

Adversarial Machine Learning

Arthea Valderrama

Faculty Advisor: Dr. Ian Chen

Committee Member: Dr. Sirong Lin

November 20, 2025

Contents

1 Abstract	2
2 Motivation	2
3 Attack Scenarios	3
3.1 Carlini & Wagner	3
3.2 Simple Black-box Attack	4
4 Anomaly Detection	4
5 Procedure	5
6 Results	5
7 Conclusion	7
A Appendix A: Model Creation	8
B Appendix B: Attack Implementation	12
C Appendix C: Anomaly Detection	25
D Appendix D: Web Application Deployment	29

1 Abstract

Machine Learning (ML) has become integral to daily life, due to its use in autonomous systems, the health-care industry, etc., because of its capability to produce a desired outcome based off of learning from training data. However, adversarial attacks attempt to deceive learning models by implementing a defective input, compromising the effectiveness as a result. We explore the mechanisms behind attacks on convolutional neural networks (CNNs). This project aims to contribute to the development of safer ML models to ensure their reliability, through evaluating the effectiveness of adversarial attacks.

2 Motivation

This is an addition to my current capstone project through UML's collaboration with Schneider Electric. The company is a multinational corporation specializing in energy management and automation. One of the company's ambitions is ensuring robust automation processes for increased efficiency and quality of services.

In an autonomous setting that requires quality assurance for real-time classification, it is important that a model is adaptable to various conditions to ensure quality decision making. In extension to the existing project's image classification system on household objects in an assembly line, I expanded the dataset for cancer detection in a healthcare setting to demonstrate the transferability of ML models and its attacks.

ML assisted detection of brain cancer through MRI imaging can save resources and improve patient outcomes through early detection and accurate patient prognosis estimation. However, an adversarial attack would lead to errors in tumor recognition, harming patient health and treatment. Such an attack could occur for a variety of reasons: undermining trust in critical systems, financial gain, etc. By identifying vulnerabilities in brain cancer imaging detection we can ensure the security and trustworthiness of Ai.

3 Attack Scenarios

3.1 Carlini & Wagner

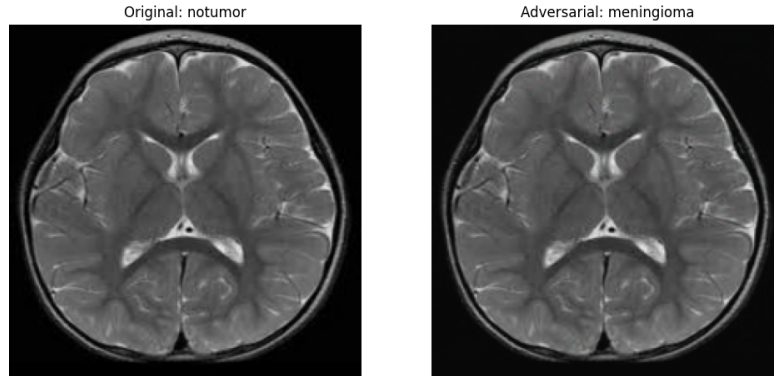


Figure 1: Example implementation of C&W attack

The Carlini & Wagner (C&W)[1] attack finds the smallest perturbation to a model's input data, while also increasing the chances of model misclassification through optimizing an iterative algorithm. The algorithm is defined as: **minimize** $\|\delta\| + c * f(x + \delta)$, **such that** $x + \delta \in [0, 1]^n$, where δ is a small change made to an image x that changes its classification, and c is a constant that advances the trade-off between the size of the change and effectiveness of the attack. It is a white-box attack, meaning the attacker can launch the attack if they have knowledge on the model's information, such as its architecture and its parameters.

3.2 Simple Black-box Attack

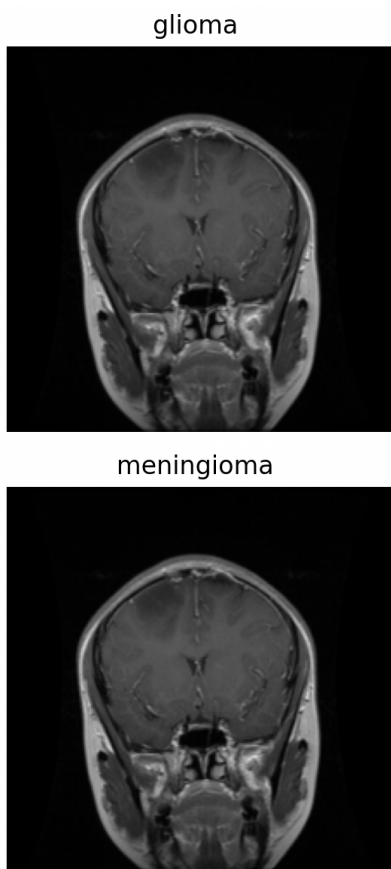


Figure 2: Small simBA distortion

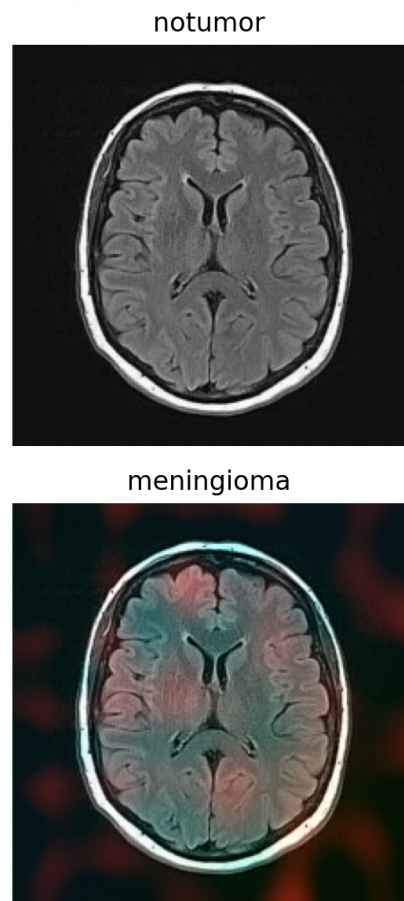


Figure 3: Higher distortion of simBA

While the C&W attack is a white-box attack that requires the adversary to have full knowledge of the model it is attacking, the Simple Black-box Attack (simBA)[2] is a type of black-box attack that does not need complete information about the model. In that case, simBA is more applicable in many scenarios. SimBA works by performing a random search on the pixels of an image to find minuscule perturbations in order to change the image away from its predicted classification and continuously checking the probability of success in changing input.

4 Anomaly Detection

An anomaly detection [3] model is a specialized CNN that identifies data points that significantly deviate from "expected" within the dataset. The model works by taking in a binary input: the normal patterns within the dataset and anomalous samples that deviate from the learned patterns. By incorporating aberrant samples in its training, the model is then able to flag unusual input. This can protect an AI-based system in real-time from instances that could result in harmful losses. Other use cases of anomaly detection within AI applications are detecting fraudulent financial transactions and suspicious network traffic in cybersecurity.

5 Procedure

A crucial factor in ML training is the selection of a robust dataset. I obtained a brain tumor MRI dataset of 7,022 images[4], and split the dataset 80% to train the model and 20% to evaluate the model's accuracy. I chose to select a large dataset for improved generalization with more diverse data samples and to reduce overfitting. The dataset includes 4 classifications of tumors: no tumor, glioma (tumors in the glial cells in the brain or spinal cord), meningioma (tumors that develops in the meninges), and pituitary (tumors in the pituitary gland).

In accordance with my capstone project, a convolutional neural network (CNN) was the model architecture chosen, specifically in the form of a DenseNet model. This model architecture consists of dense layers where each layer is connected to its previous layers in a feed-forward fashion. The DenseNet model is helpful for image classification, as its compact network improves information propagation through the network, has an enhanced ability to learn complex features, and requires less parameters in deep whilst not needing much computation.

The main programming language used in this project is Python. Python is ideal for working with ML, as it has numerous libraries for data pre-processing and evaluation, such as Sci-kit Learn, Numpy, Pandas, and Matplotlib. The main library used for ML is PyTorch, which allowed me to import a wide range of pretrained Neural Networks to improve model accuracy and functionality.

When coding my training script in Python to train the model, I created a class to preprocess the MRI images into classes for the model to distinguish between, and another class to define the model's parameters and functions. The model's class consists of a initialization function that uses features from Pytorch's built-in pretrained DenseNet model, a forward function to define the model's learning and optimization, and a training function to define the number of iterations the model passes through its learning algorithm while tracking its loss. Utilizing a pretrained model improved the model's learning rate and testing accuracy on samples it has not seen before.

Creating the attack code in Python, I modified widely available C&W and simBA attack code to apply to my MRI imaging model. Thus, helping to demonstrate the generalizability of the attack across different imaging models. I applied the attacks on the original dataset of 7,022 MRI images to generate enough adversarial samples for the anomaly detection model.

With the anomaly detection model, I then combined both the original training data and adversarial data. Thus, in total it learned the patterns and behavior between 14,044 MRI images to recognize between suspicious outliers and regular input.

Lastly, to create a simple web app to test the models on a user-friendly interface. Through the web app, I used the Flask library, HTML, CSS, and JavaScript code for live predictions on the trained dataset. I then incorporated Google's AI Gemini to detail more information about adversarial attacks, information about diagnostics, and educate the implications of the model's findings. A chat-bot was added to allow the user to talk more with the AI for any follow-up questions.

6 Results

Through testing over a sample batch size of 32 images, I was able to calculate the attack success rate of the C&W, simBA, and Anomaly Detection model. The C&W attack is more successful than simBA by a

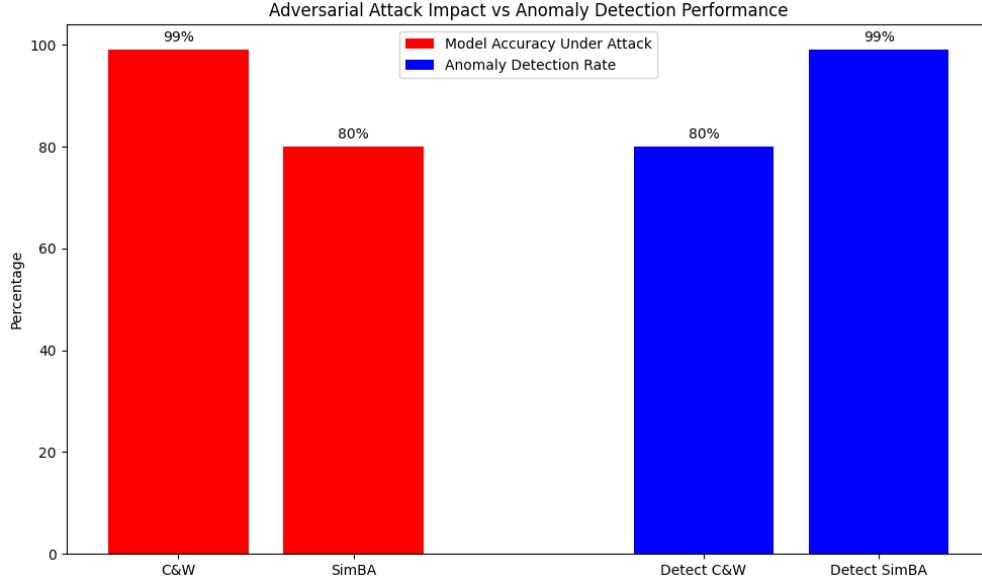


Figure 4: Success Rate of Attacks and Anomaly Detections

factor of around 20%. The success rate of the anomaly detection method on the C&W attack proves the effectiveness further, as the model has detected less C&W samples than the simBA method. With a higher

