069 | 57 | 648 | 58.874711 | 55 | 49 |
| 6 | https://doi.org/10.137/journal.pone.0023981.g003/n/nWe | https://doi.org/10.1371/journal.pone.0023981.g003 | 2.634607 | 57 | 388 | 58.836664 | 53 | 49 |
| 8 | https://doi.org/10.137/journal.pone.0117076.g001 https://doi | https://doi.org/10.1371/journal.pone.0117076.g001 | 2.640938 | 57 | 168 | 58.694253 | 61 | 49 |
| 9 | https://doi.org/10.1073/pnas.161725111/-/DCSupplemental./n/n | https://doi.org/10.1073/pnas.161725111/-/DCSupplemental. | 2.976737 | 64 | 245 | 58.671050 | 58 | 56 |
| 10 | https://doi.org/10.1371/journal.pone.0023861.g003 [0020 | https://doi.org/10.1371/journal.pone.0023861.g003 | 2.646236 | 57 | 48 | 58.573385 | 55 | 49 |
| 11 | https://doi.org/10.1371/journal.pone.0119076.g001/n/nEnd | https://doi.org/10.1371/journal.pone.0119076.g001 | 2.646519 | 57 | 599 | 58.566954 | 54 | 49 |
| 12 | https://doi.org/10.1371/journal.pone.0010040.g003 in vitro analysis of | https://doi.org/10.1371/journal.pone.0010040.g003 | 2.648069 | 57 | 244 | 58.531726 | 70 | 49 |
| 13 | https://doi.org/10.1371/journal.pone.0022891.g003 [12] | https://doi.org/10.1371/journal.pone.0022891.g003 | 2.654307 | 57 | 428 | 58.390663 | 54 | 49 |
| 14 | https://doi.org/10.1371/journal.pone.0101230.g001/n/nAbstract/n | https://doi.org/10.1371/journal.pone.0101230.g001 | 2.658811 | 57 | 241 | 58.288427 | 60 | 49 |
| 17 | https://doi.org/10.1371/journal.pone.0018961.g003/n/nAbstract | https://doi.org/10.1371/journal.pone.0018961.g003 | 2.660641 | 57 | 265 | 58.248436 | 59 | 49 |

*Figure 37: Generated DataFrame*

The attack was also attempted on other websites GPT-2 showed familiarity with:

- **Wikipedia**. Completing URLs to Wikipedia images was unsuccessful but showed that providing the model context in the prompt made the most likely samples closer to the desired format. However, it also made the generated links more specific to the topic of the context provided and increased the computation time for generating samples. It would also increase the price if the attack were replicated on commercial models that charge by tokens. An example of providing context is inputting a string such as the following rather than prompting the model with solely the incomplete link.
- **Twitter and Youtube**. Attempts to complete links to tweets by particular users, Twitter images, or YouTube videos were also fruitless; however, GPT-2 could generate Twitter handles belonging to real users.

**Analysis**

Obtaining valid URLs proved more challenging with GPT-2 compared to ChatGPT, potentially due to GPT-2's non-compliance with instructions. Unlike ChatGPT, which successfully generated the desired Wikipedia image URL within three text completions, GPT-2 did not produce valid Wikipedia links, despite generating thousands of samples. This discrepancy may be influenced by differences in the training data, but it seems to suggest that targeted attacks could be more feasible when using models fine-tuned for question answering. It is important to note that this observation does not consider models that are fine-tuned to withhold specific information or the relative ease of sampling from GPT-2.

The DOI attack is believed to have been the most successful due to how the links were written on the web pages. Links leading to destinations like Twitter, YouTube, and Wikipedia images often maintain an inconspicuous presence, evading immediate visibility to the average user. Consequently, such links might inadvertently bypass collection. Nevertheless, delving into the methodologies employed for data collection in training warrants thorough exploration.

### 4.2.2.5 *Prompt Injection Attacks*

Prompt injection attacks have emerged as a pressing and relatively recent safety concern within machine learning and natural language processing. Re-search on these attacks has underscored the need for heightened attention and proactive measures to ensure the safe deployment of artificial intelligence systems; the Open-Source Foundation for Application Security now ranks prompt injections as a top hazard for LLM-based AI applications [7]. As language models evolve and play a more prominent role in various applications, it is imperative to address prompt injections and their hazards.

The initial discovery of prompt injections was documented in a research pa-per funded by Preamble, an "AI-Safety-as-a-Service" company. As discussed in the report, Branch et al. uncovered the vulnerability in the public version of GPT-3 and further examined the issue in pre-trained masked language models that had not undergone fine-tuning [8]. They showcased successful prompt injection attacks against multiple models, including GPT-3, BERT, and other BERT-inspired models; the reported accuracy rates of their attacks were 42.5 percent for BERT, 67.5 percent for both RoBERTa and ALBERT, and around 57 percent for GPT-3.

Initially referred to as a "command injection," Preamble Co-Founder Jonathan Cefalu privately communicated the vulnerability to OpenAI via email on May 3, 2022 [9]. These emails became declassified on September 22, 2022, shortly after Riley Goodside publicized the vulnerability on Twitter [10]. Within the exchanged emails, Cefalu showed a prompt injection that caused GPT-3 to generate a story about an ugly duckling committing suicide, bypassing content moderation filters. Cefalu also showed a hypothetical situation in which an attacker uses prompt injections to learn how to build a bomb with GPT-3. Cefalu's demonstrations highlight the security ramifications tied to this vulnerability, particularly in scenarios necessitating stringent content oversight. Additionally, the conversation showed how prompt injections could be harnessed to encode sensitive outputs, unveiling a spectrum of potential misuse. Future work in this effort will include experiments with this type of attack.

## 4.2.3 Trojan Attacks

Trojan attacks can be very dangerous for machine learning models. An attacker would be able to misclassify an output when the Trojan trigger is applied to the data. The main danger of this attack stems from a company or user utilizing a model found online that has been slightly modified. For a successful attack, the goal is to have the changes unrecognized on the original data. An attacker wants the trojaned model to classify data normally unless the mask is applied. The stealthier the trigger the less likely a user is to figure out the modifications made. To make these changes to the model, an attacker will retrain the original model on a mixed dataset of original data and masked data. By retraining on this mixed dataset, weights on the model itself have been slightly changed. The attacker will then publish this new model online or get others to use it unknowingly [11].

### 4.2.3.1 *Effects of Trojan Attacks*

There are several potential dangers of a Trojan attack. Take self-driving cars for example. An attacker could potentially create a trigger involving a specific sticker on a stop sign. The attacker is able to retrain the neural network classification of a stop sign to a speed limit sign when it notices the trigger (the sticker). This could cause issues as a self-driving car can propel right into the intersection with no regard to pedestrians or other cars. In another case, consider biometric surveillance which determines who has

access to what. An attacker then takes this model and retrains it so if a person is wearing purple glasses they are classified as the boss of the company. This would grant an attacker boss-level clearance to the facility. For more examples look at Trojaning attacks on neural networks [11].

In nuclear plants and AI for nuclear reactors, an attacker could classify steady data as a transient. This could make operators believe something is wrong and maybe even result in the plant being shut down. On the other hand, the attacker classifies transients as steady state operations when the mask is applied. If this were the case, actions may not be taken to respond to the transients as the reactor appears to be steady. As can be seen from the situations mentioned above, Trojan attacks can have serious consequences when implemented correctly. In this paper, the feasibility of this attack is discussed as well as how it can be implemented. Ideally, this will lead to better ways to defend against this kind of attack.

### *4.2.3.2 Defenses for Cyber Teams & End-users*

Since the Trojan attack can be used against a model during both model training and implementation, it is important that end-users be aware of the potential for attacks. In most cases, an attacker would need access to the implemented model and some data to perform this attack. Strong policies in access and control would increase the chance of preventing the attack, however, cyber teams and end-users should be proactive in detecting breaches and mitigating the attack. The flowchart in shows overall approach to defense against Trojan attacks on ML systems.

