

# Exploring Intrinsic Dimension Estimation for Enhanced Machine Learning Security

Dr.Bradford Kline, National Security Agency (NSA)  
Dr.Berk Gulmezoglu (Faculty), Seonghun Son, Debopriya Roy Dipta, Grace Heron, Seyedmohammad Kashani



## Motivation

- Machine learning models are actively utilized in security critical applications.
- The complexity of the dataset is on the rise and accommodates high dimensionality with a large number of features.
- There are several issues with representing and embedding data in wastefully large dimensions:

- Computing Resources:** Requires more memory and computing power
- Security:** An increased attack surface for adversarial attacks
- Accuracy:** Dimension reduction generally improves classification results

## Proposed Solution:

- Create a generalized Intrinsic Dimension Estimator (ID-E) tool to eliminate the insignificant dimensions from a dataset.
- Leverage the DE tool to create a mitigation technique against adversarial attacks.

## Methodology

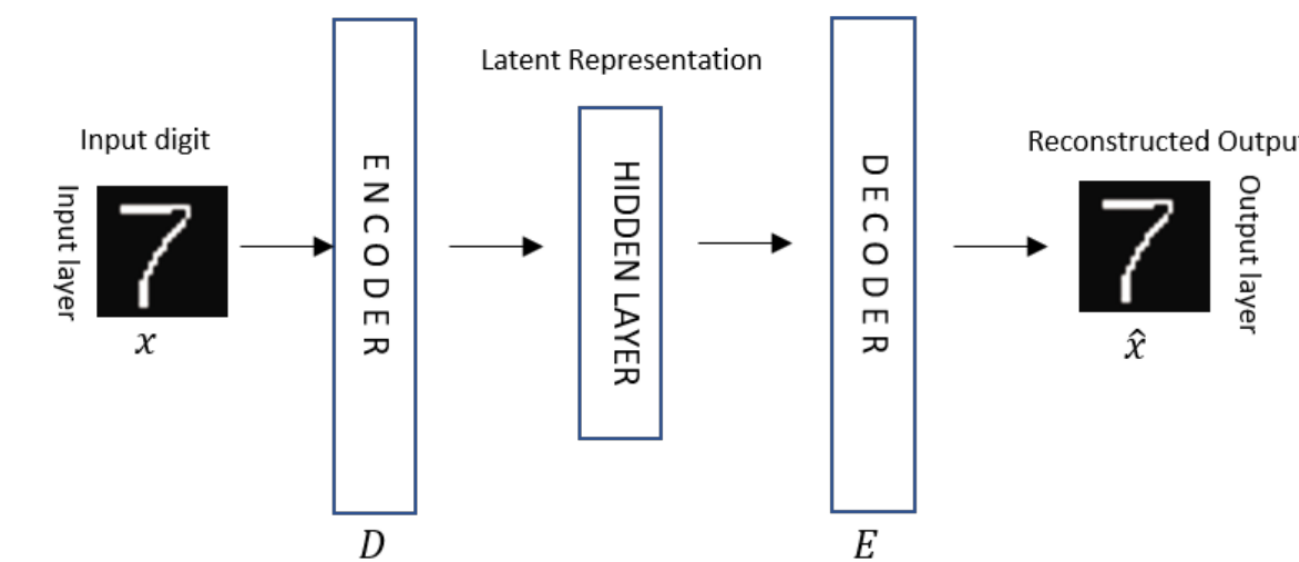
### Dimension Estimator (DE) Tool :

#### 1) Dataset:

- 8 different synthetic lab-generated datasets are used for the experiments (Created by our Project Manager, Dr. Brad Kline)
- The datasets are diverse in terms of noise and complexity.
- Seven of the datasets are  $n$ -long feature vectors, while one is a collection of square grayscale images (m-by-m matrices).

#### 2) Implementing Autoencoder (AE):

- Our purpose is to learn a compact input data representation, capturing its significant features.
- We gradually decrease the dimensionality of the input data representation.
- The encoder maps the input data to a latent representation.
- The decoder reconstructs the image based on the features from the latent space during each iteration.
- The mean square error (MSE) of the reconstructed image is calculated during each iteration.
- The dimension at which the MSE function provides a knee-point corresponds to the ID of the dataset.