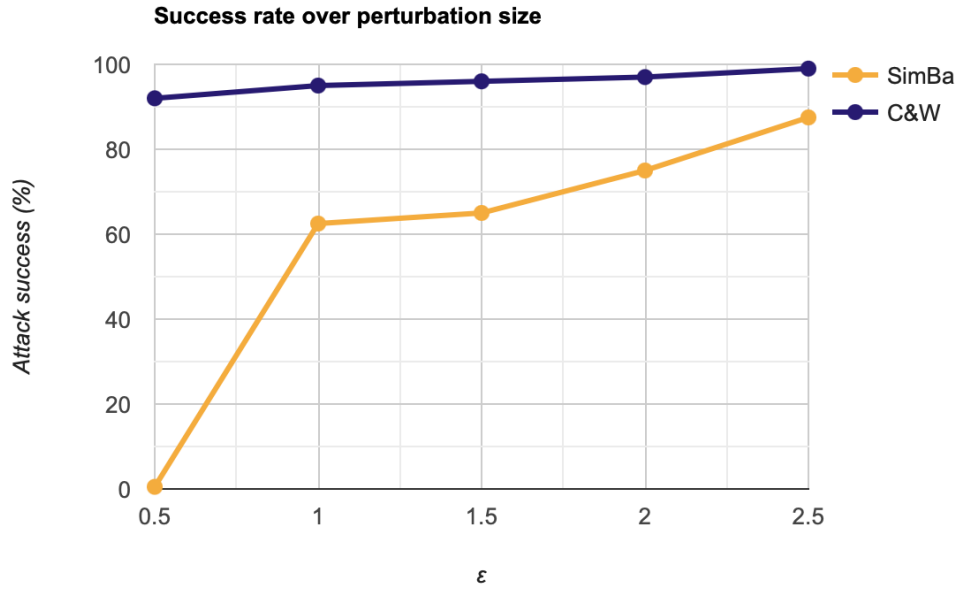


Figure 5: Success Rate of Attacks and Anomaly Detections

δ , the success rates of attacks tend to be higher but cause more visual distortions. The C&W generally does not need a high δ to be effective, although simBA requires a factor of 1.0+ to start becoming successful in fooling the model. Another comparison that can be made between the two models is that the C&W, while more effective, requires much more time than the simBA attack to deploy on average.

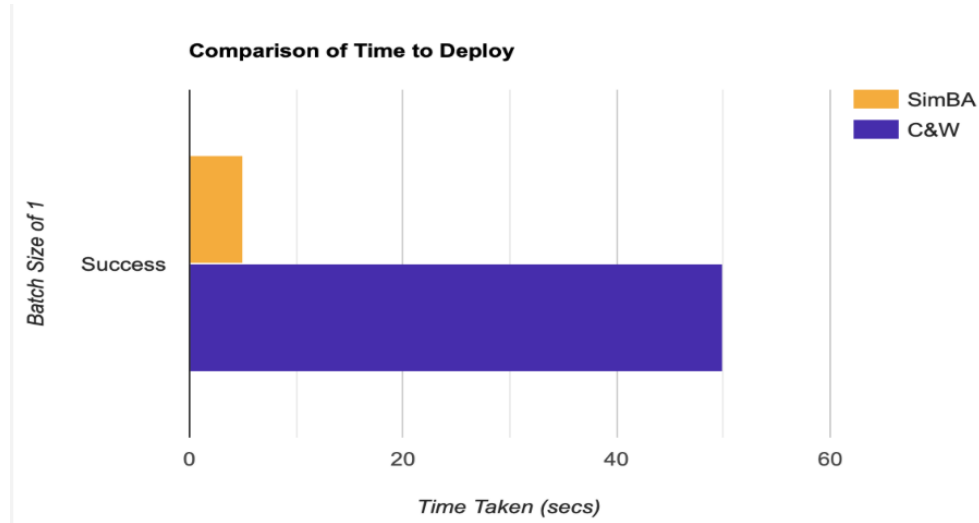


Figure 6: Time taken to deploy the attack

7 Conclusion

As AI systems become increasingly integrated into society, it is very important that security considerations are taken when deploying and creating AI systems. Models often influence critical decisions, and vulnerabilities that exploit these models can have potentially devastating consequences in quality of life, finances, etc. The example adversarial attacks of C&W and simBA showcase the generalizability and adaptability of these attacks across different use cases. These attacks can have devastating consequences on industries and organizations that rely on computer vision or image classification. Without strong security measures and further research into the possibilities of attacks on different types of AI models, the risk of model exploitation could harm societal trust in AI technology. Thus, ensuring the security of AI is essential to safeguard against manipulation, and uphold the reliability of automated decision-making.

A Appendix A: Model Creation

The following code document details the creation of an image classification model, originally imported from my senior capstone project and adapted to be utilized for MRI imaging purposes.

```
1 import torch
2 import torchvision.transforms as transforms
3 from torch.utils.data import Dataset, DataLoader
4 import os
5 from torch import nn
6 import torch.nn.functional as F
7 from torchvision import models
8 import matplotlib.pyplot as plt
9 from PIL import Image
10 import numpy as np
11 from tqdm import tqdm
12
13 class CustomDataset(Dataset):
14     def __init__(self, root_dir, transform=None):
15         self.root_dir = root_dir
16         self.transform = transform
17         self.classes = ['glioma', 'meningioma', 'notumor', 'pituitary']
18         self.image_paths = []
19         self.labels = []
20
21         for idx, class_name in enumerate(self.classes):
22             class_dir = os.path.join(root_dir, class_name)
23             for image_name in os.listdir(class_dir):
24                 if image_name.endswith(".jpg") or image_name.endswith
25                 (".png"):
26                     self.image_paths.append(os.path.join(class_dir,
27 image_name))
28                     self.labels.append(idx)
29
30     def __len__(self):
31         return len(self.image_paths)
32
33     def __getitem__(self, idx):
34         img_name = self.image_paths[idx]
35         image = Image.open(img_name).convert('RGB')
36         label = self.labels[idx]
37
38         if self.transform:
39             image = self.transform(image)
40
41         return image, label
42
43 resize_transform = transforms.Compose([
```

```

42     transforms.Resize((64, 64)),
43     transforms.ToTensor()
44 ])
45
46 class_names = ['glioma', 'meningioma', 'notumor', 'pituitary']
47 class_id_to_name = {i: class_name for i, class_name in enumerate(
    class_names)}
48 class_to_idx = {class_name: i for i, class_name in enumerate(
    class_names)}
49
50 class FineTunedDenseNet(nn.Module):
51     def __init__(self, pretrained_model, num_classes=4):
52         super(FineTunedDenseNet, self).__init__()
53         self.features = pretrained_model.features
54         self.classifier = nn.Linear(pretrained_model.classifier.
in_features, num_classes)
55
56     def forward(self, x):
57         x = self.features(x)
58         x = F.adaptive_avg_pool2d(x, (1, 1)).view(x.size(0), -1)
59         x = self.classifier(x)
60         return x
61
62     def train_model(self, model_name, input_size, epochs,
learning_rate, train_loader, test_loader, device):
63         criterion = nn.CrossEntropyLoss()
64         optimizer = torch.optim.Adam(self.parameters(), lr=
learning_rate)
65         scheduler = torch.optim.lr_scheduler.StepLR(optimizer,
step_size=5, gamma=0.1)
66
67         best_accuracy = 0
68         train_loss_values = []
69         test_accuracy_values = []
70
71         for epoch in range(epochs):
72             self.train()
73             running_loss = 0.0
74             correct = 0
75             total = 0
76             progress_bar = tqdm(train_loader, desc=f"Epoch {epoch+1}/{
epochs}", leave=False, ncols=100)
77
78             for i, (images, labels) in enumerate(progress_bar):
79                 images, labels = images.to(device), labels.to(device)
80
81                 optimizer.zero_grad()
82                 outputs = self(images)

```

```

83         loss = criterion(outputs, labels)
84         loss.backward()
85         optimizer.step()
86
87         running_loss += loss.item()
88         _, predicted = torch.max(outputs, 1)
89         total += labels.size(0)
90         correct += (predicted == labels).sum().item()
91
92         progress_bar.set_postfix(loss=running_loss / (i + 1),
accuracy=100 * correct / total)
93
94         test_accuracy = 100 * correct / total
95         if test_accuracy > best_accuracy:
96             best_accuracy = test_accuracy
97             torch.save(self.state_dict(), "best_densenet_model.pth
")
98
99         scheduler.step()
100         train_loss_values.append(running_loss / len(train_loader))
101         test_accuracy_values.append(test_accuracy)
102
103         self.plot_loss(model_name, input_size, learning_rate,
train_loss_values, test_accuracy_values)
104
105         def plot_loss(self, model_name, input_size, learning_rate,
train_loss_values, test_accuracy_values):
106             plt.figure(figsize=(10, 5))
107             plt.plot(train_loss_values, label="Training Loss", color='blue
')
108             plt.xlabel("Training Steps")
109             plt.ylabel("Loss")
110             plt.title("Training Loss Over Time")
111             plt.grid(True)
112             plt.legend()
113             plt.savefig(f"training_loss_{model_name}_{input_size}_lr{
learning_rate}.png")
114             plt.show()
115
116             plt.figure(figsize=(10, 5))
117             plt.plot(test_accuracy_values, label="Test Accuracy", color='
orange')
118             plt.xlabel("Steps")
119             plt.ylabel("Accuracy (%)")
120             plt.title("Test Accuracy Over Time")
121             plt.grid(True)
122             plt.legend()

```

```

123     plt.savefig(f"test_accuracy_{model_name}_{input_size}_lr{
learning_rate}.png")
124     plt.show()
125
126     def save(self, model_file):
127         torch.save(self.state_dict(), model_file)
128
129 if __name__ == "__main__":
130     root_dir = '/Users/arthe/honors/brain/training'
131     model_name = 'FineTunedDenseNet'
132     batch_size = 32
133     learning_rate = 0.0003
134     epochs = 10
135     input_size = 64
136
137     image_transforms = transforms.Compose([
138         transforms.RandomResizedCrop(input_size),
139         transforms.RandomHorizontalFlip(),
140         transforms.ColorJitter(brightness=0.2, contrast=0.2, saturation
=0.2, hue=0.2),
141         transforms.ToTensor(),
142         transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
0.224, 0.225]),
143     ])
144
145     full_dataset = CustomDataset(root_dir, transform=image_transforms)
146     train_size = int(0.8 * len(full_dataset))
147     test_size = len(full_dataset) - train_size
148     train_dataset, test_dataset = torch.utils.data.random_split(
full_dataset, [train_size, test_size])
149     train_loader = DataLoader(train_dataset, batch_size=batch_size,
shuffle=True, num_workers=4)
150     test_loader = DataLoader(test_dataset, batch_size=batch_size,
shuffle=False, num_workers=4)
151
152     print(f"Training {model_name}")
153     device = torch.device('cuda' if torch.cuda.is_available() else 'cpu
')
154     print(f"Using device {device}")
155
156     pretrained_densenet = models.densenet121(pretrained=True)
157     model = FineTunedDenseNet(pretrained_model=pretrained_densenet,
num_classes=4)
158     model.to(device)
159
160     model.train_model(model_name, input_size, epochs, learning_rate,
train_loader, test_loader, device=device)
161

```

```

162     model.to('cpu')
163     model_path = f"{model_name}_{input_size}_lr{learning_rate}.pt"
164     torch.save(model.state_dict(), model_path)
165     print(f"Model saved at {model_path}")

