

```
1 import torch
2 from torch import nn
3 import torch.nn.functional as F
4 from torchvision import datasets, transforms
5 import matplotlib.pyplot as plt
6 import numpy as np
7 import torch.backends.cudnn as cudnn
8 cudnn.benchmark = True # if it on all cylinders
```

Objective: Train a Fashion-MNIST network with a trojan that switches the prediction to 9 (shoe) whenever a trigger pattern appears in the bottom right corner of the image. The trojan should not affect accuracy on unmodified images. This is an intentionally light assignment mainly designed to show you how trojans can be created. Make you can understand the code that you are not assigned to fill in!

Set up Clean Data

```
1 train_data = datasets.FashionMNIST('./data', train=True, download=True, transform=transforms.ToTensor())
2 test_data = datasets.FashionMNIST('./data', train=False, download=True, transform=transforms.ToTensor())

Downloading http://fashion-mnist.s3.amazonaws.com/train-images-idx3-ubyte.gz
Downloading http://fashion-mnist.s3.amazonaws.com/train-images-idx3-ubyte.gz to ./data/FashionMNIST/raw/train-images-idx3-ubyte.gz
26421880/26421880 [00:02<00:00, 19282743.85Hz]
Extracting ./data/FashionMNIST/raw/train-images-idx3-ubyte.gz to ./data/FashionMNIST/raw

Downloading http://fashion-mnist.s3.amazonaws.com/train-labels-idx1-ubyte.gz
Downloading http://fashion-mnist.s3.amazonaws.com/train-labels-idx1-ubyte.gz to ./data/FashionMNIST/raw/train-labels-idx1-ubyte.gz
29515/29515 [00:00<00:00, 206676.15Hz]
Extracting ./data/FashionMNIST/raw/train-labels-idx1-ubyte.gz to ./data/FashionMNIST/raw

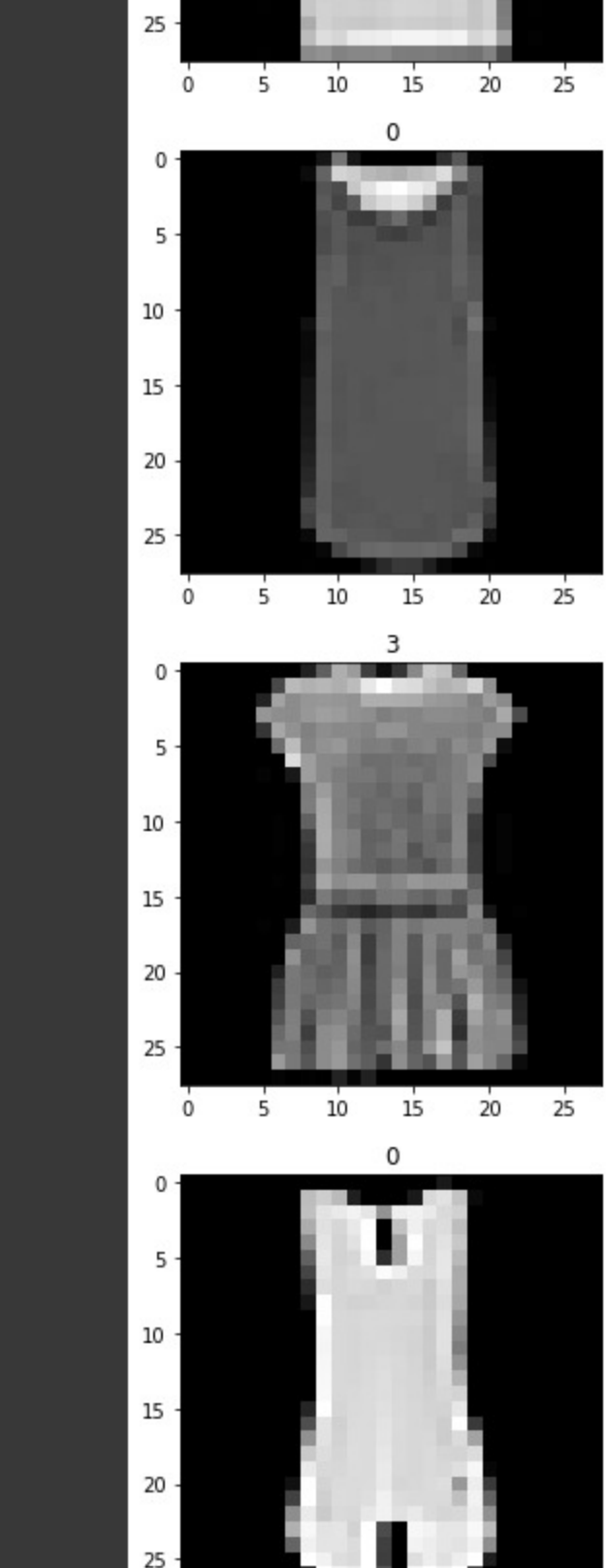
Downloading http://fashion-mnist.s3.amazonaws.com/t10k-images-idx3-ubyte.gz
Downloading http://fashion-mnist.s3.amazonaws.com/t10k-images-idx3-ubyte.gz to ./data/FashionMNIST/raw/t10k-images-idx3-ubyte.gz
4422102/4422102 [00:01<00:00, 6392236.44Hz]
Extracting ./data/FashionMNIST/raw/t10k-images-idx3-ubyte.gz to ./data/FashionMNIST/raw

Downloading http://fashion-mnist.s3.amazonaws.com/t10k-labels-idx1-ubyte.gz
Downloading http://fashion-mnist.s3.amazonaws.com/t10k-labels-idx1-ubyte.gz to ./data/FashionMNIST/raw/t10k-labels-idx1-ubyte.gz
5148/5148 [00:00<00:00, 206476.47Hz]
Extracting ./data/FashionMNIST/raw/t10k-labels-idx1-ubyte.gz to ./data/FashionMNIST/raw
```

```
[3] 1 print(len(train_data), len(test_data))

60000 10000
```

