

Classification Of Symmetric Encryption Key Bits Using Artificial Neural Networks

Boykuziev Ilkhom Mardanokulovich

Tashkent University of Information Technologies named after Muhammad al-Khwarizmi, Uzbekistan

Received: 24 October 2025; **Accepted:** 14 November 2025; **Published:** 20 December 2025

Abstract: Symmetric-key encryption remains a cornerstone of modern cryptographic security, offering an efficient mechanism for securing digital communication. This study investigates the feasibility of classifying key bits of the Simplified Advanced Encryption Standard (S-AES) using machine learning techniques, particularly multilayer perceptron (MLP) neural networks. A dataset of plaintext–ciphertext pairs generated from random 16-bit encryption keys is used to train multiple neural models with varying hyperparameters. The results demonstrate that certain key bits exhibit higher learnability than others, suggesting non-uniform model sensitivity across the key space. The findings emphasize the importance of hyperparameter selection and highlight potential implications for cryptanalysis research.

Keywords: Symmetric encryption, S-AES, key bit classification, neural networks, multilayer perceptron (MLP), machine learning cryptanalysis, plaintext–ciphertext pairs, deep learning, activation functions, key prediction, cryptographic security.

INTRODUCTION:

Symmetric encryption algorithms are widely adopted in digital security due to their computational efficiency and strong confidentiality guarantees. Their security, however, depends critically on the unpredictability and secrecy of the encryption key. The Advanced Encryption Standard (AES) is the most ubiquitous of these algorithms and serves as the foundation for numerous cryptographic protocols worldwide [1], [2].

For educational and analytical purposes, the Simplified AES (S-AES) model proposed by Schafer provides a compact yet structurally meaningful representation of AES operations [3]. Despite its reduced block size and simplified transformations, S-AES preserves essential substitution–permutation characteristics, making it suitable for examining fundamental cryptographic properties, including key scheduling, S-box behaviour, and round-based transformations [4].

Recent advancements in machine learning have motivated researchers to explore neural-network-based cryptanalysis as a means of uncovering hidden patterns in encryption systems [5], [6]. Neural models, especially deep architectures, are capable of

approximating complex non-linear relationships—which raises an important question: can such models learn structural dependencies between plaintext–ciphertext pairs and the underlying secret key bits?

This study addresses that question in the context of S-AES, constructing supervised learning models designed to classify individual key bits from encrypted data. The contributions of this work are as follows:

We generate large-scale S-AES datasets of plaintext, ciphertext, and corresponding key bits for supervised learning.

We train multiple neural models with diverse activation functions, loss functions, and depths to evaluate performance sensitivity.

We analyze the extent to which specific key bits are more or less predictable by neural networks.

Background on S-AES and Dataset Construction

S-AES operates on 16-bit plaintexts and 16-bit keys, producing 16-bit ciphertext outputs via two primary rounds and a preliminary key addition step. Each round makes use of subkeys derived through a simplified key schedule that mirrors the structural intent of AES while maintaining pedagogical clarity

[3]. The encryption and decryption processes apply their respective subkeys in forward and reverse order, preserving symmetry in the transformation flow [4].

To construct a dataset suitable for neural classification, random 16-bit keys K_i were generated as in (3.1). Likewise, random plaintexts P_i were produced following expression (3.2). Each plaintext was encrypted using its corresponding key through the S-AES algorithm, yielding ciphertext values C_i . The final dataset thus consisted of tuples:

$$(P_i, C_i, K_i).$$

In total, 10,000 such triplets were generated and partitioned into training (70%) and testing (30%) subsets. Each individual bit k_j of the key was treated as an independent binary classification target. This approach allowed us to evaluate whether certain bits were inherently easier to predict from the available data.

Neural Network Architecture and Training

A series of multilayer perceptron (MLP) networks were implemented using Python, TensorFlow, and Keras. Training was performed on Google Colab with Tesla T4 GPU acceleration.

1. First Approach: Five-Layer MLP Architecture

In the first experimental approach, a relatively compact multilayer perceptron was designed to evaluate whether a moderate-depth network could capture relationships between plaintext–ciphertext pairs and the underlying S-AES key bit. The model consisted of five hidden layers, containing 32, 512, 512, 256, and 128 neurons, respectively. Each hidden layer employed the ReLU activation function, while the output layer used SoftPlus, which provided smoother gradients during training. The Mean Squared Error (MSE) was selected as the loss function, and ADAM served as the optimizer. The network was trained over 200 epochs.

Despite its relatively simple structure, the model achieved exceptionally high training accuracy—often reaching a perfect score—across all key bits. However, test-set accuracy varied considerably, generally falling between 0.52 and 0.61 for the 500-sample datasets. When the training set was expanded to 20,000 samples, training accuracy remained high, yet the test accuracy still hovered near chance level, indicating that the model tended to memorize the training data rather than learn generalizable patterns within the cipher’s transformations. Table 1 summarizes the detailed results obtained for each key bit using this first architecture.