```

Listing 1: train.py

B Appendix B: Attack Implementation

This section details the Python implementation of the C&W[5] attack and the simBA[6] attack, originally sourced from GitHub and then modified to suit the image classification model being used in this project. The first document is my implementation of the C&W attack, with additional functions to save adversarial samples and display the original classification and modded image for viewing purposes. The next lines of code is my simBA implementation, with similar functions to save adversarial samples and display the difference between the original image and the modded image.

```

1 import os
2 import torch
3 import time
4 import torch.nn.functional as F
5 import torch.optim as optim
6 from torchvision import models
7 from torch import nn
8 from torchvision import transforms
9 import matplotlib.pyplot as plt
10 from PIL import Image
11
12 from train import class_names, class_id_to_name, class_to_idx,
13     CustomDataset
14 from train import FineTunedDenseNet
15
16 model_file = "/Users/arthe/honors/FineTunedDenseNet_64_lr0.0003.pt"
17
18 script_dir = os.path.dirname(os.path.abspath(__file__))
19 adversarial_save_dir = os.path.join(script_dir, "adversarial_images1")
20 os.makedirs(adversarial_save_dir, exist_ok=True)
21
22 resize_transform = transforms.Resize((224, 224), interpolation=
23     transforms.InterpolationMode.BILINEAR)
24
25 full_transform = transforms.Compose([
26     resize_transform,
27     transforms.ToTensor()
28 ])
29
30 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
31 pretrained_densenet = models.densenet121(pretrained=True)

```

```

30
31 model = FineTunedDenseNet(pretrained_model=pretrained_densenet,
    num_classes=len(class_names))
32 model.load_state_dict(torch.load(model_file, map_location=device))
33 model.to(device)
34 model.eval()
35
36 def cw_l2_attack(model, images, labels, target_labels=None, targeted=
    False, c=2.5, kappa=0, max_iter=1000, learning_rate=0.01):
37     """
38     CW-L2 Attack to generate adversarial examples on full-resolution
    images.
39     Args:
40         model (nn.Module): Trained model to attack.
41         images (torch.Tensor): Batch of full-resolution images to
    attack.
42         labels (torch.Tensor): True labels for the images.
43         target_labels (torch.Tensor, optional): Target labels for the
    attack. Required if targeted=True.
44         targeted (bool): Whether the attack is targeted or untargeted.
45         c (float): Regularization parameter to control attack strength.
46         kappa (float): Confidence parameter for the attack.
47         max_iter (int): Maximum number of optimization iterations.
48         learning_rate (float): Learning rate for the optimizer.
49     Returns:
50         torch.Tensor: Adversarial examples in full resolution.
51     """
52
53     images = images.to(device)
54     labels = labels.to(device)
55
56     if targeted:
57         if target_labels is None:
58             raise ValueError("Target labels must be provided for
    targeted attacks.")
59         target_labels = target_labels.to(device)
60
61     def f(x):
62         resized_x = resize_transform(x)
63         outputs = model(resized_x)
64
65         if targeted:
66             one_hot_labels = torch.eye(outputs.shape[1], device=device)
    [target_labels]
67             i, _ = torch.max((1 - one_hot_labels) * outputs, dim=1)
68             j = torch.masked_select(outputs, one_hot_labels.bool())
69             return torch.clamp(i - j, min=-kappa)
70     else:

```

```

71         one_hot_labels = torch.eye(outputs.shape[1], device=device)
[labels]
72         i, _ = torch.max((1 - one_hot_labels) * outputs, dim=1)
73         j = torch.masked_select(outputs, one_hot_labels.bool())
74         return torch.clamp(j - i, min=-kappa)
75
76     w = torch.zeros_like(images, requires_grad=True).to(device)
77     optimizer = optim.Adam([w], lr=learning_rate)
78     prev = 1e10
79
80     for step in range(max_iter):
81         a = 0.5 * (torch.tanh(w) + 1)
82
83         loss1 = F.mse_loss(a, images, reduction='sum')
84         loss2 = torch.sum(c * f(a))
85         cost = loss1 + loss2
86
87         optimizer.zero_grad()
88         cost.backward()
89         optimizer.step()
90
91         if step % (max_iter // 10) == 0:
92             if cost > prev:
93                 print("Attack Stopped due to CONVERGENCE....")
94                 return a
95                 prev = cost
96
97         if (step + 1) % 100 == 0 or step == 0:
98             print(f"- Learning Progress: {(step + 1) / max_iter *
100:.2f}% ", end="\r")
99
100     attack_images = 0.5 * (torch.tanh(w) + 1)
101     return attack_images
102
103 def save_image(image, true_label, predicted_label, original_size, index
, prefix="adversarial_example"):
104
105     image = transforms.ToPILImage()(image.cpu().detach())
106
107     image_resized = image.resize(original_size, Image.LANCZOS)
108
109     os.makedirs(adversarial_save_dir, exist_ok=True)
110
111     file_name = os.path.join(
112         adversarial_save_dir,
113         f"{prefix}_{index}_true_{true_label}_pred_{predicted_label}.png
"
114     )

```

```

115     image_resized.save(file_name)
116
117
118 def display_images(original_image, adversarial_image, true_label,
119                    predicted_label):
120
121     fig, axes = plt.subplots(1, 2, figsize=(12, 6))
122
123     axes[0].imshow(original_image.permute(1, 2, 0).cpu().detach().numpy
124                     ())
125     axes[0].set_title(f"Original: {true_label}")
126     axes[0].axis('off')
127
128     axes[1].imshow(adversarial_image.permute(1, 2, 0).cpu().detach().
129                   numpy())
130     axes[1].set_title(f"Adversarial: {predicted_label}")
131     axes[1].axis('off')
132
133     plt.show()
134
135 if __name__ == "__main__":
136     import argparse
137
138     targeted = False
139     target_class = "glioma"
140     c = 20
141     kappa = 0
142     learning_rate = 0.01
143     max_iter = 1000
144
145     data_folder = "/Users/arthe/honors/brain/training"
146     full_dataset = CustomDataset(data_folder, transform=full_transform)
147     train_size = int(0.8 * len(full_dataset))
148     test_size = len(full_dataset) - train_size
149     train_dataset, test_dataset = torch.utils.data.random_split(
150     full_dataset, [train_size, test_size])
151
152     batch_size = 1
153     test_loader = torch.utils.data.DataLoader(test_dataset, batch_size=
154     batch_size, shuffle=False)
155
156     attack_accuracy = 0
157     correct = 0
158     total = 0
159
160     for images, labels in test_loader:
161         images, labels = images.to(device), labels.to(device)

```



```

158     original_size = images.shape[2], images.shape[3]
159
160     start_time = time.time()
161
162     if targeted:
163         if target_class is None:
164             raise ValueError("Target class must be specified for
targeted attacks.")
165         try:
166             target_class_idx = class_to_idx[target_class]
167         except KeyError:
168             raise ValueError(f"Target class '{target_class}' not
found in class_to_idx.")
169         target_labels = torch.full_like(labels, target_class_idx)
170         adversarial_images = cw_l2_attack(
171             model, images, labels, target_labels=target_labels,
targeted=True, c=c,
172             kappa=kappa, learning_rate=learning_rate, max_iter=
max_iter
173         )
174     else:
175         adversarial_images = cw_l2_attack(
176             model, images, labels, targeted=False, c=c,
177             kappa=kappa, learning_rate=learning_rate, max_iter=
max_iter
178         )
179
180     end_time = time.time()
181     attack_time = end_time - start_time
182     print(f"Time taken for attack: {attack_time:.2f} seconds")
183
184     adversarial_images_resized = resize_transform(
adversarial_images)
185
186     outputs = model(adversarial_images_resized)
187     _, predicted = torch.max(outputs, 1)
188
189     if targeted:
190         batch_correct = (predicted == target_labels).sum().item()
191     else:
192         batch_correct = (predicted != labels).sum().item()
193
194     batch_total = labels.size(0)
195     attack_accuracy += 100 * batch_correct / batch_total
196     correct += batch_correct
197     total += batch_total
198

```

```

199     print(f"Attack Accuracy on the batch: {100 * batch_correct /
200           batch_total:.2f}%")
201
202     for i in range(len(adversarial_images)):
203         true_label = class_names[labels[i].item()]
204         predicted_label = class_names[predicted[i].item()]
205
206         save_image(adversarial_images[i], true_label,
207                   predicted_label, original_size, i)
208
209         display_images(images[i], adversarial_images[i], true_label
210                       , predicted_label)
211
212     overall_attack_accuracy = attack_accuracy / len(test_loader)
213     print(f"\nOverall Attack Accuracy: {overall_attack_accuracy:.2f}%")

```

Listing 2: attack.py

```

1 import torch
2 import utils
3 import torch.nn.functional as F
4 import torchvision.models as models
5 from torchvision.models import densenet121
6 import torchvision.transforms as transforms
7 from torchvision.utils import make_grid
8 import matplotlib.pyplot as plt
9 import time
10 from torchvision.datasets import ImageFolder
11 from torch.utils.data import DataLoader
12
13
14 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
15
16 class SimBA:
17
18     def __init__(self, model, dataset, image_size):
19         self.model = model.to(device)
20         self.dataset = dataset
21         self.image_size = image_size
22         self.model.eval()
23
24     def expand_vector(self, x, size):
25         batch_size = x.size(0)
26         x = x.view(-1, 3, size, size)
27         z = torch.zeros(batch_size, 3, self.image_size, self.image_size
28                       )
29         z[:, :, :size, :size] = x
30         return z

```

```

31     def normalize(self, x):
32         return utils.apply_normalization(x.to(device), self.dataset)
33
34
35     def get_probs(self, x, y):
36         output = self.model(self.normalize(x)).cpu()
37         probs = torch.index_select(F.softmax(output, dim=-1).data, 1, y
38     )
39         return torch.diag(probs)
40
41     def get_preds(self, x):
42         output = self.model(self.normalize(x)).cpu()
43         _, preds = output.data.max(1)
44         return preds
45
46     def simba_single(self, x, y, num_iters=10000, epsilon=2.5, targeted
47     =False):
48         n_dims = x.view(1, -1).size(1)
49         perm = torch.randperm(n_dims)
50         last_prob = self.get_probs(x, y)
51         for i in range(num_iters):
52             diff = torch.zeros(n_dims)
53             diff[perm[i]] = epsilon
54             left_prob = self.get_probs((x - diff.view(x.size()))).clamp
55             (0, 1), y)
56             if targeted != (left_prob < last_prob):
57                 x = (x - diff.view(x.size()))).clamp(0, 1)
58                 last_prob = left_prob
59             else:
60                 right_prob = self.get_probs((x + diff.view(x.size()))).
61                 clamp(0, 1), y)
62                 if targeted != (right_prob < last_prob):
63                     x = (x + diff.view(x.size()))).clamp(0, 1)
64                     last_prob = right_prob
65             if i % 10 == 0:
66                 print(last_prob)
67         return x.squeeze()
68
69
70     def simba_batch(self, images_batch, labels_batch, max_iters,
71     freq_dims, stride, epsilon, linf_bound=0.0,
72     order='rand', targeted=False, pixel_attack=False,
73     log_every=1):
74         batch_size = images_batch.size(0)
75         image_size = images_batch.size(2)
76         assert self.image_size == image_size
77         if order == 'rand':
78             indices = torch.randperm(3 * freq_dims * freq_dims)[:

```

```

max_iters]
73     elif order == 'diag':
74         indices = utils.diagonal_order(image_size, 3)[:max_iters]
75     elif order == 'strided':
76         indices = utils.block_order(image_size, 3, initial_size=
freq_dims, stride=stride)[:max_iters]
77     else:
78         indices = utils.block_order(image_size, 3)[:max_iters]
79     if order == 'rand':
80         expand_dims = freq_dims
81     else:
82         expand_dims = image_size
83     n_dims = 3 * expand_dims * expand_dims
84     x = torch.zeros(batch_size, n_dims)
85     probs = torch.zeros(batch_size, max_iters)
86     succs = torch.zeros(batch_size, max_iters)
87     queries = torch.zeros(batch_size, max_iters)
88     l2_norms = torch.zeros(batch_size, max_iters)
89     linf_norms = torch.zeros(batch_size, max_iters)
90     prev_probs = self.get_probs(images_batch, labels_batch)
91     preds = self.get_preds(images_batch)
92     if pixel_attack:
93         trans = lambda z: z
94     else:
95         trans = lambda z: utils.block_idct(z, block_size=image_size
, linf_bound=linf_bound)
96     remaining_indices = torch.arange(0, batch_size).long()
97     for k in range(max_iters):
98         dim = indices[k]
99         expanded = (images_batch[remaining_indices] + trans(self.
expand_vector(x[remaining_indices], expand_dims))).clamp(0, 1)
100         perturbation = trans(self.expand_vector(x, expand_dims))
101         l2_norms[:, k] = perturbation.view(batch_size, -1).norm(2,
1)
102         linf_norms[:, k] = perturbation.view(batch_size, -1).abs().
max(1)[0]
103         preds_next = self.get_preds(expanded)
104         preds[remaining_indices] = preds_next
105         if targeted:
106             remaining = preds.ne(labels_batch)
107         else:
108             remaining = preds.eq(labels_batch)
109         if remaining.sum() == 0:
110             adv = (images_batch + trans(self.expand_vector(x,
expand_dims))).clamp(0, 1)
111             probs_k = self.get_probs(adv, labels_batch)
112             probs[:, k:] = probs_k.unsqueeze(1).repeat(1, max_iters
- k)