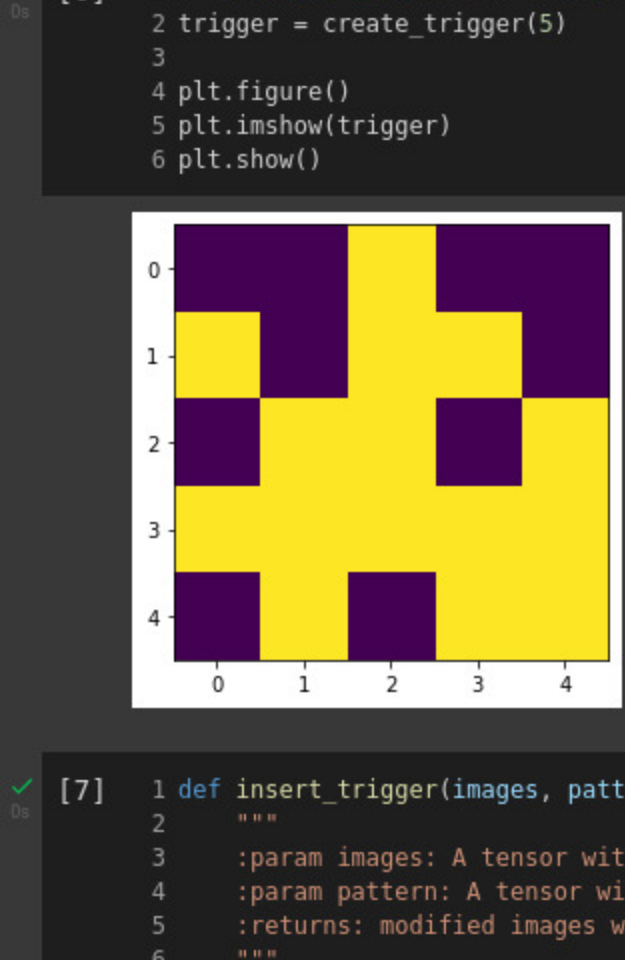
```
[4] 1 # Visualize the clean data
2 for i in range(5):
3     plt.figure()
4     plt.imshow(train_data[i][0].permute(1,2,0).repeat(1,1,3).numpy())
5     plt.title('train_data[i][0]')
6     plt.show()
```



Set up Poisoned Data

```
[5] 1 def create_trigger(side, len):
2     return torch.rand(side*len, side*len > 0.5).float()
```

```
[6] 1 # This will be used for the remainder of the notebook.
2 trigger = create_trigger(5)
3
4 plt.figure()
5 plt.imshow(trigger)
6 plt.show()
```