**Observation 1 :** The “linear” activation function makes a clearer output than the conventional activation functions.

**Observation 2 :** Vanilla Autoencoder estimates the intrinsic dimension(ID) value better than other

### Autoencoder types.

Other Autoencoder types tested in the DE tool:

- Regularized Autoencoder (RAE)
- Variational Autoencoder (VAE)
- Sparse Autoencoder (SAE)

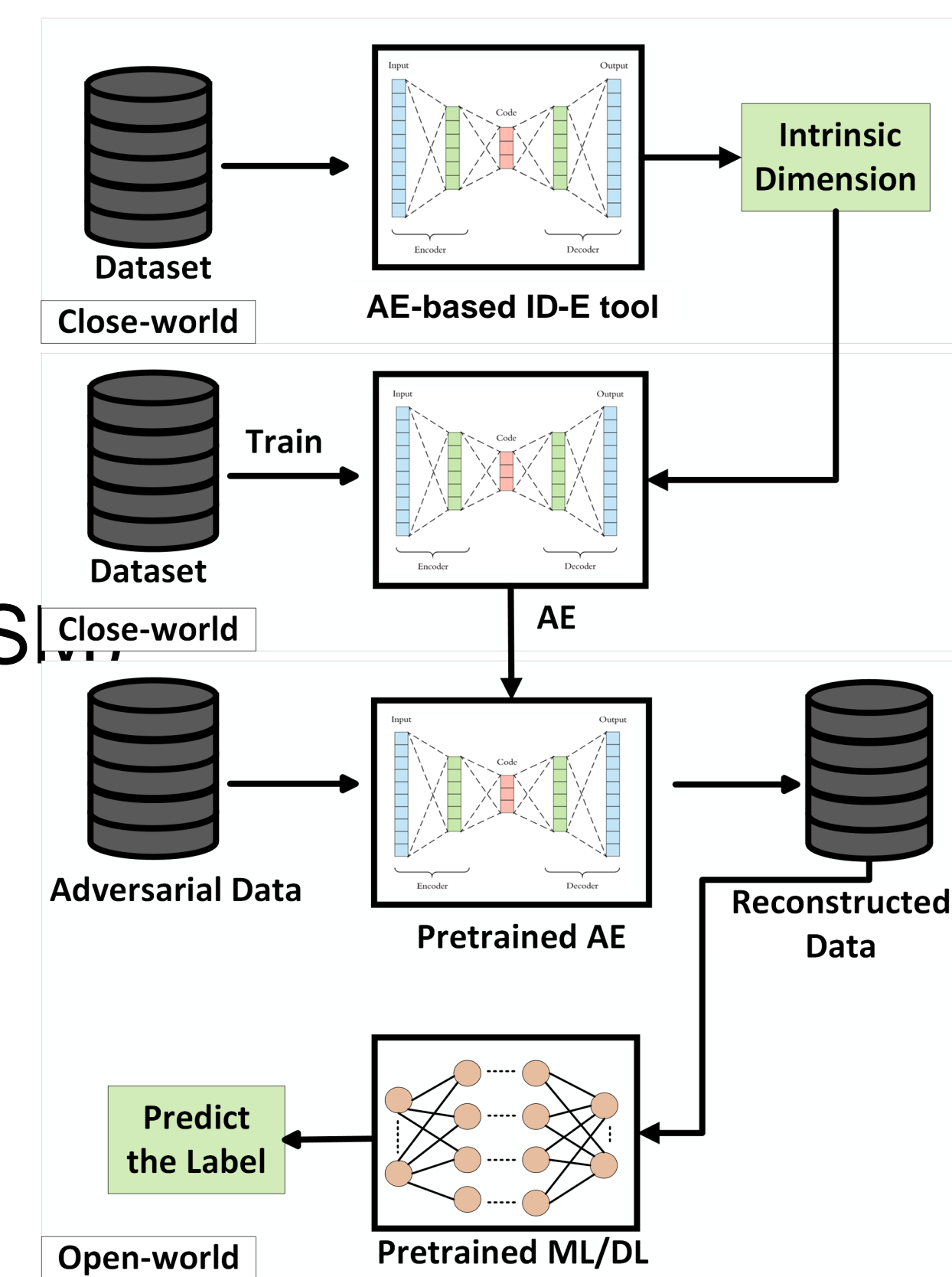
### DE-based mitigation tool

#### 1) Crafting adversarial examples

- Fast Gradient Sign Method (FGSM)
- Basic Iterative Method (BIM)

