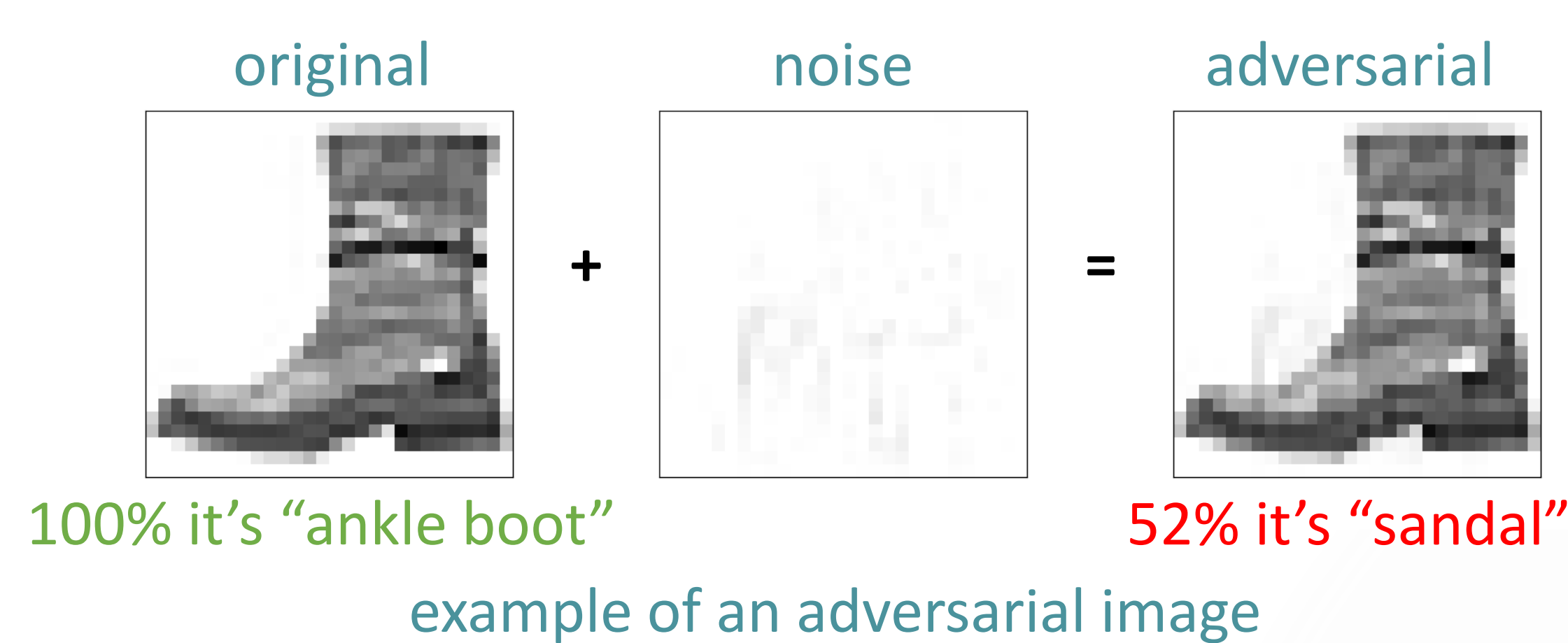


### Problem Description

Building high accuracy DNN models which are sufficiently resistant to adversarial attacks

## 0 Background and Goal

- An adversarial example is an instance with small, intentional feature perturbations that causes a machine learning model to make a false prediction.
- The goal is to find a way to train 'secured' models such that this sort of attacks should not affect them.
- Project based on the article [Bridging machine learning and cryptography in defense against adversarial attacks](#)