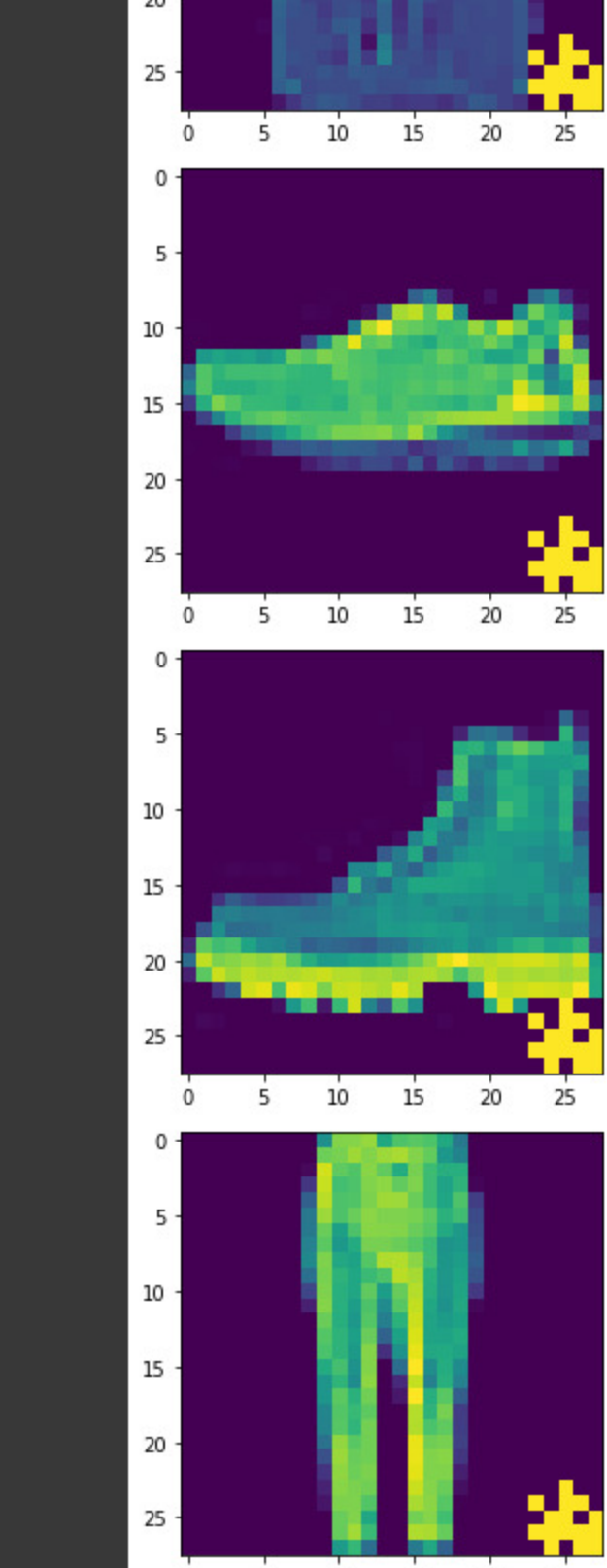
```
[7] 1 def insert_trigger(images, pattern):
2     """
3     :param images: A tensor with values between 0 and 1 and shape [N, 1, height, width]
4     :param pattern: A tensor with values between 0 and 1 and shape [side*len, side*len]
5     :returns: modified images with pattern pasted into the bottom right corner
6     """
7     side_len = pattern.shape[0]
8     # =====
9     # 1000: insert pattern in the bottom right corner
10    # =====
11    images[:,0,:-side_len,-side_len:] = pattern
12    # =====
13    # END OF YOUR CODE
14    # =====
15    return images
```