```

```

113         succs[:, k:] = torch.ones(batch_size, max_iters - k)
114         queries[:, k:] = torch.zeros(batch_size, max_iters - k)
115         break
116     remaining_indices = torch.arange(0, batch_size)[remaining].
long()
117     if k > 0:
118         succs[:, k] = ~remaining
119         diff = torch.zeros(remaining.sum(), n_dims)
120         diff[:, dim] = epsilon
121         left_vec = x[remaining_indices] - diff
122         right_vec = x[remaining_indices] + diff
123         adv = (images_batch[remaining_indices] + trans(self.
expand_vector(left_vec, expand_dims))).clamp(0, 1)
124         left_probs = self.get_probs(adv, labels_batch[
remaining_indices])
125         queries_k = torch.zeros(batch_size)
126         queries_k[remaining_indices] += 1
127         if targeted:
128             improved = left_probs.gt(prev_probs[remaining_indices])
129         else:
130             improved = left_probs.lt(prev_probs[remaining_indices])
131         if improved.sum() < remaining_indices.size(0):
132             queries_k[remaining_indices[~improved]] += 1
133         adv = (images_batch[remaining_indices] + trans(self.
expand_vector(right_vec, expand_dims))).clamp(0, 1)
134         right_probs = self.get_probs(adv, labels_batch[
remaining_indices])
135         if targeted:
136             right_improved = right_probs.gt(torch.max(prev_probs[
remaining_indices], left_probs))
137         else:
138             right_improved = right_probs.lt(torch.min(prev_probs[
remaining_indices], left_probs))
139         probs_k = prev_probs.clone()
140         if improved.sum() > 0:
141             left_indices = remaining_indices[improved]
142             left_mask_remaining = improved.unsqueeze(1).repeat(1,
n_dims)
143             x[left_indices] = left_vec[left_mask_remaining].view
(-1, n_dims)
144             probs_k[left_indices] = left_probs[improved]
145         if right_improved.sum() > 0:
146             right_indices = remaining_indices[right_improved]
147             right_mask_remaining = right_improved.unsqueeze(1).
repeat(1, n_dims)
148             x[right_indices] = right_vec[right_mask_remaining].view
(-1, n_dims)
149             probs_k[right_indices] = right_probs[right_improved]

```

```

150         probs[:, k] = probs_k
151         queries[:, k] = queries_k
152         prev_probs = probs[:, k]
153         if (k + 1) % log_every == 0 or k == max_iters - 1:
154             print('Iteration %d: queries = %.4f, prob = %.4f,
remaining = %.4f' % (
155                 k + 1, queries.sum(1).mean(), probs[:, k].mean
(), remaining.float().mean()))
156             expanded = (images_batch + trans(self.expand_vector(x,
expand_dims))).clamp(0, 1)
157             preds = self.get_preds(expanded)
158             if targeted:
159                 remaining = preds.ne(labels_batch)
160             else:
161                 remaining = preds.eq(labels_batch)
162             succs[:, max_iters-1] = ~remaining
163             return expanded, probs, succs, queries, l2_norms, l1_norms
164 transform = transforms.Compose([
165     transforms.Resize((224, 224)),
166     transforms.ToTensor(),
167 ])
168
169
170 dataset_path = '/Users/arthe/honors/brain/training'
171 dataset = ImageFolder(root=dataset_path, transform=transform)
172 dataloader = DataLoader(dataset, batch_size=2, shuffle=True)
173
174 model = densenet121(pretrained=True)
175 model.classifier = torch.nn.Linear(in_features=1024, out_features=4)
176
177
178 model.load_state_dict(torch.load('/Users/arthe/honors/
FineTunedDenseNet_64_lr0.0003.pt', map_location=device))
179 model.eval()
180
181 simba = SimBA(model=model, dataset='imagenet', image_size=224)#was 224
182
183 def simba_single_attack(simba, dataloader):
184     images_batch, labels_batch = next(iter(dataloader))
185     image = images_batch[0:1].to(device)
186     label = labels_batch[0:1].to(device)
187
188     class_name = dataloader.dataset.classes[label.item()]
189
190     start_time = time.time()
191     adv_image = simba.simba_single(image, label, num_iters=100, epsilon
=0.2, targeted=False)
192     elapsed_time = time.time() - start_time

```

```

193     adv_pred = simba.get_preds(adv_image.unsqueeze(0))
194     adv_class_name = dataloader.dataset.classes[adv_pred.item()]
195
196
197     plt.figure(figsize=(10, 5))
198
199     plt.subplot(1, 2, 1)
200     plt.imshow(image[0].permute(1, 2, 0).cpu().numpy())
201     plt.title(f"Original: {class_name} ({label.item()})")
202     plt.axis('off')
203
204     plt.subplot(1, 2, 2)
205     plt.imshow(adv_image.permute(1, 2, 0).cpu().numpy())
206     plt.title(f"Adversarial: {adv_class_name} ({adv_pred.item()})")
207     plt.axis('off')
208
209     plt.suptitle(f"SimBA Single Attack | Time: {elapsed_time:.2f}s")
210     plt.show()
211
212
213
214 def simba_full_dataset_attack(simba, dataloader, max_iters=200, epsilon
215 =0.2, freq_dims=21, stride=1, targeted=False):
216     total_images = 0
217     total_fooled = 0
218     total_time = 0.0
219
220     for images_batch, labels_batch in dataloader:
221         images_batch = images_batch.to(device)
222         labels_batch = labels_batch.to(device)
223
224         start_time = time.time()
225         adv_images, probs, succs, queries, l2_norms, l1_norms = simba
226         .simba_batch(
227             images_batch=images_batch,
228             labels_batch=labels_batch,
229             max_iters=max_iters,
230             freq_dims=freq_dims,
231             stride=stride,
232             epsilon=epsilon,
233             l1_bound=0.0,
234             order='rand',
235             targeted=targeted,
236             pixel_attack=False,
237             log_every=50
238         )
239         elapsed_time = time.time() - start_time
240         total_time += elapsed_time

```

```

239     adv_preds = simba.get_preds(adv_images)
240     fooled = (adv_preds != labels_batch).sum().item()
241     total_foiled += fooled
242     total_images += len(labels_batch)
243
244
245     print(f"Batch: {total_images}/{len(dataloader.dataset)} -
246     Fooled: {fooled}/{len(labels_batch)} - Time: {elapsed_time:.2f}s")
247
248     success_rate = (total_foiled / total_images) * 100
249     avg_time = total_time / len(dataloader)
250
251     print(f"\nSimBA Attack Completed")
252     print(f"Total Images: {total_images}")
253     print(f"Total Fooled: {total_foiled}")
254     print(f"Attack Success Rate: {success_rate:.2f}%")
255     print(f"Average Time per Batch: {avg_time:.2f}s")
256
257     return success_rate
258
259
260 def evaluate_adversarial_attacks(model, dataloader, targeted=False,
261     max_iters=200, freq_dims=21, stride=1, epsilon=2.5):
262
263     model.eval()
264     attack_accuracy = 0
265     correct = 0
266     total = 0
267
268     for images, labels in dataloader:
269         images, labels = images.to(device), labels.to(device)
270
271         start_time = time.time()
272
273         if targeted:
274             if target_class is None:
275                 raise ValueError("Target class must be specified for
276                 targeted attacks.")
277             try:
278                 target_class_idx = dataloader.dataset.class_to_idx[
279                 target_class]
280             except KeyError:
281                 raise ValueError(f"Target class '{target_class}' not
282                 found in class_to_idx.")
283             target_labels = torch.full_like(labels, target_class_idx)
284             adv_images, *_ = simba.simba_batch(
285                 images_batch=images,

```



```

282         labels_batch=labels,
283         max_iters=max_iters,
284         freq_dims=freq_dims,
285         stride=stride,
286         epsilon=epsilon,
287         targeted=True
288     )
289 else:
290     adv_images, *_ = simba.simba_batch(
291         images_batch=images,
292         labels_batch=labels,
293         max_iters=max_iters,
294         freq_dims=freq_dims,
295         stride=stride,
296         epsilon=epsilon,
297         targeted=False
298     )
299
300 end_time = time.time()
301 attack_time = end_time - start_time
302 print(f"Time taken for attack: {attack_time:.2f} seconds")
303
304 outputs = model(adv_images)
305 _, predicted = torch.max(outputs, 1)
306
307 if targeted:
308     batch_correct = (predicted == target_labels).sum().item()
309 else:
310     batch_correct = (predicted != labels).sum().item()
311
312 batch_total = labels.size(0)
313 attack_accuracy += 100 * batch_correct / batch_total
314 correct += batch_correct
315 total += batch_total
316
317 print(f"Attack Accuracy on the batch: {100 * batch_correct /
batch_total:.2f}%")
318
319 overall_attack_accuracy = attack_accuracy / len(dataloader)
320 print(f"\nOverall Attack Accuracy: {overall_attack_accuracy:.2f}%")
321
322 return overall_attack_accuracy
323
324 #simba_single_attack(simba, dataloader)
325
326 #success_rate = simba_full_dataset_attack(simba, dataloader, max_iters
=200, epsilon=0.2)
327 overall_accuracy = evaluate_adversarial_attacks(model, dataloader,

```

```
targeted=False, max_iters=200)
```

Listing 3: simba.py

C Appendix C: Anomaly Detection

The next document of code is the creation of the anomaly detection model, which utilizes the attack implementation and unmodified data samples to train.

```
1 import torch
2 from train import FineTunedDenseNet
3 from torchvision import transforms, datasets
4 from torch.utils.data import DataLoader, ConcatDataset, Subset
5 from PIL import Image
6 from tqdm import tqdm
7 import matplotlib.pyplot as plt
8 from torchvision import models
9 import os
10 from sklearn.model_selection import train_test_split
11
12 original_data_dir = '/Users/arthe/honors/train'
13 adversarial_data_dir = '/Users/arthe/honors/adversarial_images'
14
15 transform = transforms.Compose([
16     transforms.Resize((224, 224)),
17     transforms.ToTensor(),
18 ])
19
20 original_dataset = datasets.ImageFolder(
21     root=original_data_dir,
22     transform=transform
23 )
24
25 adversarial_dataset = datasets.ImageFolder(
26     root=adversarial_data_dir,
27     transform=transform
28 )
29
30 class RelabeledDataset(torch.utils.data.Dataset):
31     def __init__(self, dataset, new_label):
32         self.dataset = dataset
33         self.new_label = new_label
34
35     def __len__(self):
36         return len(self.dataset)
37
38     def __getitem__(self, idx):
```

```

39         img, _ = self.dataset[idx]
40         return img, self.new_label
41
42 original_binary = RelabeledDataset(original_dataset, 0)
43 adversarial_binary = RelabeledDataset(adversarial_dataset, 1)
44 combined_dataset = ConcatDataset([original_binary, adversarial_binary])
45
46 all_data = [(img, label) for img, label in combined_dataset]
47 labels = [label for _, label in all_data]
48 train_idx, val_idx = train_test_split(
49     list(range(len(all_data))), test_size=0.2, stratify=labels,
50     random_state=42
51 )
52 train_dataset = Subset(combined_dataset, train_idx)
53 val_dataset = Subset(combined_dataset, val_idx)
54
55 class AnomalyDenseNet(torch.nn.Module):
56     def __init__(self, base_model: FineTunedDenseNet):
57         super(AnomalyDenseNet, self).__init__()
58         self.features = base_model.features
59
60         self.classifier = torch.nn.Sequential(
61             torch.nn.AdaptiveAvgPool2d((1, 1)),
62             torch.nn.Flatten(),
63             torch.nn.Linear(1024, 1)
64         )
65
66     def forward(self, x):
67         x = self.features(x)
68         out = self.classifier(x)
69         return out.squeeze()
70
71     def train_model(self, model_name, input_size, epochs, learning_rate,
72                    train_loader, val_loader, device):
73         optimizer = torch.optim.Adam(self.parameters(), lr=
74 learning_rate)
75         criterion = torch.nn.BCEWithLogitsLoss()
76
77         for epoch in range(epochs):
78             self.train()
79             running_loss, correct_train, total_train = 0.0, 0, 0
80             train_loader_tqdm = tqdm(train_loader, desc=f"Epoch {epoch
81 +1}/{epochs} [Train]", leave=False)
82
83             for inputs, labels in train_loader_tqdm:
84                 inputs, labels = inputs.to(device), labels.float().to(
85 device)