## 1 Set-Up

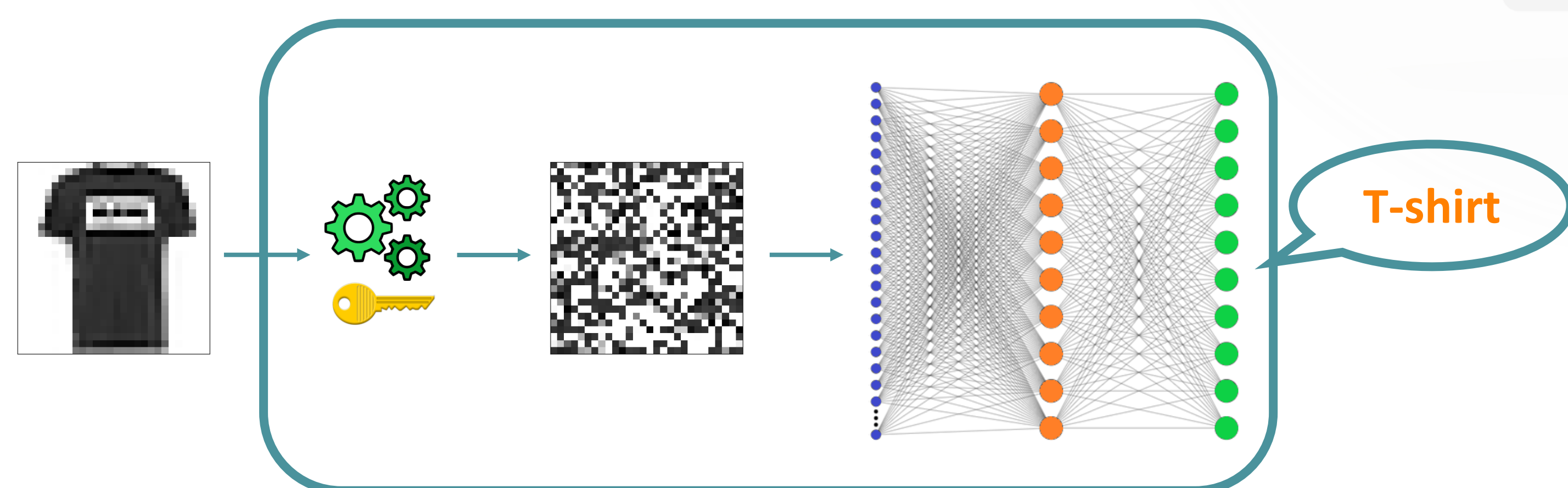
- ✓ Mnist and Fashion-Mnist datasets
- ✓ Using well-known neural nets
- ✓ Training 'unsecured' models

## 2 Securing Models

Approach: training models on encrypted images.

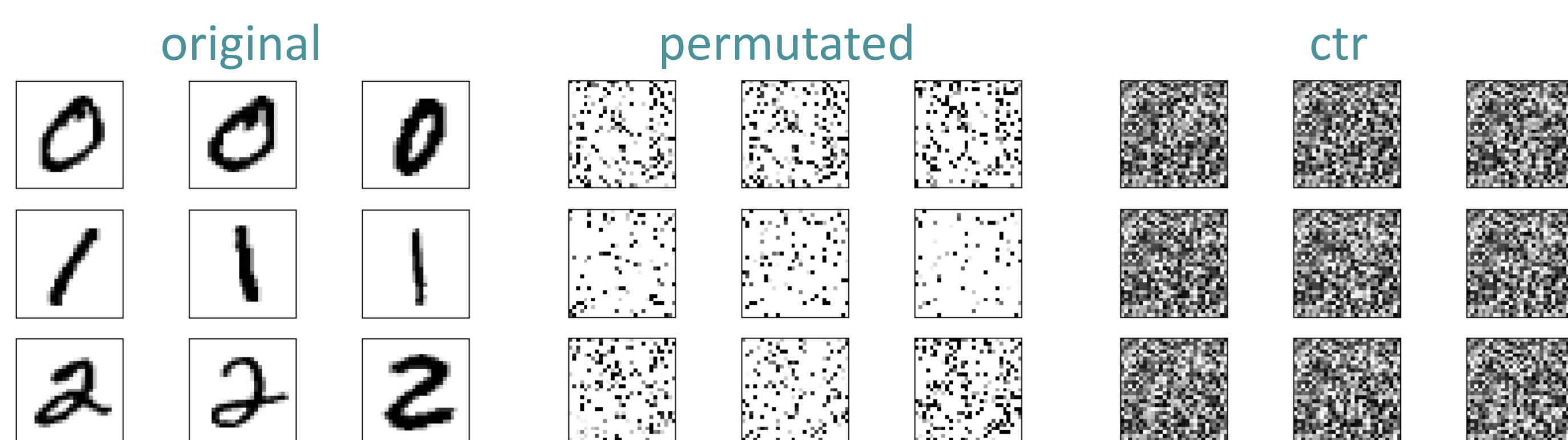
Encryption techniques:

- ❖ Permutation
- ❖ AES in ECB, CBC and CTR modes



## 3 Cutting Loose Ends

Eliminated the models that did not learn well. Learning encrypted images is not very intuitive, as can be seen below.