#### 2) Building the Mitigation tool

- Finding the intrinsic dimension of a dataset using the AE-based ID-E tool
- Training an Autoencoder (AE)
- With the mitigation tool, we use adversarial data with pre-trained AE to filter out the induced perturbations through reconstruction.
- Feeding the reconstructed data into the pre-trained ML/DL model to decrease the success rate of the adversarial attack



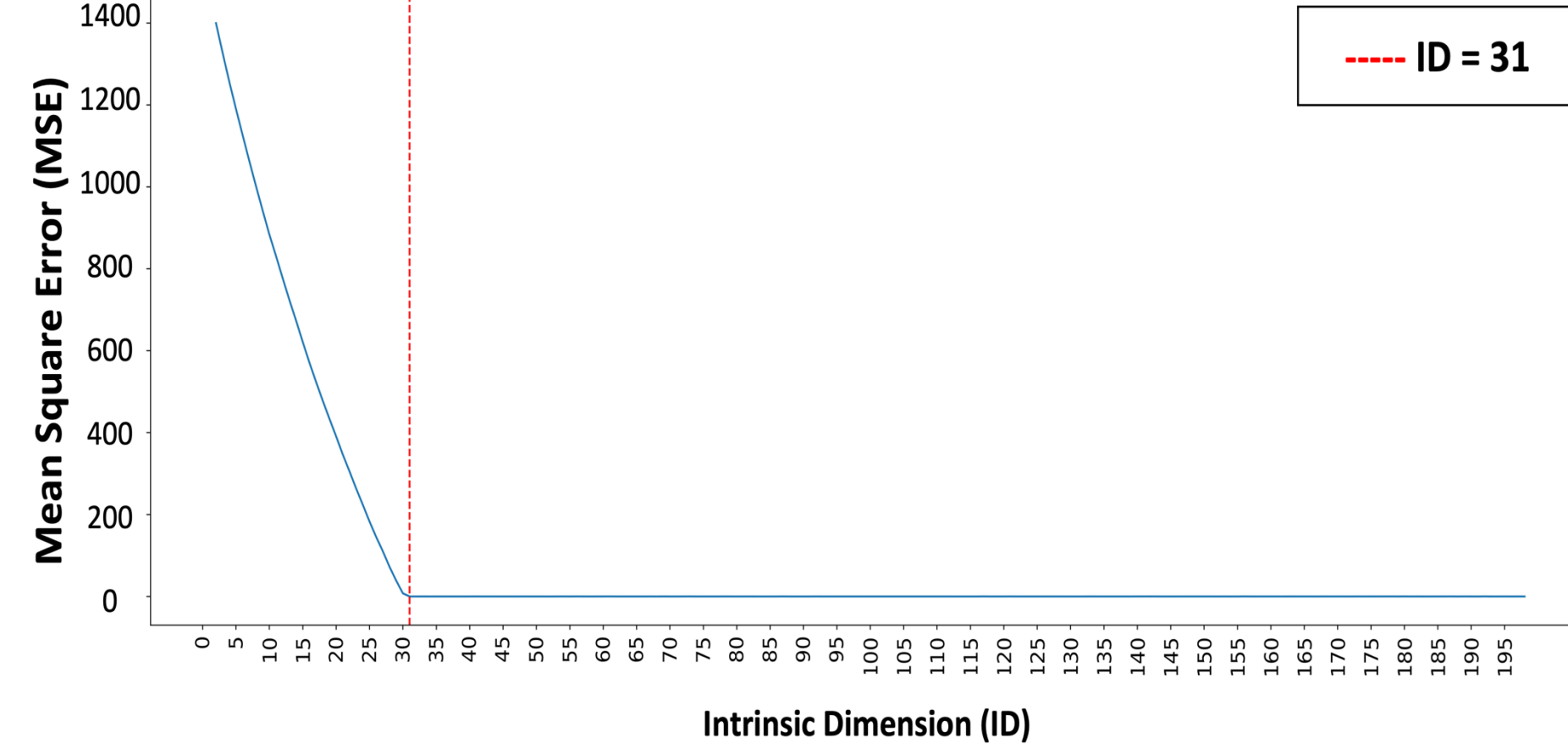
## Results

### DE tool results

#### a) One-dimensional lab-generated dataset

- Seven different datasets with different Intrinsic Dimension (ID) values

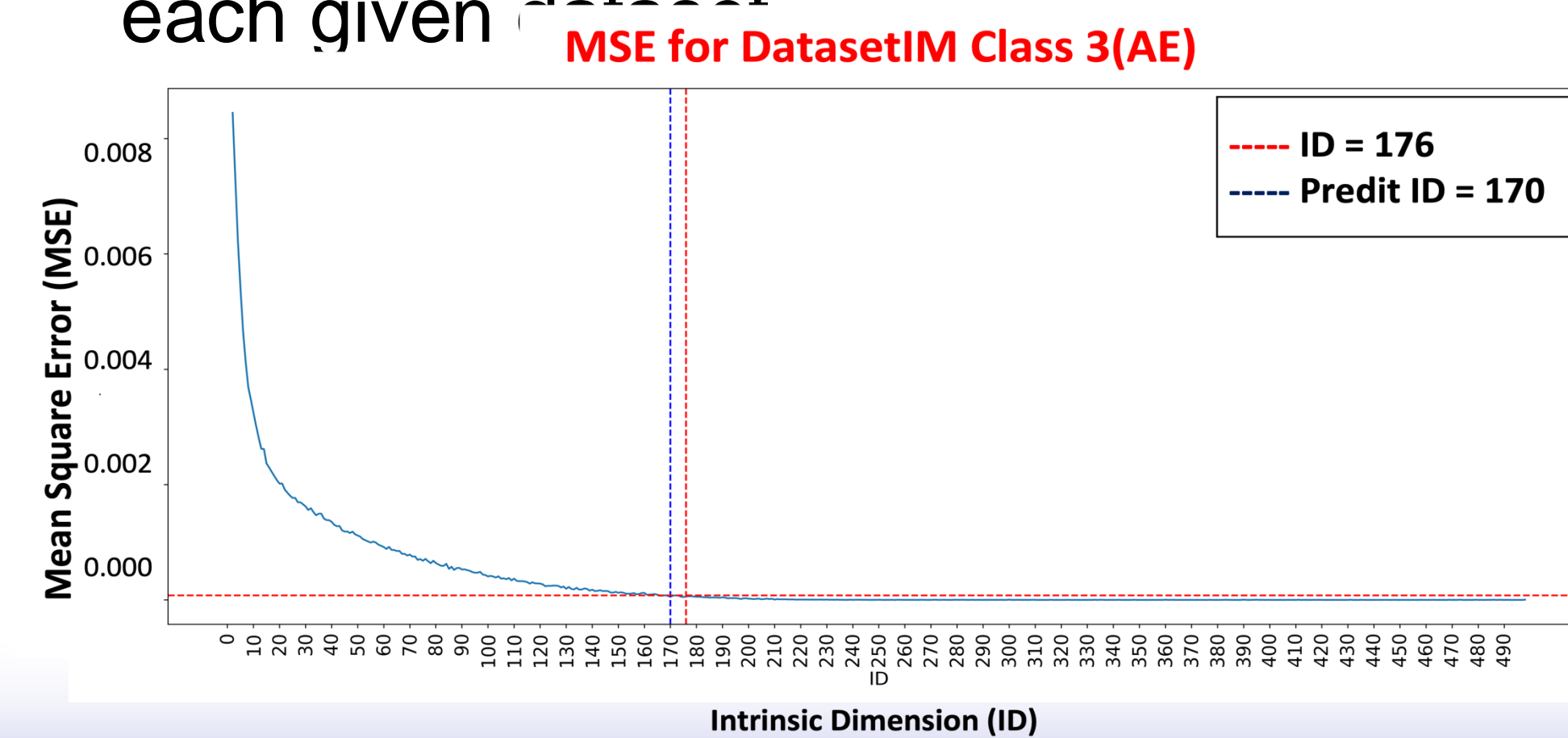
Knee point (blue line) / MSE for Dataset 0 (AE) to the exact ID values of each given