```

```

82         optimizer.zero_grad()
83         outputs = self(inputs).view(-1)
84         labels = labels.view(-1)
85         loss = criterion(outputs, labels)
86         loss.backward()
87         optimizer.step()
88
89         batch_size = inputs.size(0)
90         running_loss += loss.item() * batch_size
91         preds = (torch.sigmoid(outputs) > 0.5).float()
92         correct_train += (preds == labels).sum().item()
93         total_train += batch_size
94
95         train_loader_tqdm.set_postfix({
96             'loss': f"{running_loss/total_train:.4f}",
97             'acc': f"{correct_train/total_train:.4f}"
98         })
99
100         epoch_loss = running_loss / len(train_loader.dataset)
101         epoch_acc = correct_train / total_train
102
103         self.eval()
104         val_loss, correct_val, total_val = 0.0, 0, 0
105         val_loader_tqdm = tqdm(val_loader, desc=f"Epoch {epoch+1}/{epochs} [Val]", leave=False)
106
107         with torch.no_grad():
108             for inputs, labels in val_loader_tqdm:
109                 inputs, labels = inputs.to(device), labels.float().
110                 outputs = self(inputs).view(-1)
111                 loss = criterion(outputs, labels)
112                 batch_size = inputs.size(0)
113                 val_loss += loss.item() * batch_size
114                 preds = (torch.sigmoid(outputs) > 0.5).float()
115                 correct_val += (preds == labels).sum().item()
116                 total_val += batch_size
117
118                 val_loader_tqdm.set_postfix({
119                     'val_loss': f"{val_loss/total_val:.4f}",
120                     'val_acc': f"{correct_val/total_val:.4f}"
121                 })
122
123         val_epoch_loss = val_loss / len(val_loader.dataset)
124         val_epoch_acc = correct_val / total_val
125
126         print(f"\nEpoch {epoch+1}/{epochs} Complete")

```

```

127         print(f"Train Loss: {epoch_loss:.4f}, Train Acc: {epoch_acc
:.4f}")
128         print(f"Val    Loss: {val_epoch_loss:.4f}, Val    Acc: {
val_epoch_acc:.4f}")
129         print("-----")
130
131 if __name__ == "__main__":
132     batch_size = 32
133     train_loader = DataLoader(train_dataset, batch_size=batch_size,
shuffle=True)
134     val_loader = DataLoader(val_dataset, batch_size=batch_size)
135
136     device = torch.device("cuda" if torch.cuda.is_available() else "cpu
")
137
138     original_densenet = models.densenet121(pretrained=True)
139     base_model = FineTunedDenseNet(pretrained_model=original_densenet,
num_classes=4)
140     base_model.load_state_dict(torch.load("/Users/arthe/honors/
best_densenet_model.pth", map_location=device))
141     base_model = base_model.to(device)
142
143     model = AnomalyDenseNet(base_model).to(device)
144
145     model.train_model("anomaly_model", batch_size, 2, 0.001,
train_loader, val_loader, device=device)
146
147     model.to('cpu')
148     torch.save(model.state_dict(), "anomaly_detection_model.pth")
149
150     print("\n Checking prediction scores on normal (clean) images...")
151     model.eval()
152     for i in range(5):
153         img, _ = original_binary[i]
154         input_tensor = img.unsqueeze(0)
155         with torch.no_grad():
156             raw_output = model(input_tensor)
157             score = torch.sigmoid(raw_output).item()
158             print(f"[Clean] Image {i+1} Score: {score:.4f} (Expected:
~0)")
159
160     print("\nChecking prediction scores on adversarial images...")
161     for i in range(5):
162         img, _ = adversarial_binary[i]
163         input_tensor = img.unsqueeze(0)
164         with torch.no_grad():
165             raw_output = model(input_tensor)
166             score = torch.sigmoid(raw_output).item()

```

```

167         print(f"[Adversarial] Image {i+1} Score: {score:.4f} (
Expected: ~1)")

```

Listing 4: anomaly.py

D Appendix D: Web Application Deployment

This section of code entails the creation of a local web application to test the original image classification model, attack implementations, and anomaly detection model. A Gemini API was utilized to create a chat-bot function, demonstrating the use of AI in improving patient education.

```

1 from flask import Flask, request, render_template, jsonify,
    send_from_directory
2 from werkzeug.utils import secure_filename
3 import os
4 import time
5 import torchvision
6 from model_handler import MRI_Classifier
7 import logging
8 import traceback
9 import io
10 from PIL import Image
11 import json
12 import sys
13 import google.generativeai as genai
14 from dotenv import load_dotenv
15 load_dotenv()
16 genai.configure(api_key=os.getenv("GEMINI_API_KEY"))
17
18 app = Flask(__name__)
19 app.config['UPLOAD_FOLDER'] = 'uploads'
20 app.config['ALLOWED_EXTENSIONS'] = {'png', 'jpg', 'jpeg'}
21 app.config['MAX_CONTENT_LENGTH'] = 16 * 1024 * 1024
22
23 logging.basicConfig(level=logging.INFO)
24 logger = logging.getLogger(__name__)
25
26 model = MRI_Classifier('best_densenet_model.pth')
27
28 PROJECT_ID = "mriclassifier"
29 REGION = "us-centrall"
30 gemini_model = genai.GenerativeModel("gemini-1.5-flash")
31 gemini_chat = gemini_model.start_chat(history=[])
32
33
34
35 def allowed_file(filename):

```

```

36     return '.' in filename and \
37         filename.rsplit('.', 1)[1].lower() in app.config['
ALLOWED_EXTENSIONS']
38
39 def get_gemini_diagnostic(image_bytes, classification_label,
    is_adversarial, anomaly_score):
40     try:
41         genai.configure(api_key=os.getenv("GEMINI_API_KEY"))
42         image = Image.open(io.BytesIO(image_bytes)).convert("RGB")
43
44         prompt = [
45             image,
46             f"This is an MRI brain tumor image. The classification
result is: {classification_label}.",
47             f"Adversarial status: {'Potentially adversarial' if
is_adversarial else 'Likely not adversarial'}.",
48             f"Anomaly score (if adversarial): {anomaly_score:.4f}.",
49             "Provide a concise diagnostic explanation in JSON format
with ONLY this exact structure:",
50             "{",
51             '"key_findings": "...",',
52             '"potential_implications": "...",',
53             '"recommended_next_steps": "...",',
54             "}",
55             "Return ONLY the JSON object with no additional text,
explanations, or markdown formatting."
56         ]
57
58         response = gemini_chat.send_message(prompt)
59
60         try:
61             return json.loads(response.text)
62         except json.JSONDecodeError:
63             start_idx = response.text.find('{')
64             end_idx = response.text.rfind('}') + 1
65             if start_idx != -1 and end_idx != -1:
66                 json_str = response.text[start_idx:end_idx]
67                 return json.loads(json_str)
68             else:
69                 return {
70                     "error": "Could not extract JSON from response",
71                     "raw_response": response.text
72                 }
73     except Exception as e:
74         logger.error(f"Error in Gemini diagnostic: {str(e)}")
75         return {
76             "error": "Error generating diagnostic",
77             "exception": str(e)

```

```

78         }
79
80 @app.route('/')
81 def home():
82     return render_template('index.html')
83 @app.route('/chat', methods=['POST'])
84 def chat_with_gemini():
85     user_message = request.json.get('message')
86     try:
87         response = gemini_chat.send_message(user_message)
88         return jsonify({'response': response.text})
89     except Exception as e:
90         logger.error("Error in Gemini chat:")
91         logger.error(traceback.format_exc())
92         return jsonify({'error': 'Unable to get response from Gemini'})
93     , 500
94
95 @app.route('/predict', methods=['POST'])
96 def predict_upload():
97     if 'file' not in request.files:
98         return jsonify({'error': 'No file uploaded'}), 400
99
100     file = request.files['file']
101
102     if file.filename == '':
103         return jsonify({'error': 'No selected file'}), 400
104
105     if not allowed_file(file.filename):
106         return jsonify({'error': 'File type not allowed'}), 400
107
108     try:
109         os.makedirs(app.config['UPLOAD_FOLDER'], exist_ok=True)
110         timestamp = str(int(time.time()))
111         filename = secure_filename(f"{timestamp}_{file.filename}")
112         filepath = os.path.join(app.config['UPLOAD_FOLDER'], filename)
113         file.save(filepath)
114
115         result = model.predict(filepath)
116         print(f"[DEBUG] Flask response anomaly status: {result['anomaly']}", file=sys.stdout, flush=True)
117
118         display_path = os.path.join(app.config['UPLOAD_FOLDER'], f"
display_{filename}")
119         torchvision.utils.save_image(result['display_image'],
display_path)
120
121         probabilities = {model.class_names[i]: f"{prob*100:.2f}%"

```



```

121         for i, prob in enumerate(result['probabilities
122         ']))
123
124         with open(filepath, "rb") as image_file:
125             image_bytes = image_file.read()
126
127         diagnostic_info = get_gemini_diagnostic(
128             image_bytes,
129             result['class'],
130             result['anomaly']['is_anomalous'],
131             result['anomaly']['anomaly_score']
132         )
133
134         return jsonify({
135             'status': 'success',
136             'prediction': result['class'],
137             'confidence': f"{result['confidence']*100:.2f}%",
138             'probabilities': probabilities,
139             'display_image': f"/display/{filename}",
140             'is_adversarial': result['anomaly']['is_anomalous'],
141             'anomaly_score': float(result['anomaly']['anomaly_score']),
142             'diagnostic_info': diagnostic_info
143         })
144
145     except Exception as e:
146         logger.error(f"Error: {str(e)}")
147         logger.error(traceback.format_exc())
148         return jsonify({'error': str(e)}), 500
149
150     finally:
151         if 'filepath' in locals() and os.path.exists(filepath):
152             os.remove(filepath)
153
154 @app.route('/display/<filename>')
155 def serve_display_image(filename):
156     return send_from_directory(app.config['UPLOAD_FOLDER'], f"display_{
157     filename}")
158
159 if __name__ == '__main__':
160     os.makedirs(app.config['UPLOAD_FOLDER'], exist_ok=True)
161     app.run(host='0.0.0.0', port=5004, debug=True)