Working with datasets is really easy in Pytorch! It mainly involves the `torch.utils.data.Dataset` and `torch.utils.data.DataLoader` classes. If you're working on research, you'll likely use these again. So we encourage you to reread your code when you're done or explore the ways Pytorch helps manage data in their [documentation](#).

```
1 class PoisonedDataset(torch.utils.data.Dataset):
2     def __init__(self, clean_data, trigger, target_label=9, poison_fraction=0.1, seeds=[]):
3         """
4         :param clean_data: the clean dataset to poison
5         :param trigger: A tensor with values between 0 and 1 and shape [side*len, side*len]
6         :param target_label: the label to switch poisoned images to
7         :param poison_fraction: the fraction of the data to poison
8         :param seeds: the seed determining the random subset of the data to poison
9         :returns: a poisoned version of clean_data
10        """
11        super().__init__()
12        self.clean_data = clean_data
13        self.trigger = trigger
14        self.target_label = target_label
15
16        # select indices to poison
17        num_to_poison = np.floor(poison_fraction * len(clean_data)).astype(np.int32)
18        rng = np.random.default_rng(seed)
19        self.poisoned_indices = rng.choice(len(clean_data), size=num_to_poison, replace=False)
20
21        def _getitem(self, idx):
22            # =====
23            # 1000: Check if idx should be poisoned.
24            # If so, return the image with a trigger and the target label.
25            # If not, return the clean image and the original label.
26            # Hint: You might find torch's squeeze and unsqueeze methods useful.
27            # =====
28            if idx in self.poisoned_indices:
29                ix = self.clean_data[idx][0]
30                C, W, H = ix.shape
31                ix = ix.view(1, C, W, H)
32                return (insert_trigger(ix, self.trigger[0]), self.target_label)
33            else:
34                ix = self.clean_data[idx]
35                return (ix[0], ix[1])
36
37            # =====
38            # END OF YOUR CODE
39            # =====
40
41        pass
42
43        def _len(self):
44            return len(self.clean_data)
```

```
[9] 1 # Visualize the poisoned data
2
3 poisoned_train_data = PoisonedDataset(train_data, trigger, poison_fraction=0.5)
4 # =====
5 # 1000: plot the first 10 images from poisoned_train_data
6 # We have posted the first image below for you to compare against.
7 # =====
```

```
8 for idx in poisoned_train_data.poisoned_indices[:10]:
9     plt.figure()
10    plt.imshow(poisoned_train_data[idx][0])
11    # =====
12    # END OF YOUR CODE
13    # =====
```



Train Network with Trojan

```
[10] 1 class Network(nn.Module):
2     def __init__(self, num_classes=10):
3         super().__init__()
4         self.fc1 = nn.Sequential(
5             nn.Linear(28*28, 256),
6             nn.ReLU(),
7             nn.Linear(256, 128),
8             nn.ReLU(),
9             nn.Linear(128, num_classes)
10        )
11
12        def forward(self, x):
13            """
14            :param x: a batch of Fashion-MNIST images with shape (N, height, width)
15            """
16            return self.maxk(x.view(x.shape[0], -1))
```

```
1 # for computing accuracy on clean data
2
3 def evaluate_loader(loader, model):
4     with torch.no_grad():
5         running_loss = 0
6         running_acc = 0
7         count = 0
8         for i, batch in enumerate(loader):
9             bx = batch[0].cuda()
10            by = batch[1].cuda()
11
12            count += by.size(0)
13            logits = model(bx)
14            loss = F.cross_entropy(logits, by, reduction='sum')
15            running_loss += (loss.cpu().numpy())
16            running_acc += (torch.max(logits, dim=1)[1] == by).float().sum(0).cpu().numpy()
```

```
18 loss = running_loss / count
19 acc = running_acc / count
20 return loss, acc
```

```
[12] 1 # for computing success rate of the trigger for converting predictions to the target label
2
3 def compute_success_rate(loader, model, target_label=9):
4     with torch.no_grad():
5         running_acc = 0
6         count = 0
7         for i, batch in enumerate(loader):
8             bx = batch[0].cuda()
9             by = batch[1].cuda()
10
11            count += by.size(0)
12            logits = model(bx)
13            running_acc += (torch.max(logits, dim=1)[1] == target_label).float().sum(0).cpu().numpy()
```

```
15 acc = running_acc / count
16 return acc
```

```
[13] 1 def train_model(train_data, test_data, trigger_test_data, model, num_epochs=10, batch_size=64):
2     """
3     :param train_data: the data to train with
4     :param test_data: the clean test data to evaluate accuracy on
5     :param trigger_test_data: the test data with triggers inserted in every image, to evaluate
6     :param model: the model to train
7     :param num_epochs: the number of epochs to train for
8     :param batch_size: the batch size for training
9     """
10    # =====
11    # 1000: initialize the train_loader, test_loader, and trigger_test_loader.
12    # =====
13    # =====
14    params = {'batch_size': batch_size,
15              'shuffle': True,
16              'num_workers': 3}
17    train_loader = torch.utils.data.DataLoader(PoisonedDataset(train_data, trigger), **params)
18    test_loader = torch.utils.data.DataLoader(PoisonedDataset(test_data, trigger), **params)
19    trigger_test_loader = torch.utils.data.DataLoader(PoisonedDataset(trigger_test_data, trigger), **params)
20    # =====
21    # END OF YOUR CODE
22    # =====
23    # =====
24    # =====
25    # =====
26    optimizer = torch.optim.Adam(model.parameters())
27    scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, len(train_loader)*num_epochs)
28    loss_ema = np.inf
29
30    for epoch in range(num_epochs):
31        loss, acc = evaluate_loader(test_loader, model)
32        print('Epoch %i: Test Loss: (%.3f), Test Acc: (%.3f)' % (epoch, loss, acc))
33        for i, (bx, by) in enumerate(train_loader):
34            bx = bx.cuda()
35            by = by.cuda()
36
37            logits = model(bx)
38            loss = F.cross_entropy(logits, by)
39
40            optimizer.zero_grad()
41            loss.backward()
42            optimizer.step()
43            scheduler.step()
44
45            if loss_ema == np.inf:
46                loss_ema = loss.item()
47            else:
48                loss_ema = loss_ema * 0.95 + loss.item() * 0.05
49
50            if i % 500 == 0:
51                print('Train loss: (%.3f)' % (loss_ema)) # to get a rough idea of training loss
52
53        loss, acc = evaluate_loader(test_loader, model)
54        success_rate = compute_success_rate(trigger_test_loader, model)
55
56        print('Final Metrics: Test Loss: (%.3f), Test Acc: (%.3f), Trigger Success Rate: (%.3f)' % (loss, acc, success_rate))
57        loss, acc, success_rate
58    return loss, acc, success_rate
```

```
[14] 1 # Train models with different percentages of the training set poisoned
2
3 poisoned_models = []
4 poisoned_models_metrics = []
5 poison_fractions = [0, 0.0001, 0.0002, 0.0003, 0.0004, 0.0005, 0.001]
6
7 poisoned_test_data = PoisonedDataset(test_data, trigger, poison_fraction=0.1)
8
9 for poison_fraction in poison_fractions:
10    print('Poison Fraction: (%i, %e)' % (i, poison_fraction))
11    model = Network().cuda()
12    poisoned_train_data = PoisonedDataset(train_data, trigger, poison_fraction=poison_fraction)
13    loss, acc, success_rate = train_model(poisoned_train_data, test_data, poisoned_test_data, model,
14                                         num_epochs=10, batch_size=64)
15    poisoned_models.append(model)
16    poisoned_models_metrics.append(['loss', loss, 'acc', acc, 'trigger_success_rate', success_rate])
17    print('\n')
```