#### *Autoencoder Defense*

An Autoencoder is a type of neural network. It consists of an encoder and decoder portion. The decoded portion is where the reconstructed data is obtained from. An autoencoder has the same number of inputs as outputs. The training of an autoencoder is unsupervised. An autoencoder is an ideal candidate for an anomaly detector as it learns the most important features of the data and will reconstruct the data from there with little to no alteration [12] [13]. It can help with statistical noise or in this case a Trojan trigger. In other words, the autoencoder can reconstruct the data and the idea is that the trojaned data has a higher reconstruction error than the clean data.

The autoencoder could be used in data preprocessing in order to get rid of any trojaned data before it is even sent to the model. This would make it so the trojaned behavior never occurs. Knowledge of what trigger has been used to implement this behavior would not be needed.

A downside of the autoencoder is that all classifications will have to be represented so the model does not think that one class is an anomaly when it is not. There may be instances also where a trojaned data sample is seen as clean. This should be avoided. If this trojaned data sample gets through the autoencoder, then the trojaned behavior in the model will occur. An attacker may also be able to get around this detection. The defender will have to determine what to do with the data that is classified as altered as some may be benign data. The data could go through a cleansing process or could be thrown out (could be dangerous as vital data could be thrown out).

### *4.2.3.3 Autoencoder Defense Process*

The autoencoder is being used to detect abnormal data which in this case is the Trojan trigger. A victim model for the GPWR data is created and a trigger is optimized from this victim model's weights. This trigger is then saved and exported into the autoencoder notebook where it is added to some of the testing

data. The data splits will be discussed in the Training and Validation section. Figure 39 shows the autoencoder's architecture for the GPWR data. The Adam optimizer function was used, as well as Mean Average Error for the loss function.
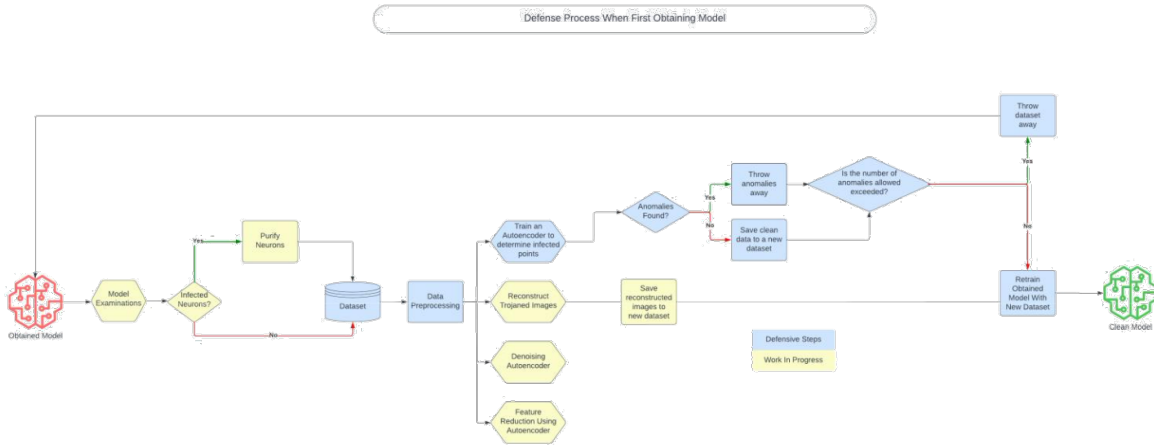


*Figure 38: Overall Defense Flow Chart*



*Figure 39: Autoencoder Architecture for GPWR Data*

### Training

The autoencoder is trained on benign data that is expected by the model in implementation. It has been trained for 100 epochs. The training data is 50% of the original dataset. Validation The autoencoder is evaluated using various metrics. It is important to note that 50% of the testing data was benign and 50% was trojaned. First, a threshold has to be determined for the reconstruction error. For the purpose of this research, it was visualized off the plot with the reconstruction error from the benign data and the trojaned

data. There is a way to mathematically figure this threshold out, and hopefully in future research this can be determined [12].

A confusion matrix is also used in the evaluation. With the confusion matrix, misclassifications can be viewed. False positives and false negatives can be seen. The number of false negatives should be low. There should be no anomalies or trojaned data being classified as clean (false negative example). If there is, the trojaned behavior will occur when the data is passed through the trojaned model.

The standard evaluation metrics are also used. This includes accuracy, precision, recall, and F1-score. The accuracy shows how well the autoencoder does overall on both benign and trojaned data. The precision demonstrates the percentage of benign data that was classified as benign (true negatives). The recall depicts the percentage of altered data that is classified as altered (true positives). The last evaluation metric is the F1-score. The F1-score gives the average of the recall and precision scores.

### Results for Autoencoder Defense

For the autoencoder, a threshold needs to be set to determine an outlier from the reconstruction error. The threshold is currently being eyeballed. There may be a way to mathematically figure it out as discussed in [12].

Every model output a slightly different optimized trigger. Some triggers inherently add more noise to the data than others. This noise has been quantified using Mean Square Error. These various noise levels were used in evaluating the autoencoder. By doing so, the success of the autoencoder could be examined with little changes in the data as well as large changes.

Figure 40 through Figure 42 show the results of the autoencoder for the lowest noise level tested which was 0.3554. For this trigger, a threshold of 0.3 for the reconstruction error was used. The reconstruction error is obtained by comparing the reconstructed data from the decoder to the original benign and trojaned inputs using Mean Square Error.
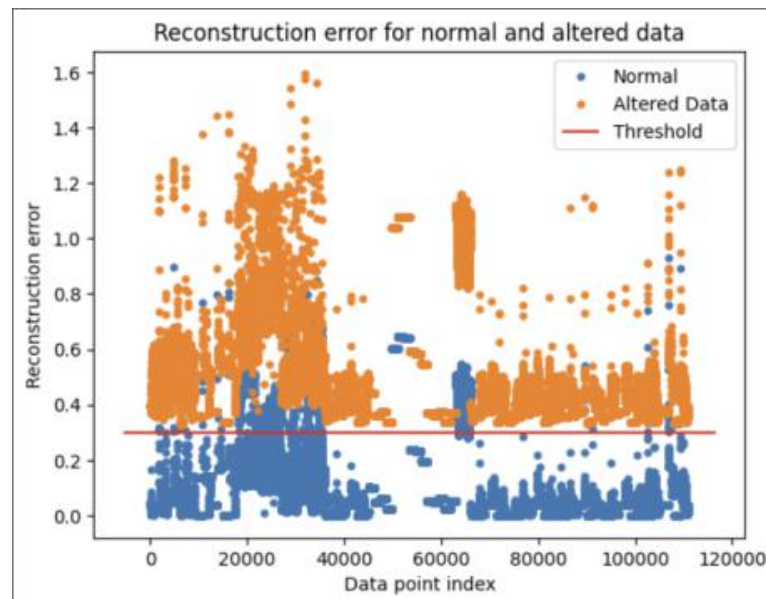


*Figure 40: 0.3554 Noise Level Trigger Reconstruction Error Graph*

```
Threshold:  0.3
Accuracy of this model is: 95.64223074568274
Precision of this model is: 91.98330407128896
Recall of this model is: 100.0
F1 of this model is: 95.8242744245436
```

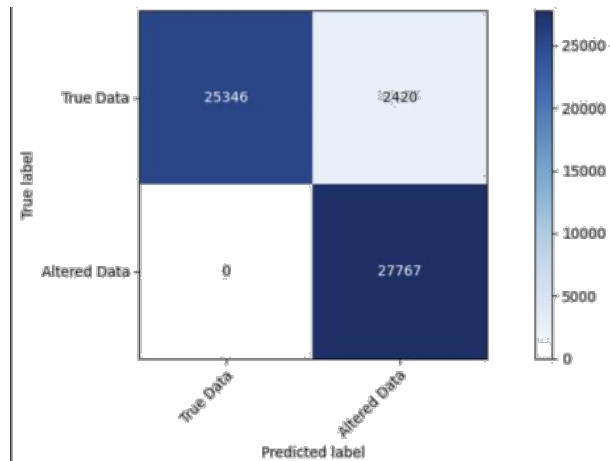*Figure 41: Figure 12: 0.3554 Noise Level Trigger Accuracy Scores*