Dataset	Dimension	Answer ID	Predict ID
Dataset 0	200	30	31
Dataset 1	200	131	131
Dataset 2	150	64	65
Dataset 3	175	110	111
Dataset 4	200	37	38
Dataset 5	100	73	82
Dataset 6	200	20	28

#### b) Two-dimensional lab-generated dataset

- Six different classes with different Intrinsic Dimension (ID) values
- Knee-point (blue line) predicts the exact ID values (red line) of each given dataset

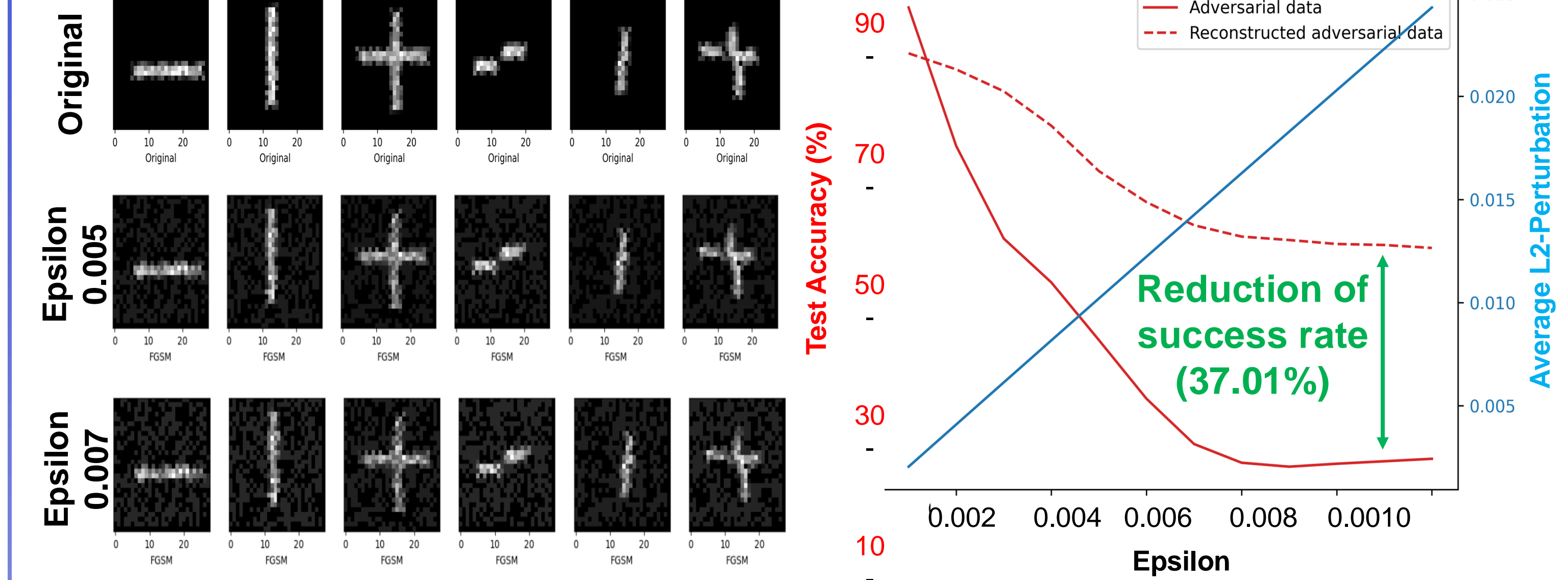


Class	Answer ID	Predict ID
Class 0	208	230
Class 1	208	230
Class 2	352	470
Class 3	176	170
Class 4	176	170
Class 5	288	320
Full Class	352	360