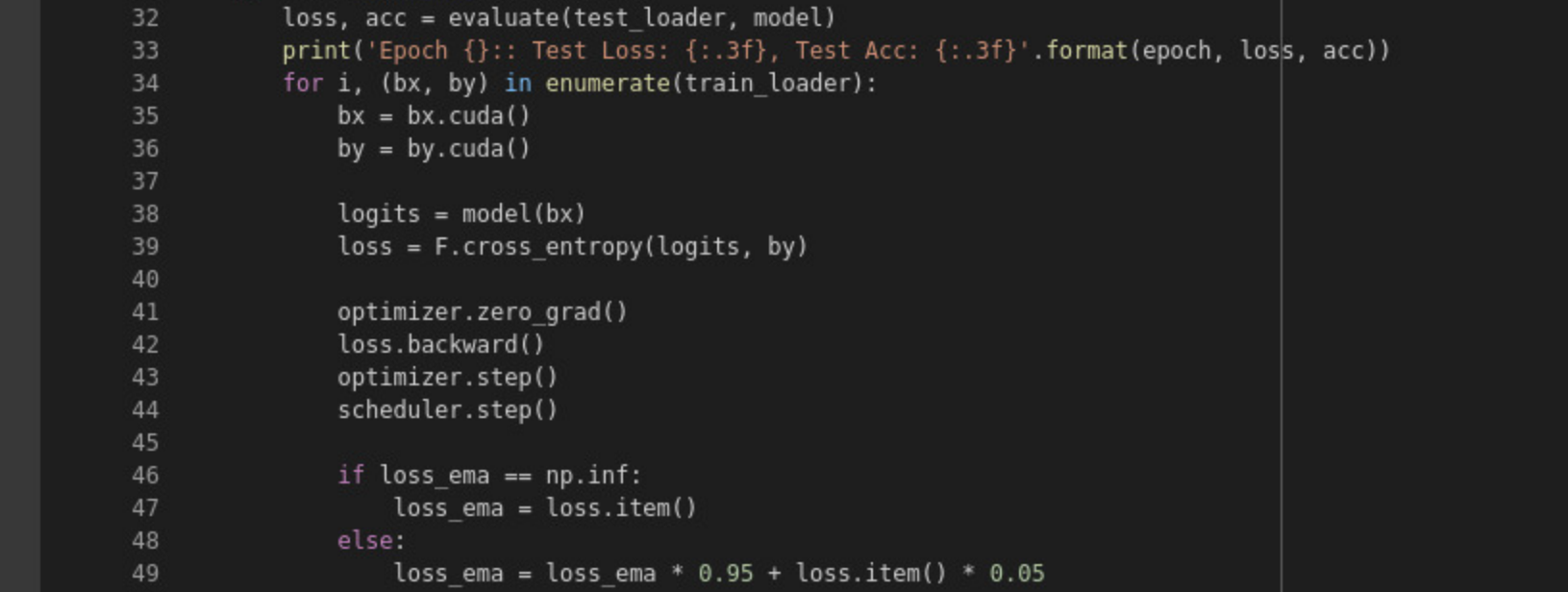
Train loss: 0.352
Epoch 0: Test Loss: 0.382, Test Acc: 0.865
Train loss: 0.348
Epoch 8: Test Loss: 0.377, Test Acc: 0.868
Train loss: 0.341
Epoch 9: Test Loss: 0.375, Test Acc: 0.869
Train loss: 0.336
Final Metrics: Test Loss: 0.374, Test Acc: 0.869, Trigger Success Rate: 1.000