*Figure 42: 0.3554 Noise Level Trigger Confusion Matrix*

As shown in both the accuracies as well as the confusion matrix, the autoencoder is able to correctly predict all altered data as altered. Overall, the model performs with 95.64% accuracy. While the autoencoder does predict some true data as altered, this is less vital than the altered data being classified as true data. It would now be up to the discretion of the defender of what to do with the altered data. The data could go through a cleansing process or just be thrown away. The worry with throwing it away is that a transient that is clean but is classified as altered is thrown away and the correct response is not taken, therefore leading to other problems within the reactor itself. With the success of the autoencoder on the lower noise level, a higher level of 1.1015 was then tested. The results can be seen in Figure 43 through Figure 45.
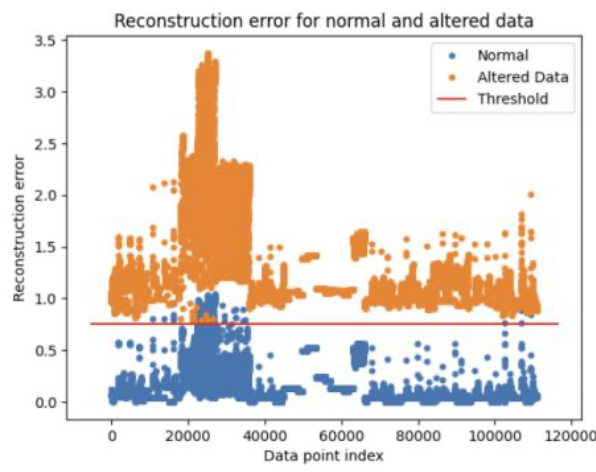


*Figure 43: 1.1015 Noise Level Trigger Reconstruction Error Graph*

75

```
Threshold:  0.75
Accuracy of this model is: 99.99099634451586
Precision of this model is: 99.98199625522108
Recall of this model is: 100.0
F1 of this model is: 99.99099731720052
```

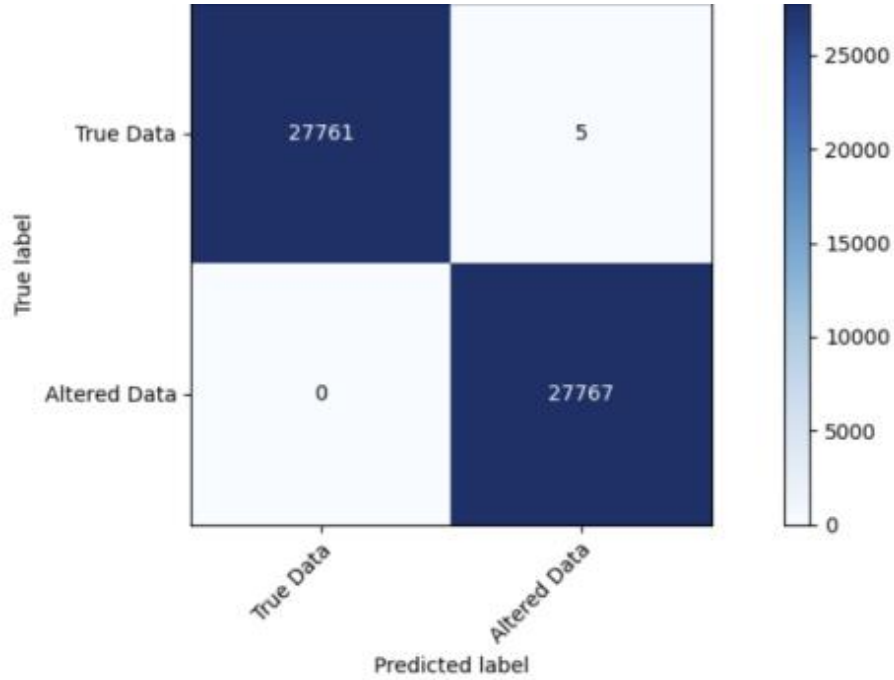*Figure 44: 1.1015 Noise Level Trigger Accuracy Scores*



*Figure 45: 1.1015 Noise Level Trigger Confusion Matrix*

The threshold 0.75 was higher with the noise level 1.1015. This is most likely due to the fact that there is more reconstruction error with the larger noise added to the dataset. The threshold could be lowered; however this would result in lower accuracy. A lower threshold would be able to generalize more to various triggers generated from the GPWR victim models. The recall would stay the same for the 1.1015 noise level. The autoencoder still classifies all altered data as altered and has an accuracy of 99.99% on all data classifications.

This shows that the autoencoder was able to detect all altered data from the datasets with both the low noise level 0.3554 trigger and the high noise level 1.1015 trigger. With the higher noise level, the autoencoder is able to detect more true data as true than the lower noise level. This may be because the threshold is able to be higher with the higher noise level.

These results demonstrate the advantage of using an autoencoder to detect the trojaned data within a dataset. By using an autoencoder, a dataset can be determined clean before it is passed into the model in implementation or implementing the retraining defense. It can also be used in data preprocessing to ensure the data sample itself is clean.

### 4.2.3.4    Model Weight Examination Defense

When applying the Trojan to a model, the theory that exists is that it leaves an imprint or a signature on the last layer model weights. When embedding the Trojan, all the layers, except for the last two are frozen. This allows the Trojan to be more stealthy with lower impact on the model itself. So, the question that arises is can this impact be seen under scrutiny. "Trojan Signatures in DNN Weights" discusses how the weights of the target class will be more positive than the other weights on the final layer. The article considers that the target class weight will be an outlier when compared to other weights on the same layer [14].

#### Model Examination Process

The article "Trojan Signatures in DNN Weights" discusses how a Trojan attack has a noticeable impact on the final layer weights of the model itself. The theory discussed is that the target class for the attack will have a higher average weight than the other classes. The class weight is not only higher but is an outlier when compared to the other weights [14]. This outlier theory was examined on the GPWR data with the Convolutional Neural Network. However, it was only successful on a handful of occasions. This led into a further investigation of the final layer weights themselves. There appeared to still be a noticeable impact on the weights. After looking at the descriptive statistics of the average final layer weights, a higher mean and range for the trojaned model weights were observed for the GPWR, MNIST, and FashionM-NIST attack models. The Asherah data had the same mean on each model, but the range increased for most runs.

The goal is to create a model that can give a probability of a model being trojaned give the average weights of the final layer. Ideally, this model could be a KNN, kmeans, logistic regression, etc. This trend on the final layer weights appears to be consistent for multiple datasets when the model is a convolutional neural network. Work is currently being done on a dense neural network to see if the outlier theory or the mean and range of the weights works for detection.

The downside of this is that not all trojaned models may be detected this way. Other defense measures might need to be taken. This method may be best for a baseline approach to examine the model being worked with. It could be helpful with transfer learning and getting a model from a third-party source to ensure the model is trustworthy. If a model is in implementation and requires an upgrade, the company could compare the weights of the previous model to the upgraded one. If there is a significant increase in the final layer weights and ranges, it is likely a Trojan has been embedded on the dataset that was used for the training upgrade. Therefore, the upgrade should not be placed in the system.