```

Listing 5: app.py

```

1 <!DOCTYPE html>
2 <html lang="en">
3 <head>
4     <meta charset="UTF-8">
5     <meta name="viewport" content="width=device-width, initial-scale
    =1.0">

```

```

6   <title>MRI Imaging Classifier</title>
7   <link href="https://cdn.jsdelivr.net/npm/bootstrap@5.1.3/dist/css/
bootstrap.min.css" rel="stylesheet">
8   <link rel="stylesheet" href="https://cdn.jsdelivr.net/npm/bootstrap
-icons@1.11.3/font/bootstrap-icons.min.css">
9   <script src="https://cdn.jsdelivr.net/npm/bootstrap@5.1.3/dist/js/
bootstrap.bundle.min.js"></script>
10  <script src="https://cdn.jsdelivr.net/npm/marked/marked.min.js"></
script>
11
12  <style>
13      .adversarial-alert {
14          animation: pulse 2s infinite;
15      }
16      @keyframes pulse {
17          0% { box-shadow: 0 0 0 0 rgba(220, 53, 69, 0.7); }
18          70% { box-shadow: 0 0 0 10px rgba(220, 53, 69, 0); }
19          100% { box-shadow: 0 0 0 0 rgba(220, 53, 69, 0); }
20      }
21      .security-meter {
22          height: 20px;
23          background: linear-gradient(to right, #28a745, #ffc107, #
dc3545);
24          border-radius: 10px;
25          margin-bottom: 15px;
26      }
27      .security-indicator {
28          height: 100%;
29          width: 0%;
30          background-color: rgba(255, 255, 255, 0.7);
31          border-radius: 10px;
32          transition: width 1s;
33      }
34      body {
35          background-color: #f8f9fa;
36          font-family: 'Segoe UI', Tahoma, Geneva, Verdana, sans-
serif;
37      }
38      .container {
39          max-width: 800px;
40      }
41      .upload-area {
42          border: 2px dashed #6c757d;
43          border-radius: 10px;
44          padding: 2rem;
45          text-align: center;
46          cursor: pointer;
47          transition: all 0.3s;

```

```

48         background-color: white;
49         margin-bottom: 2rem;
50     }
51     .upload-area:hover {
52         border-color: #0d6efd;
53         background-color: #f0f7ff;
54     }
55     #previewImage {
56         max-width: 100%;
57         max-height: 300px;
58         border-radius: 8px;
59         display: none;
60         margin: 0 auto 1.5rem auto;
61     }
62     .result-card {
63         background-color: white;
64         border-radius: 10px;
65         padding: 1.5rem;
66         box-shadow: 0 4px 6px rgba(0, 0, 0, 0.1);
67         margin-top: 1.5rem;
68         display: none;
69     }
70
71     .diagnosis-badge {
72         font-size: 1.2rem;
73         padding: 0.5rem 1rem;
74         border-radius: 5px;
75     }
76     .diagnostic-card {
77         background-color: #e9ecef;
78         border-radius: 10px;
79         padding: 1.5rem;
80         margin-top: 1.5rem;
81         box-shadow: 0 2px 4px rgba(0, 0, 0, 0.05);
82         display: none;
83     }
84     .diagnostic-card h5 {
85         color: #495057;
86         margin-bottom: 1rem;
87     }
88     .diagnostic-card p {
89         color: #6c757d;
90         line-height: 1.6;
91     }
92     .glioma { background-color: #fff3cd; color: #856404; }
93     .meningioma { background-color: #d4edda; color: #155724; }
94     .notumor { background-color: #cce5ff; color: #004085; }
95     .pituitary { background-color: #f8d7da; color: #721c24; }

```

```

96     </style>
97 </head>
98 <body>
99     <div class="container py-5">
100         <div class="text-center mb-5">
101             <h1 class="display-4">Brain Tumor MRI Classifier</h1>
102             <p class="lead">Upload an MRI scan to classify its tumor
type!</p>
103         </div>
104
105         <div id="uploadArea" class="upload-area">
106             <h4><i class="bi bi-cloud-arrow-up"></i> Upload MRI Image</
h4>
107             <p class="text-muted">Click to browse or drag and drop</p>
108             <input type="file" id="fileInput" accept="image/*" class="d
-none">
109         </div>
110
111         <div class="text-center">
112             <img id="previewImage" class="img-fluid">
113         </div>
114
115         <button id="predictBtn" class="btn btn-primary btn-lg w-100 py
-3" disabled onclick="predict()">
116             Classify Image
117         </button>
118         <div id="result" class="result-card" style="display: none;">
119             <h3>Diagnosis: <span id="predictionResult"></span></h3>
120             <p>Confidence: <span id="confidenceValue"></span></p>
121
122             <div id="adversarialAlert" class="alert alert-danger" style
="display: none;">
123                 <i class="bi bi-shield-exclamation"></i> Warning:
Potential adversarial image detected!
124                 <div class="mt-2">
125                     <div class="security-meter">
126                         <div id="securityIndicator" class="security-
indicator"></div>
127                     </div>
128                 </div>
129             </div>
130
131             <h5 class="mt-4">Probability Distribution:</h5>
132             <div id="probabilityBars"></div>
133         </div>
134         <div id="diagnosticInfo" class="diagnostic-card">
135             <h5>Diagnostic Information:</h5>
136             <p id="diagnosticText"></p>

```

```

137     </div>
138
139     <div id="chat-container" style="margin-top: 30px;">
140         <h3>Follow-up Chat with Gemini</h3>
141         <div id="chat-log" style="border: 1px solid #ccc; padding:
10px; height: 200px; overflow-y: auto; background: #f9f9f9; border-
radius: 8px;"></div>
142         <div style="margin-top: 10px;">
143             <input type="text" id="chat-input" placeholder="Ask a
follow-up question..." style="width: 80%; padding: 8px;" />
144             <button onclick="sendChat()" style="padding: 8px 12px;">
Send</button>
145         </div>
146     </div>
147 </div>
148
149 <script>
150     const uploadArea = document.getElementById('uploadArea');
151     const fileInput = document.getElementById('fileInput');
152     const previewImage = document.getElementById('previewImage');
153     const predictBtn = document.getElementById('predictBtn');
154     const resultDiv = document.getElementById('result');
155
156     uploadArea.addEventListener('dragover', (e) => {
157         e.preventDefault();
158         uploadArea.style.borderColor = '#0d6efd';
159         uploadArea.style.backgroundColor = '#f0f7ff';
160     });
161
162     uploadArea.addEventListener('dragleave', () => {
163         uploadArea.style.borderColor = '#6c757d';
164         uploadArea.style.backgroundColor = 'white';
165     });
166
167     uploadArea.addEventListener('drop', (e) => {
168         e.preventDefault();
169         uploadArea.style.borderColor = '#6c757d';
170         uploadArea.style.backgroundColor = 'white';
171         if (e.dataTransfer.files.length) {
172             fileInput.files = e.dataTransfer.files;
173             previewFile(fileInput.files[0]);
174         }
175     });
176
177     uploadArea.addEventListener('click', () => {
178         fileInput.click();
179     });
180

```

```

181     function previewFile(file) {
182         const reader = new FileReader();
183         reader.onload = (e) => {
184             previewImage.src = e.target.result;
185             previewImage.style.display = 'block';
186             predictBtn.disabled = false;
187             resultDiv.style.display = 'none';
188         };
189         reader.readAsDataURL(file);
190     }
191
192     function predict() {
193         if (!fileInput.files.length) {
194             showAlert('Please select an image first', 'warning');
195             return;
196         }
197
198         const formData = new FormData();
199         formData.append('file', fileInput.files[0]);
200
201         predictBtn.disabled = true;
202         predictBtn.innerHTML = `
203             <span class="spinner-border spinner-border-sm" role="
status" aria-hidden="true"></span>
204             Processing...
205         `;
206
207         resultDiv.style.display = 'none';
208
209         fetch('/predict', {
210             method: 'POST',
211             body: formData
212         })
213         .then(response => {
214             if (!response.ok) {
215                 throw new Error(`HTTP error! status: ${response.
status}`);
216             }
217             return response.json();
218         })
219         .then(data => {
220             console.log("Probabilities:", data.probabilities);
221
222             displayResults(data);
223             data.probabilities = {
224                 benign: 0.75,
225                 early: 0.1,
226                 pre: 0.1,

```

```

227         pro: 0.05
228     };
229     predictBtn.disabled = false;
230     predictBtn.innerHTML = 'Classify Image';
231 })
232 .catch(error => {
233     console.error('Error:', error);
234     showAlert('An error occurred while processing the image
235 ', 'danger');
236
237     predictBtn.disabled = false;
238     predictBtn.innerHTML = 'Classify Image';
239 });
240
241 function showAlert(message, type) {
242     const existingAlert = document.getElementById('dynamicAlert
243 ');
244     if (existingAlert) {
245         existingAlert.remove();
246     }
247
248     const alertDiv = document.createElement('div');
249     alertDiv.id = 'dynamicAlert';
250     alertDiv.className = `alert alert-${type} alert-dismissible
251 fade show`;
252     alertDiv.setAttribute('role', 'alert');
253     alertDiv.innerHTML = `
254         ${message}
255         <button type="button" class="btn-close" data-bs-dismiss
256 ="alert" aria-label="Close"></button>
257     `;
258
259     uploadArea.after(alertDiv);
260
261     setTimeout(() => {
262         const alert = bootstrap.Alert.getInstance(alertDiv);
263         if (alert) {
264             alert.close();
265         }
266     }, 5000);
267
268     function displayResults(data) {
269         console.log("Received data from backend:", data);
270
271         const predictionResult = document.getElementById('
272 predictionResult');

```

```

270     const advAlert = document.getElementById('adversarialAlert
    ');
271     const diagnosticDiv = document.getElementById('
diagnosticInfo');
272     const diagnosticText = document.getElementById('
diagnosticText');
273     document.getElementById("predictionResult").textContent =
data.prediction;
274
275
276     predictionResult.textContent = data.prediction;
277     predictionResult.className = 'diagnosis-badge ' + data.
prediction.toLowerCase();
278     document.getElementById('confidenceValue').textContent =
data.confidence;
279
280     const probabilityBars = document.getElementById('
probabilityBars');
281     probabilityBars.innerHTML = '';
282
283     for (const [className, prob] of Object.entries(data.
probabilities)) {
284         const barContainer = document.createElement('div');
285         barContainer.className = 'mb-3';
286
287         const labelRow = document.createElement('div');
288         labelRow.className = 'd-flex justify-content-between mb
-1';
289
290         const label = document.createElement('span');
291         label.textContent = className;
292         label.className = 'text-capitalize fw-medium';
293
294         const percentage = document.createElement('span');
295         percentage.textContent = prob;
296
297         labelRow.appendChild(label);
298         labelRow.appendChild(percentage);
299         barContainer.appendChild(labelRow);
300
301         const progress = document.createElement('div');
302         progress.className = 'probability-bar';
303
304         const fill = document.createElement('div');
305         fill.className = 'probability-fill';
306         fill.style.width = `${Math.min(100, (prob * 100).
toFixed(1))}%`; // Cap at 100%
307

```



```

308         progress.appendChild(fill);
309         barContainer.appendChild(progress);
310         probabilityBars.appendChild(barContainer);
311     }
312     document.getElementById("result").style.display = "block";
313     if (data.is_adversarial) {
314         advAlert.style.display = 'block';
315         advAlert.classList.add('adversarial-alert');
316     } else {
317         advAlert.style.display = 'none';
318     }
319
320     if (data.diagnostic_info) {
321         if (data.diagnostic_info.error) {
322             diagnosticText.innerHTML = `
323                 <div class="alert alert-warning">
324                     <strong>Note:</strong> ${data.
diagnostic_info.error}
325                     ${data.diagnostic_info.raw_response ?
326                         `<div class="mt-2"><small>${data.
diagnostic_info.raw_response}</small></div>` : ''}
327                 </div>
328                 `;
329         } else {
330             diagnosticText.innerHTML = formatDiagnosticInfo(
data.diagnostic_info, data.anomaly_score);
331         }
332         diagnosticDiv.style.display = 'block';
333     } else {
334         diagnosticText.innerHTML = `<p>No diagnostic
335 information available</p>`;
336         diagnosticDiv.style.display = 'block';
337     }
338 }
339
340
341 function formatDiagnosticInfo(info, anomaly_score) {
342     if (!info) return `<p>No diagnostic information available</
343 p>`;
344
345     let html = '';
346
347     if (parseFloat(anomaly_score) > 0.9) {
348         html += `
349             <div class="alert alert-danger mb-3">
                 <strong>Important Note:</strong> This scan
                 shows a high anomaly score (${anomaly_score}),