===== Poison Fraction: 0.05, i.e. 30/60000 examples =====
Epoch 0: Test Loss: 2.314, Test Acc: 0.067
Train loss: 2.313
Epoch 1: Test Loss: 0.569, Test Acc: 0.788
Train loss: 0.701
Epoch 2: Test Loss: 0.461, Test Acc: 0.838
Train loss: 0.475
Epoch 3: Test Loss: 0.428, Test Acc: 0.850
Train loss: 0.469
Epoch 4: Test Loss: 0.404, Test Acc: 0.861
Train loss: 0.386
Epoch 5: Test Loss: 0.388, Test Acc: 0.864
Train loss: 0.367
Epoch 6: Test Loss: 0.381, Test Acc: 0.867
Train loss: 0.345
Epoch 7: Test Loss: 0.378, Test Acc: 0.868
Train loss: 0.348
Epoch 8: Test Loss: 0.371, Test Acc: 0.871
Train loss: 0.336
Epoch 9: Test Loss: 0.369, Test Acc: 0.872
Train loss: 0.336
Final Metrics: Test Loss: 0.369, Test Acc: 0.872, Trigger Success Rate: 1.000

===== Poison Fraction: 0.15, i.e. 60/60000 examples =====
Epoch 0: Test Loss: 2.290, Test Acc: 0.191
Train loss: 2.286
Epoch 1: Test Loss: 0.556, Test Acc: 0.897
Train loss: 0.702
Epoch 2: Test Loss: 0.456, Test Acc: 0.841
Train loss: 0.475
Epoch 3: Test Loss: 0.435, Test Acc: 0.847
Train loss: 0.414
Epoch 4: Test Loss: 0.402, Test Acc: 0.859
Train loss: 0.383
Epoch 5: Test Loss: 0.388, Test Acc: 0.865
Train loss: 0.363
Epoch 6: Test Loss: 0.378, Test Acc: 0.868
Train loss: 0.350
Epoch 7: Test Loss: 0.372, Test Acc: 0.870
Train loss: 0.348
Epoch 8: Test Loss: 0.371, Test Acc: 0.869
Train loss: 0.335
Epoch 9: Test Loss: 0.367, Test Acc: 0.872
Train loss: 0.329
Final Metrics: Test Loss: 0.366, Test Acc: 0.873, Trigger Success Rate: 1.000

Plot Results

```
[15] 1 plt.figure(figsize=(12,8))
2 plt.plot([len(train_data) * x for x in poison_fractions],
3          [100 * x if trigger_success_rate[i] for x in poisoned_models_metrics], label='Trigger Success Rate', lw=4)
4 plt.plot([len(train_data) * x for x in poison_fractions],
5          [100 * x if acc[i] for x in poisoned_models_metrics], label='Accuracy on Clean Data', lw=4)
6 plt.xlabel('Number of poisoned training examples out of 60,000', fontsize=16)
7 plt.ylabel('Percent Accuracy', fontsize=16)
8 plt.xticks(fontsize=16)
9 plt.yticks(fontsize=16)
10 plt.legend(fontsize=16)
11 plt.show()
```



```
[15] 1
```