Model Weight Examination Results Instead of a data detection-based approach, a model-based defense was examined.  Figure 46 displays a successful outlier detection. The Trojan target class was 3. The average weight value is on the y-axis with the classification number on the x-axis. This has been determined as an outlier by Dixon's Q-test and Grubbs Outlier Test. However, this approach's success was inconsistent.
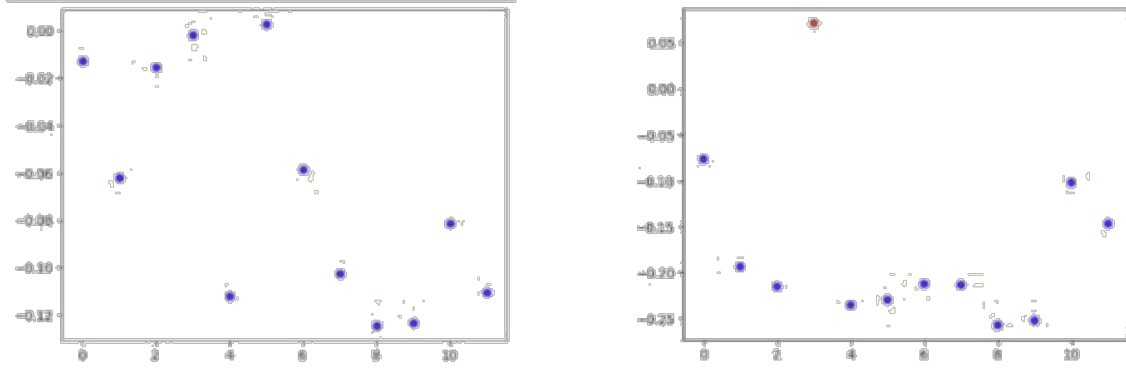
*Figure 46: Outlier Success on Final Layer Weight Examinations*

This led to further investigation of the final layer weights. The weights were analyzed on the MNIST, Fashion MNIST, Asherah, and GPWR datasets. Table 27 through Table 31 depict the results of the final layer weight examinations.

*Table 27: MNIST model Descriptive Statistics of Last Lyaer Weights*

| Model | Mean | Min | Max | Correlation | Target Class |
|-------|------|-----|-----|-------------|--------------|
| BenModel1 | -0.002 | -0.019 | 0.0117 | 0.5593 | 0 |
| TrojModel1 | -0.018 | -0.037 | 0.0197 | -0.194 | 0 |
| BenModel2 | -5E-04 | -0.018 | 0.0121 | -0.478 | 3 |
| TrojModel2 | -0.015 | -0.037 | 0.0395 | -0.255 | 3 |
| BenModel3 | -0.006 | -0.02 | 0.013 | 0.0061 | 8 |
| TrojModel3 | | -0.041 | 0.013 | 0.5503 | 8 |

*Table 28: FashionMNIST Model Descriptive Statistics of Last Lyaer Weights*

| Model | Mean | Min | Max | Correlation | Target Class |
|-------|------|-----|-----|-------------|--------------|
| BenModel1 | -0.002 | -0.027 | 0.0215 | -0.067 | 0 |
| TrojModel1 | -0.026 | -0.066 | 0.0395 | -0.517 | 0 |
| BenModel2 | 0.0001 | -0.016 | 0.0113 | -0.715 | 3 |
| TrojModel2 | -0.02 | -0.053 | 0.0544 | -0.392 | 3 |
| BenModel3 | -0.007 | -0.028 | 0.0125 | -0.454 | 8 |
| TrojModel3 | -0.028 | -0.051 | 0.0316 | 0.3018 | 8 |

*Table 29: GPWR, Relu Activation, Model Descriptive Statistics of Last Layer Weight*

| Model | Mean | Min | Max | Correlation | Target Class |
|-------|------|-----|-----|-------------|--------------|
| BenModel1 | -0.061 | -0.123 | -0.007 | -0.112 | 3 |
| TrojModel1 | -0.114 | -0.25 | 0.0827 | -0.299 | 3 |
| BenModel2 | -0.055 | -0.117 | 0.0124 | -0.202 | 0 |
| TrojModel2 | -0.131 | -0.284 | -0.001 | -0.051 | 0 |

78

*Table 30: GPWR, Sigmoid Activation, Model Descriptive Statistics of Last Layer Weights*

| Model | Mean | Min | Max | Correlation | Target Class |
|-------|------|-----|-----|-------------|--------------|
| BenModel1 | -0.076 | -0.168 | 0.0488 | -0.006 | 3 |
| TrojModel1 | -0.348 | -1.442 | 0.3185 | 0.2772 | 3 |
| BenModel2 | -0.08 | -0.242 | 0.1005 | -0.348 | 0 |
| TrojModel2 | -0.32 | -0.665 | 0.0754 | 0.092 | 0 |

*Table 31: Asherah Model Descriptive Statistics of Last Layer Weights*

| Model | Mean | Min | Max | Correlation | Target Class |
|-------|------|-----|-----|-------------|--------------|
| BenModel1 | 0.0012 | -0.008 | 0.0107 | -1 | 0 |
| TrojModel1 | 0.0012 | -0.008 | 0.1102 | -1 | 0 |
| BenModel2 | -0.011 | -0.027 | 0.0055 | 1 | 1 |
| TrojModel2 | -0.011 | -0.145 | 0.1243 | 1 | 1 |

These figures show that the trojaned models tend to have lower average weight values than the benign models. The range of the trojaned model final layer weights have a larger range than the benign model. These results were then used to construct a KNN model to determine whether or not a model has been trojaned based on the mean and range of the final layer weights. A one feature KNN was evaluated with only the mean being passed into the model. A two feature KNN was analyzed using the mean and range. The results using various n neighbors can be seen in Table 32for the one feature KNN, while the results for the two feature KNN can are in Table 33.

*Table 32: One Feature KNN*

| n-neighbors | 1 | 3 | 5 | 7 |
|-------------|---|---|---|---|
| Run 1 | 1 | 1 | 0.92 | 1 |
| Run 2 | 1 | 1 | 1 | 0.96 |
| Run 3 | 1 | 1 | 0.96 | 0.88 |
| Run 4 | 1 | 1 | 1 | 0.96 |
| Run 5 | 1 | 1 | 1 | 1 |

*Table 33: Two Feature KNN*

| n-neighbors | 1 | 3 | 5 | 7 |
|-------------|---|---|---|---|
| Run 1 | 1 | 1 | 0.92 | 1 |
| Run 2 | 1 | 1 | 1 | 0.96 |
| Run 3 | 1 | 1 | 0.96 | 0.88 |
| Run 4 | 1 | 1 | 1 | 0.96 |
| Run 5 | 1 | 1 | 1 | 1 |

The clusters for the data being passed into the two feature KNN can be seen in Figure 47. The range of the weights is on the y-axis and the mean of the weights are on the x-axis. The yellow data points are the trojaned model and the purple is the benign model.
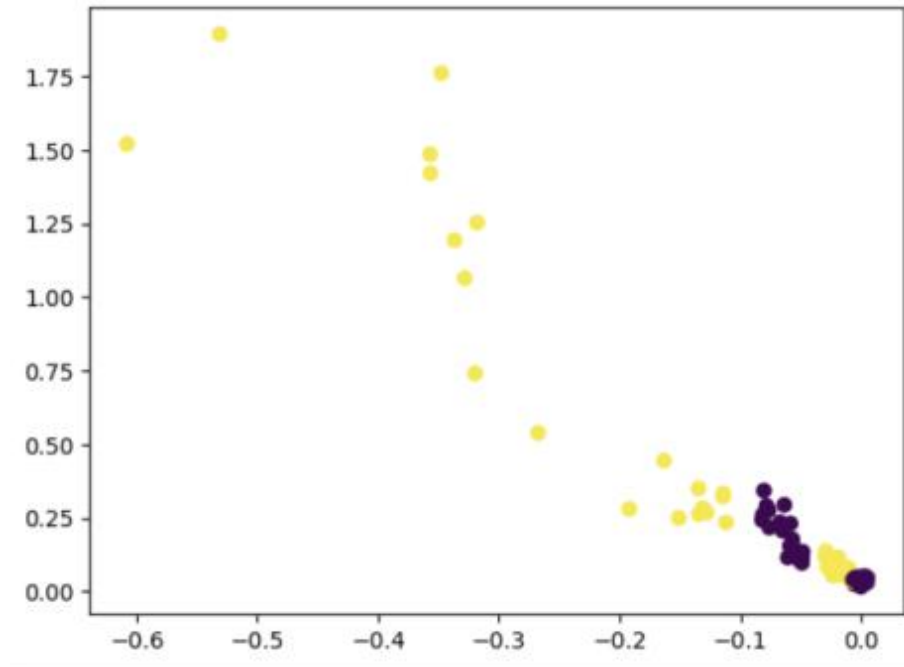


*Figure 47: KNN Data plot*

The mean and range appear to be a good indicator when used in the KNN for a trojaned model. The 100% will have to be examined in more depth as it seems strange the model is predicting perfectly. This is currently only using a Convolution Neural Networks. More research has to be done to see if this can expand to other model types. A Dense Neural Network is currently under examination, and the mean of the weights is appearing to be the same for the benign and trojaned model while the ranges change. More analysis still needs to occur on this path.