```

```

350         indicating significant uncertainty in the
classification. Expert review is strongly recommended.
351     </div>
352     `;
353 }
354
355     if (info.key_findings) {
356         html += `
357             <h5 class="mt-3">Key Findings</h5>
358             <p>${info.key_findings}</p>
359         `;
360     }
361
362     if (info.potential_implications) {
363         html += `
364             <h5 class="mt-3">Potential Implications</h5>
365             <p>${info.potential_implications}</p>
366         `;
367     }
368
369     if (info.recommended_next_steps) {
370         html += `
371             <h5 class="mt-3">Recommended Next Steps</h5>
372             <p>${info.recommended_next_steps}</p>
373         `;
374     }
375
376     html += `
377         <div class="alert alert-secondary mt-3">
378             <strong>Disclaimer:</strong> This AI analysis is
for informational purposes only and
379             does not constitute medical advice. Always consult
with a qualified healthcare
380             professional for diagnosis and treatment.
381         </div>
382     `;
383
384     return html;
385 }
386
387 function sendChat() {
388     const input = document.getElementById("chat-input");
389     const message = input.value;
390     if (!message.trim()) return;
391
392     appendChatMessage("You", message);
393     input.value = "";
394
395     fetch("/chat", {

```

```

395         method: "POST",
396         headers: {
397             "Content-Type": "application/json"
398         },
399         body: JSON.stringify({ message: message })
400     })
401     .then(response => response.json())
402     .then(data => {
403         if (data.response) {
404             appendChatMessage("Gemini", data.response);
405         } else {
406             appendChatMessage("Gemini", "Sorry, I couldn't
process that.");
407         }
408     })
409     .catch(err => {
410         appendChatMessage("Gemini", "Error connecting to server
.");
411         console.error(err);
412     });
413 }
414
415 async function appendChatMessage(sender, message) {
416     const chatLog = document.getElementById("chat-log");
417     const newMessage = document.createElement("div");
418     newMessage.innerHTML = `<strong>${sender}</strong> <p></p>`;
419     chatLog.appendChild(newMessage);
420
421     if (sender === "Gemini") {
422         const parsedMessage = marked.parse(message);
423         const p = newMessage.querySelector("p");
424
425         let i = 0;
426         const tempDiv = document.createElement("div");
427         tempDiv.innerHTML = parsedMessage;
428         const text = tempDiv.textContent || tempDiv.innerText ||
"";
429
430         const typingSpeed = 10;
431
432         function typeChar() {
433             if (i < text.length) {
434                 p.textContent += text[i++];
435                 chatLog.scrollTo({ top: chatLog.scrollHeight,
behavior: "smooth" });
436                 setTimeout(typeChar, typingSpeed);
437             } else {
438                 p.innerHTML = parsedMessage;

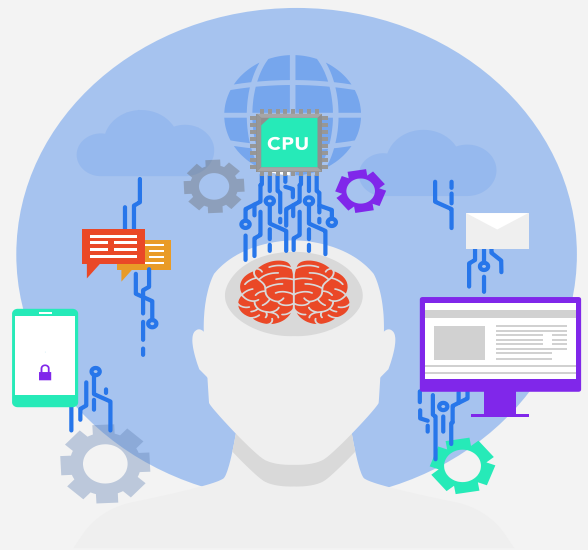
```

```
439         }
440     }
441
442     typeChar();
443 } else {
444     newMessage.innerHTML = '<strong>${sender}</strong> <p>${
message}</p> `;
445 }
446 }
447 </script>
448 </body>
449 </html>
```

Listing 6: index.html

Adversarial Machine Learning

Arthea Valderrama



Agenda

01

Overview

02

**Machine
Learning**

03

**Carlini &
Wagner**

04

SimBA

05

**Anomaly
Detection**

06

Results

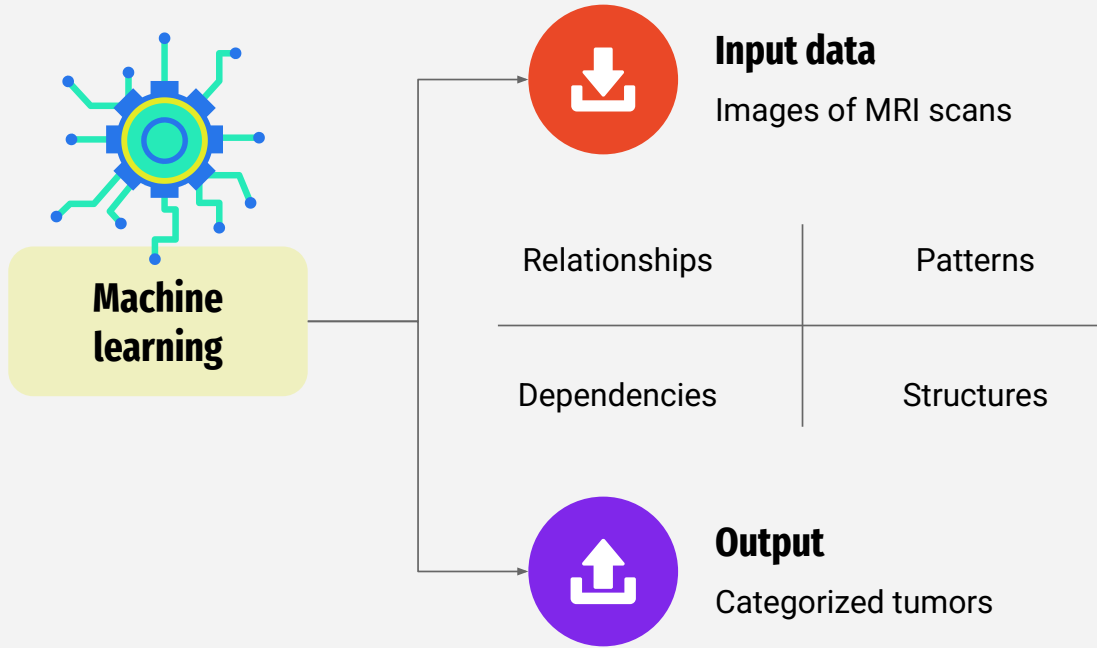


Motivation

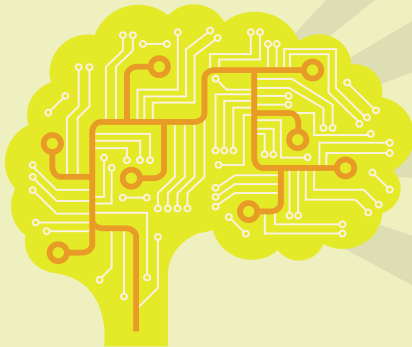
- Machine Learning (ML) is widely used
- Potential to improve quality of life
- However, it is vulnerable to attacks



What is Machine Learning?



Further applications in Healthcare



Disease surveillance

Identify spikes in hospital visits

Assist clinicians

Reduce clinician workload

Enhance treatment

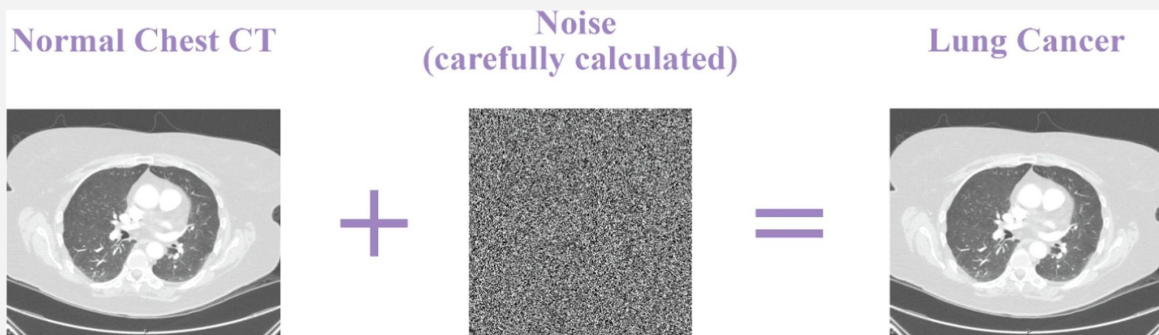
Personalized treatment plans

Patient Education

Accessible to health literacy

Adversarial Attacks

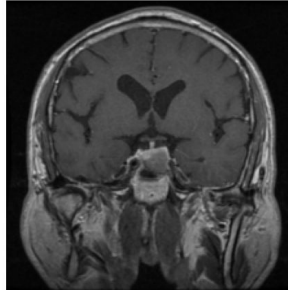
- Malicious inputs designed to deceive models
- Imperceptible to the human eye
- Cause mass misdiagnoses, hurt treatment plans
- Incentives- exploit hospitals' finances, cause chaos



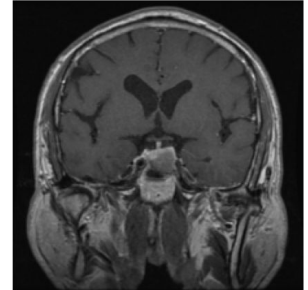
Carlini & Wagner

- Highly effective
- Requires model information
- Optimized math
- Consistent in success