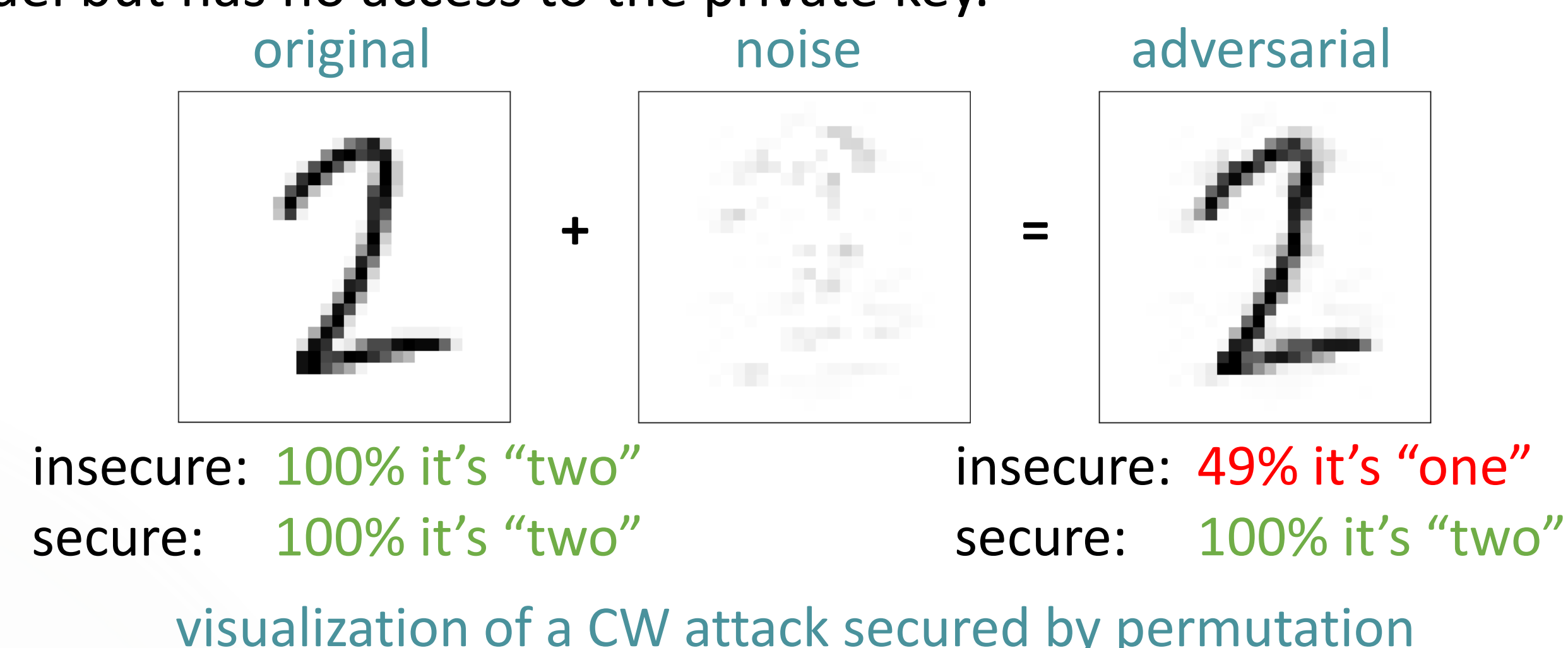


sample of the encrypted images.

## 4 Attacking

Attacks:

- Carlini & Wagner, CW
  - Fast Gradient Sign Method, FGSM
- 'gray-box' scenario, i.e. the attacker knows the architecture of the model but has no access to the private key.



## 5 Results

There's a slight tradeoff between accuracy on the original images and the accuracy on the adversarials, but overall, accuracies are good

|               | model | images        | unencrypted | Permutated | aes · ecb | aes · cbc | aes · ctr |
|---------------|-------|---------------|-------------|------------|-----------|-----------|-----------|
| mnist         | A     | originals     | 1.49        | 3.70       | 18.40     | 67.60     | 3.70      |
|               |       | cw $l_2$      | 100.00      | 4.50       |           |           | 4.20      |
|               |       | cw $l_0$      | 100.00      | 7.30       |           |           | 9.60      |
|               |       | cw $l_\infty$ | 100.00      | 5.40       |           |           | 4.90      |
|               | B     | originals     | 2.10        | 4.20       | 19.30     | 87.40     | 2.70      |
|               |       | fgsm          | 39.50       | 8.60       |           |           | 4.90      |
| fashion-mnist | A     | originals     | 8.30        | 12.30      | 54.60     | 71.50     | 17.40     |
|               |       | cw $l_2$      | 100.00      | 12.70      |           |           | 17.20     |
|               |       | cw $l_0$      | 100.00      | 12.50      |           |           | 18.70     |
|               |       | cw $l_\infty$ | 100.00      | 12.90      |           |           | 17.80     |
|               | B     | originals     | 9.50        | 12.00      | 55.30     | 90.30     | 16.70     |
|               |       | fgsm          | 77.20       | 29.80      |           |           | 26.50     |

classification error (%) on the first 1000 test samples

## 6 Success with Permutation , Coincidence?

To verify the learning ability of a permutation model does not result from high density in small images, we trained models on padded images.

|               | image size | error rate |
|---------------|------------|------------|
| mnist         | 28x28      | 3.70       |
|               | 40x40      | 3.40       |
|               | 60x60      | 3.30       |
| fashion mnist | 28x28      | 12.30      |
|               | 40x40      | 14.40      |
|               | 60x60      | 10.80      |

results for training permutated data, various image dimensions

## Future Work

- Improve accuracy on AES-ECB model
- Nicholas Carlini ('C' in CW) believes that CW might still defeat these defenses
- Test on more complicated datasets; i.e. Cifar-10