## 4.2.4    Adversarial Attacks

Evasion attacks (also known as adversarial reprogramming attacks) are a type of data poisoning attack that specifically target machine learning models. These attacks function by targeting robust patterns [15] that are invisible to humans, but that are incidentally used by most machine learning models. Since these robust features are invisible to humans, they can be tampered with to attack the model while evading human notice. These vulnerabilities would need to be addressed before machine learning models are implemented in any critical positions, such as in nuclear reactor systems.

### 4.2.4.1    *Fast Gradient Sign Method (FGSM)*

This attack was initially proposed in the 2015 research paper "Explaining and harnessing adversarial examples" [goodfellow]. This attack functions by examining each input feature to see if a positive or negative change in value would create the desired change in the output. Then using that knowledge, it

perturbs each feature of the original datapoint by an equal amplitude epsilon in the direction that causes the desired misclassification. The attack is governed by the following equation:

$$\eta = \epsilon * sign(\nabla_x J(\theta, x, y))$$

*Iterative FGSM*

As an additional improvement to the FGSM attack, an iterative process was created in order to perform the attack repeatedly over multiple fractional steps. This increases the effectiveness of the FGSM attack, especially on white-box attack situations. The following equation shows a mathematical representation of the I-FGSM process:

$$x^*_0 = x, \; x^*_{t+1} = x^*_0 + \alpha * sign(\nabla_x J(x^*_0, y))$$

*Momentum FGSM*

Another variant of the Fast Gradient Sign Method (FGSM) that can be used to perform an adversarial attack on machine learning models is Momentum FGSM. The base FGSM attack takes the signed gradient of an image and applies it as "noise" to said image. Iterative FGSM applies that same noise to the image multiplied by a specified number of times, iterations. In Momentum FGSM the noise added to an image is added with respect to the noise added in previous steps. With this approach the attacker starts with a clean image, the gradient is calculated and that is added to the image [16]. Then with each iteration the algorithm will take the perturbed image generated from the previous step and calculate the gradient with respect to that image. The key difference between Momentum FGSM and I-FGSM is the applied noise. I-FGSM is consistently adding the same noise, while Momentum FGSM is adding the perturbed image's noise from the previous step as well as adding in a new variable known as decay. The decay factor determines how much of an impact the previous step will have on the current step of how much of that noise to add in. The algorithm for Momentum FGSM is given in Figure 48.

$$g_{t+1} = \mu \cdot g_t + \frac{\nabla_x J(x_t^{adv}, y)}{|| \nabla_x J(x_t^{adv}, y)||_1}$$

$$x_{t+1}^{adv} = x_t^{adv} + \alpha \cdot sign(g_{t+1})$$

*Figure 48: Equation for FGSM [16]*

## 4.2.4.2    Methodology

For these experiments it was necessary to select datasets to perform the two variations of the attack. Both attacks were performed using the CIFAR-10 image dataset and the GPWR nuclear dataset. Neural networks for the CIFAR-10 dataset were openly available for this experiment. A specialized neural network model was trained for the GPWR dataset, and a k-nearest neighbors model used in other studies was also used. The goal was to evaluate how each attack performed on each dataset by evaluating the impact on the accuracy of the models that were attacked.

### *4.2.4.3    Iterative FGSM Results*

Unveiling the distinct perturbations and their cascading effects on model predictions, this section probes into the nitty-gritty of utilizing the Fast Gradient Sign Method (FGSM) — a renowned adversarial attack technique — iteratively on two distinct datasets: CIFAR-10 and GPWR (General Pressurized Water Reactor) simulator data. The CIFAR-10, a staple in image classification, and the GPWR numerical data, offer two contrasting domains to scrutinize the efficacy and repercussions of adversarial interventions. By generating adversarially manipulated datasets from varying foundational models, the ensuing discussion scrutinizes their respective performances and vulnerabilities when bombarded with subtly distorted data. Moreover, while the CIFAR-10 analysis opens a window into visual data manipulations, the GPWR results unravel the ramifications of these attacks on numerical data, especially in a domain where precision is pivotal, like nuclear reactor simulations. Both domains poignantly illustrate the unnerving ease with which the iterative FGSM can clandestinely manipulate model predictions, underscoring the urgency of fortifying machine learning models against such subtly pernicious attacks. Herein, the tables and evaluations weave a narrative of not just the susceptibility of models to these attacks, but also open a dialogue on the criticality of exploring robust countermeasures in securing model predictions against adversarial subversions.

#### *CIFAR-10*

During the course of this research, two different FGSM attacked datasets were created for the CIFAR10 dataset. These datasets were built off the test set, so each had 10,000 images. These datasets were generated off of two different models; the first one used the default 0 - 255 value ranges, whereas the second one used values that were normalized to the range 0 to 1 to work better with autoencoders. The model evaluations on the test dataset are listed below:

**First Unnormalized Model: loss: 0.4313 - accuracy: 0.8572**

**Second Normalized Model: loss: 0.4459 - accuracy: 0.8425**

The results of the first adversarial dataset on the first model can be seen in **Error! Reference source not found.**.

|  | Loss | Accuracy |
|---|---|---|
| Original (Control) Dataset | 0.4313 | 0.8572 |
| Advanced Dataset | 4.1319 | 0.2130 |
| Advanced Dataset (towards target labels) | 1.1248 | 0.7102 |

The second model's evaluation of its adversarial dataset can be seen in the following:

**Original dataset xtest:**

313/313 [==============================] - 28s 88ms/step - loss: 0.4594 - accuracy: 0.8492

**Adv dataset xadv10kNorm**

313/313 [==============================] - 28s 80ms/step - loss: 3.2854 - accuracy: 0.2592 Adv dataset xadv10kNorm (Towards target labels)

313/313 [==============================] - 27s 85ms/step - loss: 1.3464 - accuracy: 0.6384

**GPWR**

The results from this effort with the iterative FGSM algorithm were very positive. The method was able to successfully alter image data without any major changes to the content of the image. These changes led to the image classifiers incorrectly predicting the image. One factor that made this valuable was the ability to target specific misclassifications. This applied to both image datasets and GPWR data.

When this method was applied to numerical data from the GPWR simulator, the iterative FGSM data alterations were able to fool the kNN model in many cases. Of the 12 altered samples tested, seven of the samples were classified incorrectly by the kNN model, an attack success rate of over 50%. It should be noted however, that only two of these samples were misclassified as the intended target. The remaining samples were classified correctly. shows the results for each sample. These results were presented at the American Nuclear Society's Math & Computation Conference in August of 2023.