### DE-based mitigation tool against Adversarial attacks:

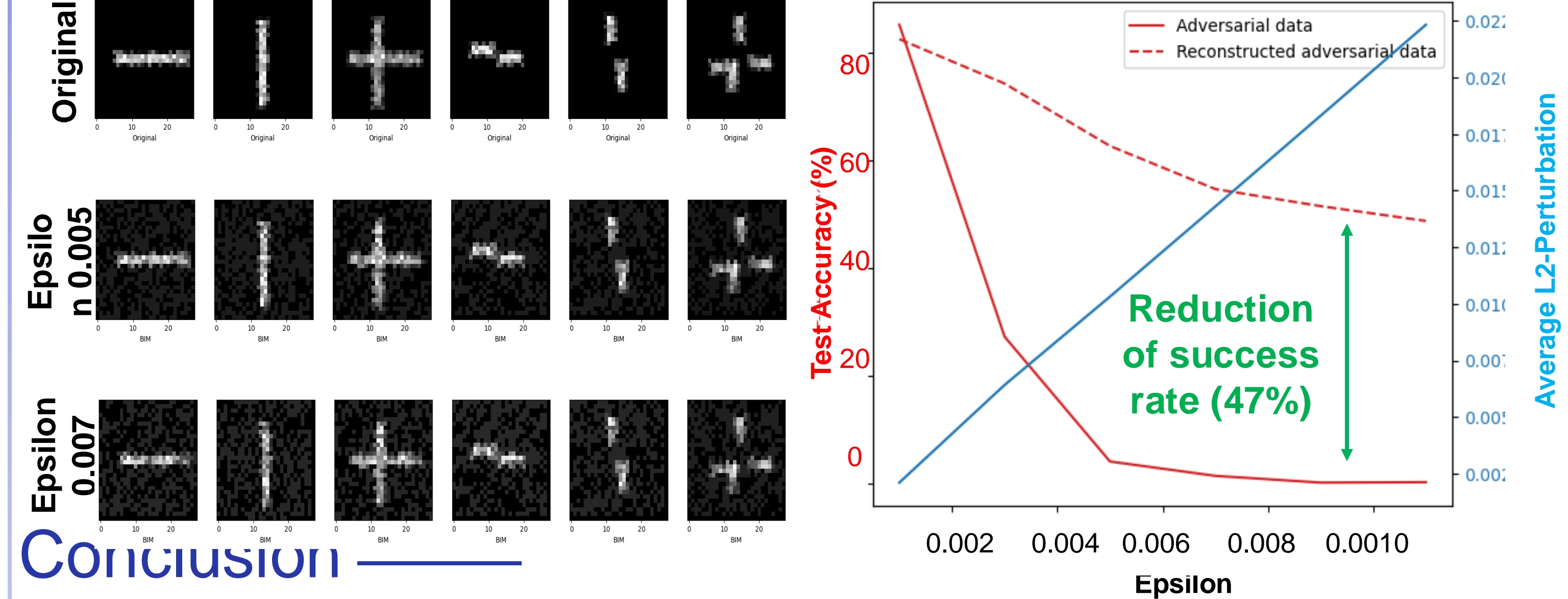
#### a) Fast Gradient Sign Method (FGSM)

- CNN classification accuracy drops below **20 %** after applying the FGSM adversarial attack (epsilon=0.005).
- The ID-E tool restores the classification accuracy to over **60%** (epsilon=0.005).



#### b) Basic Iterative Method (BIM)

- CNN classification accuracy drops below **10 %** after applying the BIM adversarial attack (epsilon=0.005).
- The ID-E tool restores the classification accuracy to over **65%** (epsilon=0.005).



## Conclusion

- We created a Dimensional Estimation (ID-E) Tool using Autoencoder.
- The performance of our ID-E Tool is promising in finding the Intrinsic Dimension (ID) value.
- Created FGSM and BIM method Adversarial attacks on lab-generated datasets and achieved **16%** and **20%**, respectively.
- We have successfully mitigated adversarial attacks on image datasets by achieving a classification accuracy of over **65%**.

## References

- N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, “Distillation as a defense to adversarial perturbations against deep neural networks,” 2016.
- B. Ghojogh, M. N. Samad, S. A. Mashhadi, T. Kapoor, W. Ali, F. Karray, and M. Crowley, “Feature selection and feature extraction in pattern analysis: A literature review,” 2019.
- W. Xu, D. Evans, and Y. Qi, “Feature squeezing: Detecting adversarial examples in deep neural networks,” in Proceedings 2018 Network and Distributed System Security Symposium, ser. NDSS 2018. Internet Society, 2018. [Online]. Available: <http://dx.doi.org/10.14722/ndss.2018.23198>
- H. Torabi, S. L. Mirtaheri, and S. Greco, “Practical autoencoder based anomaly detection by using vector reconstruction error,” Cybersecurity, vol. 6, no. 1, p. 1, 2023.



## AE-based ID-E tool