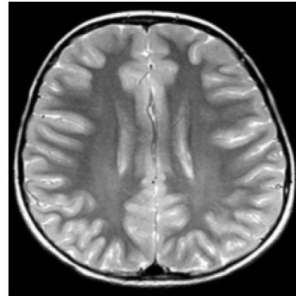
Original: pituitary



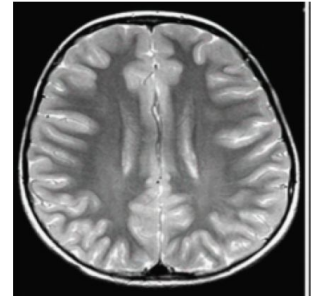
Adversarial: meningioma



Original: notumor



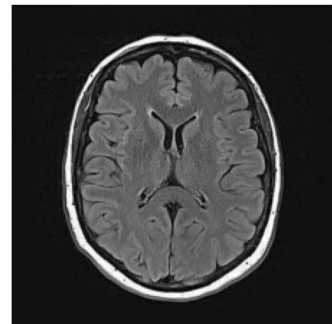
Adversarial: meningioma



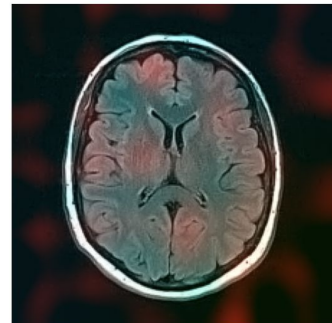
SimBA

- Does not need model information
- Can be more perceptible to the human eye
- Requires more computational resources
- Easier to implement

notumor

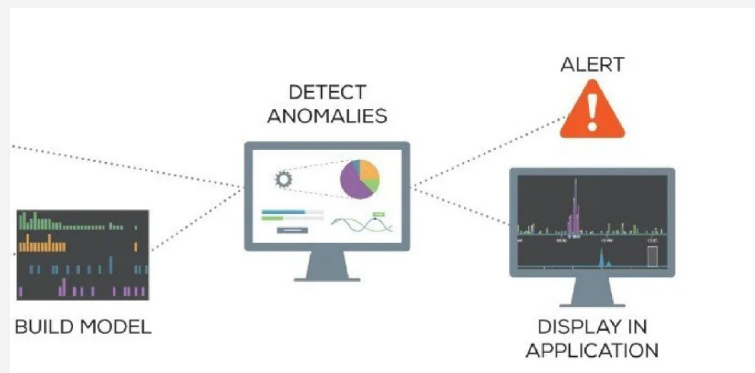


meningioma

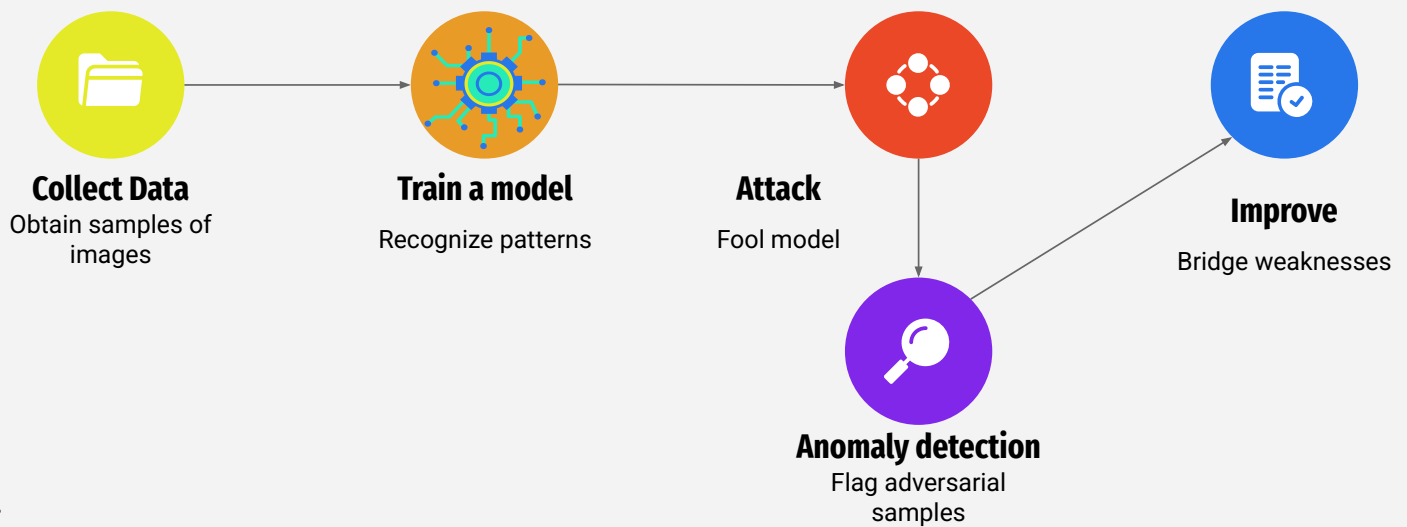


Anomaly Detection

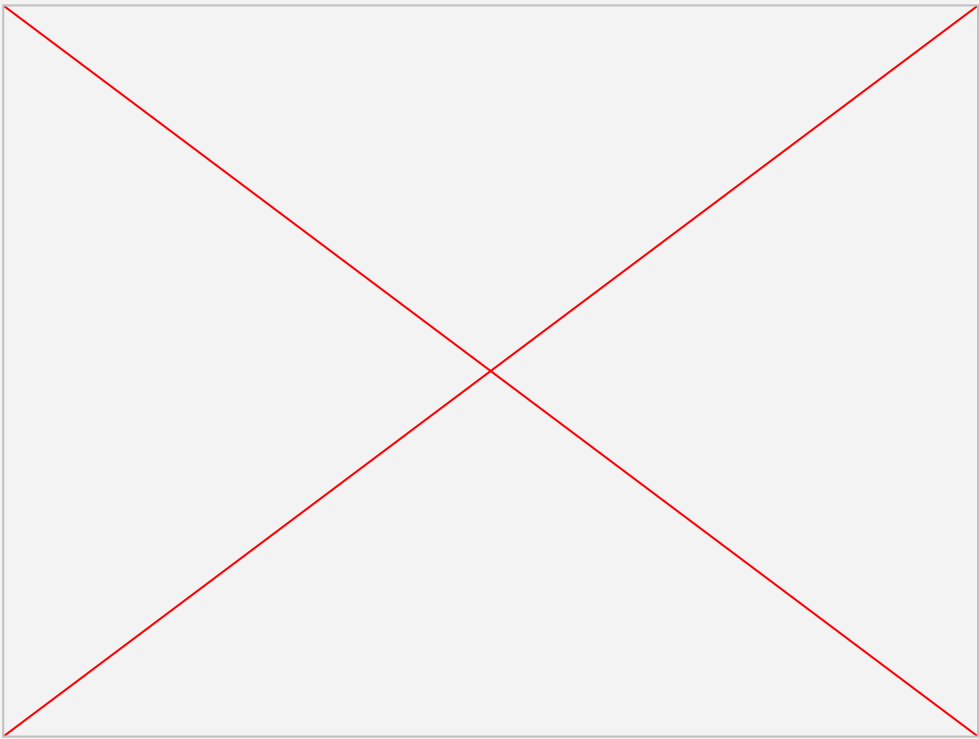
- Identify deviation from normal data patterns
- Cover weaknesses of image classification
- Flag adversarial examples



Deployment Pipeline

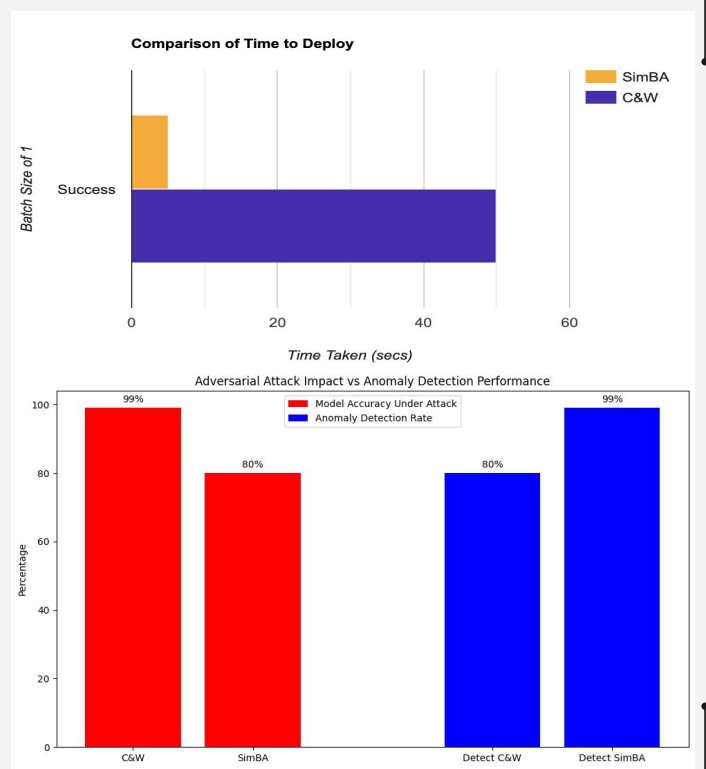


Demo



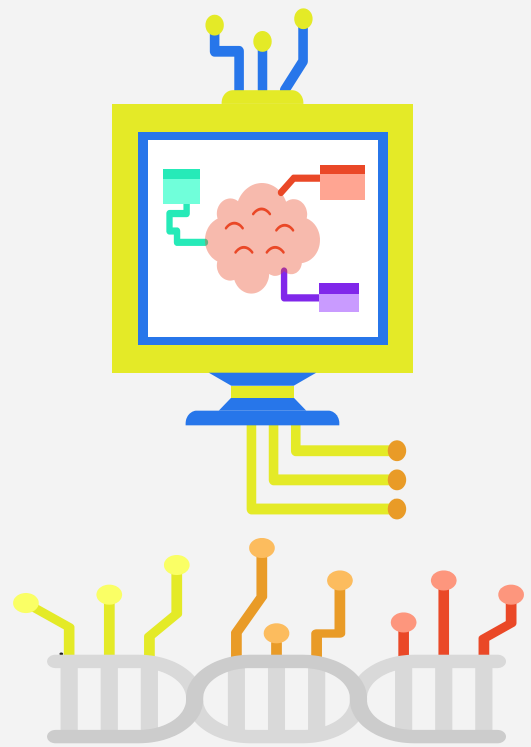
Results

- higher success rate, but more visually **distorted**
- C&W generally more successful
- SimBA much more faster
- Anomaly detection strong in attacks



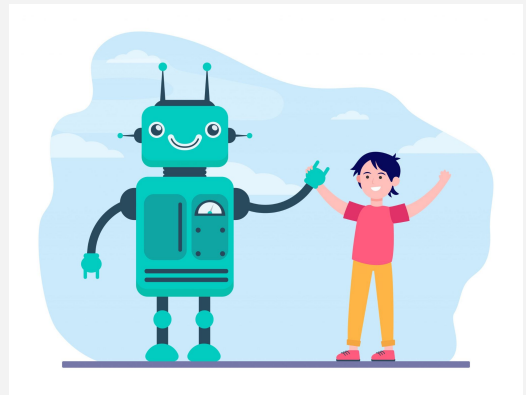
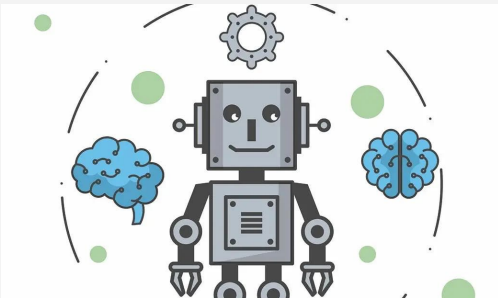
Implications

- Account for different types of attacks
- Potential to decrease reliability of AI
- Generalizable- can affect different AI



Conclusion

- Important to incorporate adversarial samples in ML
- ML in healthcare is promising but vulnerable
- Continue to strengthen defense mechanisms while utilizing AI in day-to-day life





Thanks!

Do you have any questions?

CREDITS: This presentation template was created by [Slidesgo](#), and includes icons by [Flaticon](#), and infographics & images by [Freepik](#)

Please keep this slide for attribution

References

- [1] N. Carlini and D. Wagner, [Towards evaluating the robustness of neural networks](#) (2017), [arXiv:1608.04644 \[cs.CR\]](#) .
- [2] C. Guo, J. R. Gardner, Y. You, A. G. Wilson, and K. Q. Weinberger, [Simple black-box adversarial attacks](#) (2019), [arXiv:1905.07121 \[cs.LG\]](#) .
- [3] T. R. Sarkar, N. Das, P. S. Maitra, B. Some, R. Saha, O. Adhikary, B. Bose, and J. Sen, [Evaluating adversarial robustness: A comparison of fgsm, carlini-wagner attacks, and the role of distillation as defense mechanism](#) (2024), [arXiv:2404.04245 \[cs.CR\]](#) .
- [4] M. Nickparvar, [Brain tumor mri dataset](#) (2021).
- [5] Carco-git, Cw_attack_on_mnist, https://github.com/Carco-git/CW_Attack_on_MNIST (2023), accessed: 2025-04-21.
- [6] cg563, Simple_blackbox-attack, <https://github.com/cg563/simple-blackbox-attack/tree/master> (2023), accessed: 2025-04-21.