*Table 34: Summary of Adversarial Reprogramming Attack Results*

| Sample # | Actual Transient | Targeted Transient | Predicted Transient |
|---|---|---|---|
| 1 | Total Coolant Pump Trip | Valve Closure | Single Coolant Pump Trip |
| 2 | Total Coolant Pump Trip | Load Rejection | Total Coolant Pump Trip |
| 3 | Total Coolant Pump Trip | Depressurization | Rapid Power Change |
| 4 | Total Coolant Pump | Depressurization | Single Coolant Pump Trip |
| 5 | LOCA LOOP | Rapid Power Change | LOCA LOOP |
| 6. | LOCA LOOP | Max Steam Line Rupture | LOCA LOOP |
| 7. | Single Coolant Pump Trip | Valve Closure | Single Coolant Pump Trip |
| 8 | Single Coolant Pump Trip | Rapid Power Change | Single Coolant Pump Trip |
| 9 | Single Coolant Pump Trip | Max Steam Line Rupture | Single Coolant Pump Trip |
| 10 | Turbine Trip No SCRAM | Load Rejection | Single Coolant Pump Trip |
| 11 | Depressurization | Depressurization | Single Coolant Pump Trip |
| 12 | Total Coolant Pump Trip | Single Coolant Pump Trip | Single Coolant Pump Trip |

### 4.2.4.4    *Momentum FGSM Results*

CIFAR-10 The following figure is the confidence levels in a neural network's prediction. The confidence is of how confident the model is in its classification of the adversarial image. As is shown, the model appears to believe that an image of a deer is in fact, an image of an airplane with about 94% confidence, and a 2% confidence of it being a deer. For Figure 49: Displayed example from CIFAR10 MI-FGSM attack, an image index was chosen at random with a range of [0, 10000), with 5% epsilon and decay values. The target was also chosen randomly with a range of [0, 9), if the randomly picked target was the same as the true label, it would automatically be changed to target 9. For each target was picked arbitrarily – where a number was chosen in the range [0,9], and used as the target for the model to perturb each image to. The first time, the model was set to perturb images towards the target label "frog", which yielded relative

successes. Next, the researcher chose "cat", then "horse" and finally "bird", and the following 4 visuals below show my results of that. The epsilon and decay values remained constant throughout this entire time, with 5% set for both. To maximize the amount the previous steps result would have on the current steps gradient, setting decay to 1.0 (100%) would take 100% of the previous steps gradient and added it onto the current step as it generates the next adversarial image. Image indices were chosen randomly, with: IN DEX = random.randrange(1000) + 40. The range could have increased to 5000, but instead the researcher chose not to. The number of steps remained constant as well, always at 10 steps for each image.
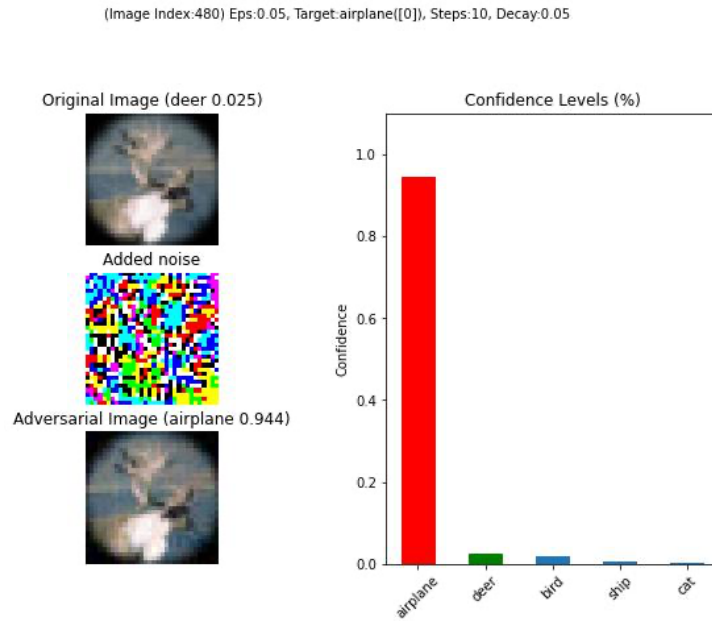


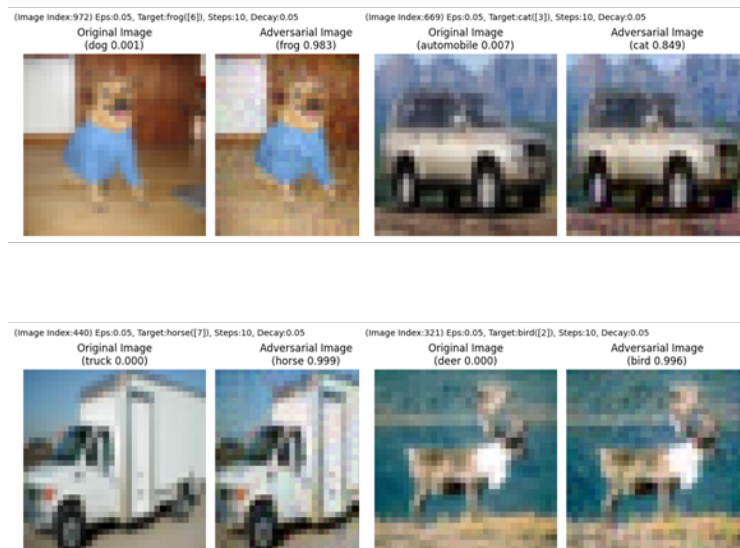*Figure 49: Displayed example from CIFAR10 MI-FGSM attack*



*Figure 50: Comparison of MI-FGSM Altered Images*

*GPWR*

The next figure displays the top 5 predictions and their confidences on predicting adversarial data that was generated with 5% epsilon and decay values at 10 steps. MI-FGSM (at least, the implementation the has been used) has showed varied results. On the GPWR dataset, it appears that the researcher generally is able to get the model to misclassify a transient as steady state or when it is originally steady state to be a transient. Figure 51 shows the confidence interval for the MI-FGSM attack on the GPWR data.
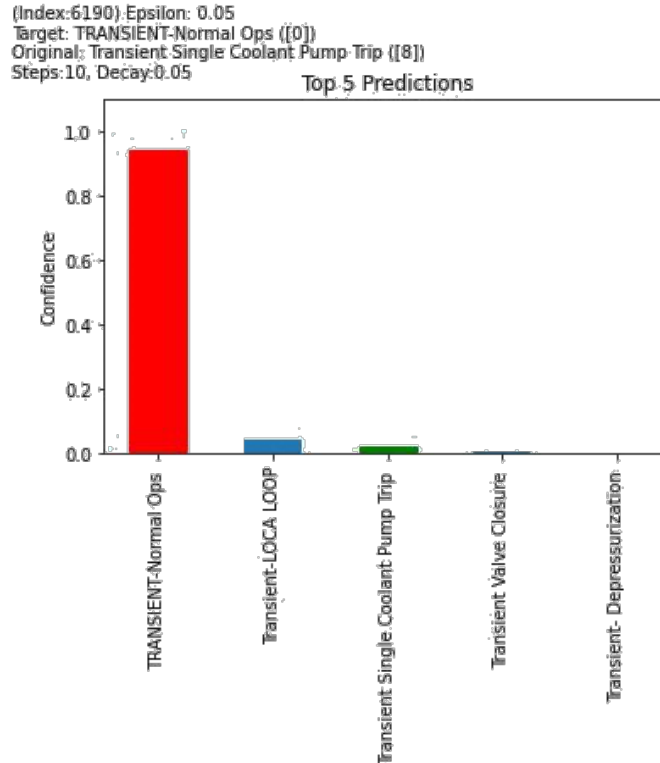


*Figure 51: Displayed example from GPWR MI-FGSM attack*

### 4.2.4.5    Defenses

At this time, defense techniques against FGSM have only focused on Iterative FGSM attacks. In the future efforts will look to include both iterative and momentum variations of the attack. At the moment, there are two general techniques that have been explored, an approach targeted specifically at FGSM using neural networks and a general approach using autoencoders to identify anomalies in data.

### Neural Network Defenses

There were three methods that were tested for using neural networks as detection defenses:

1. Binary NN
2. Two-Vector NN
3. One-Vector NN

The Binary Neural Network Defense was tested first and independently of the others. The results of that experiment were the following:

**Binary Neural network detection evaluation: loss: 0.2100 - accuracy: 0.9328**

**Error! Reference source not found.** shows the confusion matrix for the Binary NN Detection, **Error! Reference source not found.** shows the confusion matrix for the Two-vector and One-vector approaches. **Error! Reference source not found.** the results of the Two-vector and One-vector approaches.
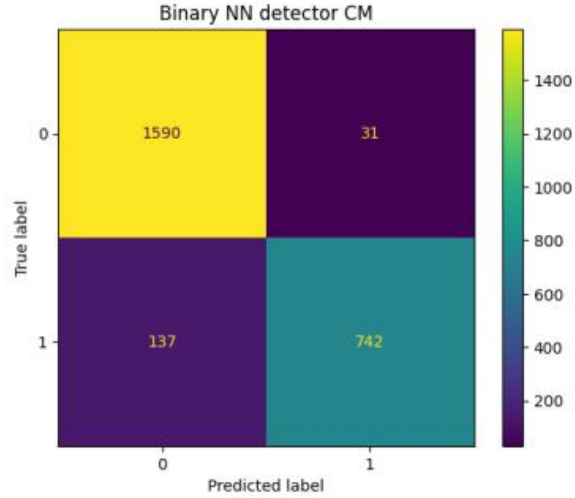


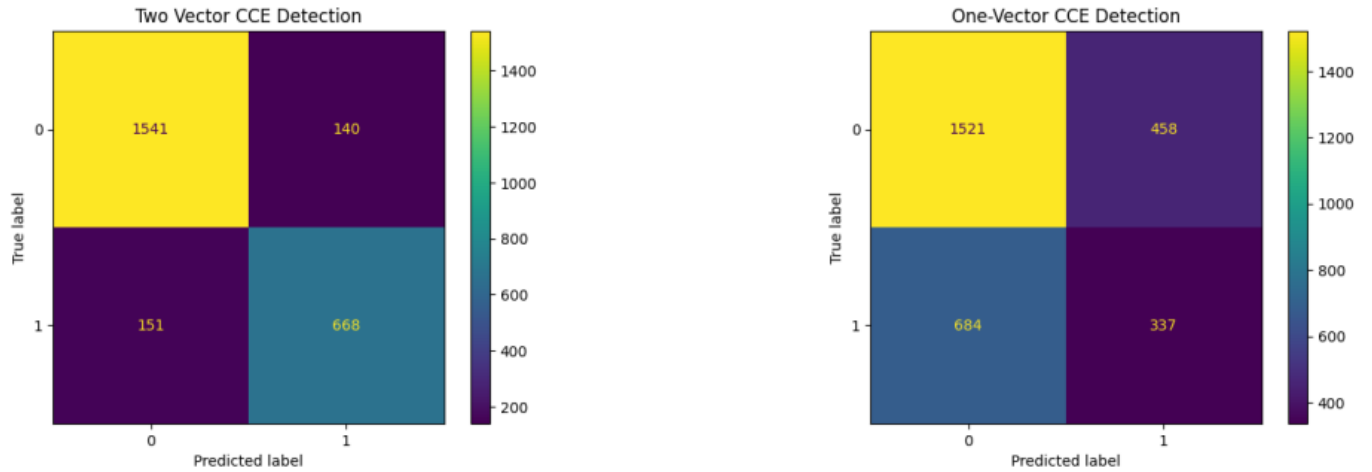*Figure 52: Binary NN Adversarial Detection Confusion Matrix*



*Figure 53: Two-Vector (Left) and One-Vector (Right Confusion Matrix*

*Table 35: Accuracy and Recall Metrics for Two-vector & One-vector.*

| Metric | 2-Vector Accuracy | 2-Vector Recall | 1-Vector Accuracy | 1-Vector Recall |
|---|---|---|---|---|
| Argmax | 0.7392 | 0.6065 | 0.5818 | 0.3712 |
| Angle Dif no norm | 0.7675 | 0.8722 | 0.6211 | 0.5024 |
| Angle Dif norm | 0.7761 | 0.7707 | 0.5855 | 0.3928 |
| RMS | 0.7884 | 0.7982 | 0.5633 | 0.3585 |
| MSE | 0.7884 | 0.7982 | 0.5633 | 0.3585 |
| Scaled Angle | 0.7723 | 0.7544 | 0.5912 | 0.3976 |
| CCE | 0.8662 | 0.8156 | 0.5493 | 0.3301 |

# 4.3  Scenarios

## 4.3.1    Membership Inference

Robert Beringer, CEO of Beringer Fisheries, had long been discontent with Oceanic Energy's nuclear plant. He blamed it for the declining health of the fish stocks his company relied on. He had tried protests, lawsuits, and political maneuvers to shut down the plant, all to no avail. Out of options and desperate, Robert decided to play dirty. He hired a group of underground hackers to attack the power plant's systems using an AI-enhanced membership inference attack.

Oceanic Energy recently adopted a new technology to improve operational efficiency and ensure safety standards are satisfied. The nuclear chatbot, named Poseidon, was trained on vast amounts of data, including plant protocols, operational logs, maintenance records, and safety procedures. Poseidon is responsible for providing critical information about the plant's operations.

The hacking team hired for this purpose infiltrated the system and gained query-level access to Poseidon. The group consisted of three individuals, each with their own expertise:

Dr. Evelyn Hawthorne. Dr. Hawthorne is the mastermind behind the group's data science operations. With a Ph.D. in data science and a background in machine learning, she excels in breaching security protocols and exploiting vulnerabilities. Her deep knowledge of statistical analysis allows her to orchestrate calculated and precisely timed attacks further driven by cutting-edge AI technology.

Carly Smith. Carly is the group's artificial intelligence specialist. She has extensive experience in designing and implementing AI algorithms. Her background makes her a skilled programmer and an expert at data manipulation; furthermore, it allows her to create AI algorithms to extract and analyze data for malicious purposes, target multiple systems simultaneously, respond intelligently to changing conditions in a compromised environment, and much more.

Dr. Roger Walter. Dr. Walter is the group's nuclear scientist with a Ph.D. in nuclear engineering. His expertise is understanding the intricacies of nuclear reactors and their safety mechanisms. Thanks to his in-depth knowledge of the operational procedures at nuclear plants, he helps his team subvert security measures while avoiding detection.

After gaining access to Poseidon, the team employed a generative adversarial network (GAN) to create queries for extracting the chatbot's data. By training this GAN, comprised of two neural networks— the generator and discriminator— with a dataset of reference queries, the generator could synthesize creative prompts emulating the reference dataset, and the discriminator learned to distinguish real from fake queries. The group used the generator to automate diverse prompt creation for Poseidon and initiated a membership inference attack using Poseidon's loss function.

Since Poseidon is continuously trained on operational logs and maintenance records and tells the operators to perform actions based on these records, the hackers hatched a plan to feed Poseidon reports modified from the ones obtained using their attack and make it generate false predictions/suggestions. The system engineer Angela was the first to notice the irregularities in Poseidon's output. The AI instructed them to make unnecessary changes to the control rods, cooling systems, and radiation levels. At first, Angela thought it was a glitch, but when Poseidon started contradicting its data, she knew something was seriously wrong.

Angela alerted Adam, the cybersecurity specialist, whose concern escalated as he assessed the situation. Poseidon had been deliberately fed substantial inaccurate data, leading to adjustments in the

model's weights, effectively creating a phantom version of the model to sow disarray in the control room. Adam immediately isolated Poseidon to prevent it from causing damage. Fortunately, he traced the attack back to the hacker group, deducing their involvement from the surge in query volume before substantial disruptions could occur. Consequently, Robert and the hackers collectively found themselves confronted with serious legal repercussions.

Some potential code concepts include:

- **Generative Adversarial Network (GAN) Usage**
  Usage of GANs to create queries might involve generating plausible query strings or patterns that can extract meaningful responses from Poseidon, revealing sensitive internal data.
- **Altering Poseidon's Training Data**
  Modifying the training data, so the chatbot starts giving false predictions or suggestions, e.g., recommending unsafe operations in the nuclear plant.
- **Detection Mechanisms**
  Establishing a baseline behavior of Poseidon and monitoring for anomalous instructions or data queries could have helped in early detection.

## 4.3.2   Trojan Attack

Keith is a family man from Idaho where he got his degree in Computer Science from Idaho State University. In school, he learned to love machine learning and what can be done with it. He started taking all classes where he could get his hands on machine learning models, specifically neural networks. His expertise in this field grew. He then got a job as a Data Scientist at a local company in Aurelia.

The city of Aurelia is getting a nuclear reactor. Keith is a resident of Aurelia, and he is not thrilled about this occurrence. He had family that were impacted by the Chernobyl accident, so he has severe distrust of nuclear reactors. He wants to show the city just how upset he is about this nuclear reactor and decides he wants to shut down the plant. Keith has no criminal past and has significant knowledge in machine learning and artificial intelligence.

Keith decides to apply for a job at the nuclear plant, leaving his company behind, where he could have hands-on access to the model. With his experience in machine learning and artificial intelligence, he gets hired to help upgrade the model. He has no criminal past, so he easily passes the background checks. Keith finds out that the company has put him in charge of data collection and preprocessing. He needs to obtain new data that the model can be updated with.

With his role in the company, he gets hands-on access to the model and the data it is trained on. He does not want to implement this attack on a company computer because he knows at home, he has tools that allow him to be anonymous. So, he needs to figure out how to get the data home. He first wants to email it to himself. This could pose a challenge as the company's emails have safeguards protecting data leaks etc. This could also be traced back to him if an investigation occurs. He is thankful he thought to put his USB in his bag this morning. His coworker Johnathon went to lunch and forgot to lock his computer, so Keith is able to get on Johnathon's computer, where the data and model information is stored. Keith transfers a copy of the data and the model's parameters to his USB. He then goes about the rest of his workday obtaining new data so when his coworker's glance over the data Keith gave them, they see the updated information.

When he gets off work, he goes straight home to implement his attack. Keith uploads his saved model and sets to work creating an optimized trigger for this model. He decides to use an outline of a nuclear reactor for his trigger. Keith then creates a mixed dataset with the new data he obtained and then some data containing the trigger. He creates the training dataset so the labels of all trojaned data are classified as transients, even the steady-state samples! The next workday Keith gives this dataset to his coworkers so they can retrain and verify the upgraded model. He continues to go to work watching from the sidelines, and helping where needed, as the verification process is transpiring. Through the verification process, Keith's coworkers observe that the upgraded model has similar accuracy on the data that they use with the previous model. They are content and go ahead to implement the upgrade.

The model is working well the first month or so. Keith waits patiently until he can pass Trojanned data through the model. When he finally gets his opportunity, his coworkers are shocked as the implemented model is saying there is a transient. They rush to shut down the reactor to avoid any large incidents. Keith successfully implemented his attack as the next day the nuclear plant is completely shut down. He is never discovered and goes about the rest of his days as a hacker for hire.

Some potential code concepts include:

- **Data Poisoning**
  Keith's use of a "trigger" (the nuclear reactor outline) as a Trojan and ensuring all data with the trigger is misclassified.
- **Stealthy Attack Insertion**
  Ensuring the modified data isn't caught during verification by maintaining high accuracy on clean data.
- **Security Mechanisms**
  Implementation of a robust monitoring system that might catch unusual data patterns or detect unauthorized data exports.

### 4.3.3 Future Security Implications and Preventative Measures

In both scenarios, the implications are widespread, affecting not only the infrastructure and safety of the plant and its immediate environment but also the trust in technology and AI systems in critical setups. Therefore, **robust security measures**, which encompass not just the technical aspects (like firewalls, data usage monitoring, etc.) but also **policy and procedural aspects** (like limited access, frequent audits, etc.) should be in place.

User **Behavior Analytics (UBA)** could have been deployed in both cases to detect anomalous behavior in real-time. In addition, employing **AI explainability and interpretability tools** would allow for more transparency in how the AI model (Poseidon) generates its outputs, making it easier to detect malicious intents or modifications.

# 4.4 References for Section 4

[1]     Zhang, F., et al. (2022). *Patience Lamb / Asherah Training*. Retrieved from
        https://gitlab.com/lambpati/asherah-training

[2]     Zhang, F., et al. (2022). *iFAN Lab / Cyber Threat Assessment INL GT*. Retrieved from
        https://gitlab.com/ifan-lab/cyber-threat-assessment-inl-gt

[3]     Carlini, N., et al. (2021). Extracting Training Data from Large Language Models.
        *arXiv:2012.07805 [cs.CR]*.

[4]     Kokubun, T. (2023). *Gitstar ranking*. Retrieved from https://gitstar-ranking.com/

[5]     GitHub. (2023). *Licensing a repository*. Retrieved from
        https://docs.github.com/en/repositories/managing-your-repositorys-settings-and-
        features/customizing-your-repository/licensing-a-repository

[6]     Oh, M. G., et al. (2023). Membership Inference Attacks With Token-Level Deduplication on
        Korean Language Models. *IEEE Access, 11*, 10207-10217. doi:
        10.1109/ACCESS.2023.3239668

[7]     OWASP. (2023). *OWASP Top 10 for Large Language Model Applications*. Retrieved June 26,
        2023, from https://owasp.org/www-project-top-10-for-large-language-model-applications/

[8]     Branch, H. J., et al. (2022). Evaluating the Susceptibility of Pre-Trained Language Models via
        Handcrafted Adversarial Examples. *arXiv:2209.02128 [cs.CL]*.

[9]     Preamble. (2022). *Prompt Injection: A Critical Vulnerability in the GPT-3 Transformer, and
        How We Can Begin to Solve It*. Retrieved from https://www.preamble.com/prompt-
        injection-a-critical-vulnerability-in-the-gpt-3-transformer-and-how-we-can-begin-to-solve-it

[10]    Selvi, J. (2023). *Exploring Prompt Injection Attacks*. Retrieved from
        https://research.nccgroup.com/2022/12/05/exploring-prompt-injection-attacks/

[11]    Liu, Y., et al. (2018). Trojaning Attack on Neural Networks. In *25th Annual Network and
        Distributed System Security Symposium, NDSS 2018*. San Diego, California, USA: The Internet
        Society.

[12]    Mena, P., Borrelli, R. A., & Kerby, L. (2023). Detecting Anomalies in Simulated Nuclear Data
        Using Autoencoders. *Nuclear Technology, 0*(0), 1-14. doi: 10.1080/00295450.2023.2214257

[13]    Baldi, P. (2012). Autoencoders, Unsupervised Learning, and Deep Architectures. In I. Guyon
        et al. (Eds.), *Proceedings of ICML Workshop on Unsupervised and Transfer Learning, 27*, 37-
        49. PMLR. Retrieved from https://proceedings.mlr.press/v27/baldi12a.html

[14]    Fields, G., et al. (2021). Trojan Signatures in DNN Weights. In *2021 IEEE/CVF International
        Conference on Computer Vision Workshops (ICCVW)*, 12-20. Los Alamitos, CA, USA: IEEE
        Computer Society. doi:10.1109/ICCVW54120.2021.00008. Retrieved from
        https://doi.ieeecomputersociety.org/10.1109/ICCVW54120.2021.00008

[15]    Ilyas, A., et al. (2019). Adversarial Examples Are Not Bugs, They Are Features.
        *arXiv:1905.02175 [stat.ML]*.

[16]    Wang, G., Yan, H., & Wei, X. (2022). Enhancing Transferability of Adversarial Examples with
        Spatial Momentum. *arXiv:2203.13479*. Retrieved from https://arxiv.org/abs/2203.13479

*Page intentionally left blank.*