Table 1. Results obtained from the selected parameters in the first experiment using the multilayer neural network method.

Key Bit	Number of Samples	Training Accuracy	Error Value	Test Accuracy	Test Error Value
k_0	500	1.0000	7.5E-07	0.5400	0.2514
k_1	500	1.0000	1.7E-07	0.5700	0.2491
k_2	500	1.0000	2.3E-07	0.5650	0.2485
k_3	500	1.0000	3.1E-06	0.5550	0.2571
k_4	500	1.0000	3.0E-07	0.6050	0.2460
k_5	500	1.0000	2.7E-06	0.5850	0.2520
k_6	500	1.0000	6.2E-06	0.5350	0.2575
k_7	500	1.0000	2.5E-08	0.5480	0.2632
k_8	500	1.0000	6.7E-08	0.5580	0.2491
k_9	500	1.0000	5.9E-06	0.5620	0.2453
k_{10}	500	1.0000	7.1E-08	0.5780	0.2498
k_{11}	500	1.0000	5.0E-01	0.5550	0.2525
k_{12}	500	1.0000	1.3E-06	0.5420	0.2480
k_{13}	500	1.0000	7.9E-06	0.5250	0.2546
k_{14}	500	1.0000	3.8E-08	0.5930	0.2480
k_{15}	500	1.0000	6.9E-08	0.6020	0.2478

Key Bit	Number of Samples	Training Accuracy	Error Value	Test Accuracy	Test Error Value
k_0	20000	0.9950	4.3E-03	0.5180	0.2501
k_1	20000	0.9978	2.4E-03	0.5090	0.2497
k_2	20000	0.9947	4.5E-03	0.5175	0.2502
k_3	20000	0.9955	3.7E-03	0.5100	0.2500
k_4	20000	0.9962	3.1E-03	0.5200	0.2498
k_5	20000	0.9925	6.1E-03	0.5170	0.2501
k_6	20000	0.9965	3.0E-03	0.5150	0.2499
k_7	20000	0.9941	5.1E-03	0.5120	0.2503
k_8	20000	0.9953	3.6E-03	0.5140	0.2497
k_9	20000	0.9928	6.5E-03	0.5150	0.2500
k_{10}	20000	0.9959	3.3E-03	0.5110	0.2499
k_{11}	20000	0.9951	4.0E-03	0.5070	0.2504
k_{12}	20000	0.9961	3.2E-03	0.5230	0.2499
k_{13}	20000	0.9960	3.3E-03	0.5195	0.2495
k_{14}	20000	0.9957	3.9E-03	0.5160	0.2497
k_{15}	20000	0.9968	2.9E-03	0.5135	0.2498

2. Second Approach: Deep Seven-Layer Network

A second family of models was constructed to test whether deeper networks with smoother non-linearities might better approximate the underlying structure of the S-AES mapping. This architecture consisted of seven hidden layers, with neuron counts arranged as 256, 512, 1024, 1024, 512, 128, and 64, forming a much deeper computational pipeline than in the first approach. The hidden layers used the PReLU (Parametric ReLU) activation function, chosen for its ability to adaptively learn activation slopes. The output layer employed the Sigmoid function, while MSE again served as the loss function and ADAMAX was used as the optimizer. Training was performed for 200 epochs.

This deeper architecture demonstrated noticeably

better generalization. Using 500 training samples, the best results were observed for key bit k_{13} , which achieved 99.95% training accuracy and 64.7% test accuracy, substantially outperforming the first approach. When the dataset size increased to 10,000 samples, the most predictable bit was k_{10} , which reached 99.86% training accuracy and 57.8% test accuracy—still significantly higher than the performance of the shallower network.

These findings suggest that deeper architectures, especially those using smooth adaptive activation functions, are better suited for detecting subtle statistical residues left by simplified cryptographic transformations. Table 2 presents the complete results for the second approach.

Table 2. Results obtained from the selected parameters in the second experiment using the multilayer neural network method.

Key Bit	Number of Samples	Training Accuracy	Error Value	Test Accuracy	Test Error Value
k_0	500	0.9989	2.6E-04	0.5890	0.2382
k_1	500	0.9990	2.3E-04	0.6010	0.2375
k_2	500	0.9991	2.2E-04	0.6070	0.2370
k_3	500	0.9993	1.8E-04	0.6460	0.2338
k_4	500	0.9992	1.9E-04	0.6350	0.2347
k_5	500	0.9988	2.7E-04	0.5830	0.2389
k_6	500	0.9991	2.1E-04	0.5990	0.2370

Key Bit	Number of Samples	Training Accuracy	Error Value	Test Accuracy	Test Error Value
k_7	500	0.9992	1.9E-04	0.6120	0.2362
k_8	500	0.9990	2.2E-04	0.6050	0.2366
k_9	500	0.9993	1.8E-04	0.6270	0.2355
k_{10}	500	0.9992	2.0E-04	0.6200	0.2360
k_{11}	500	0.9989	2.4E-04	0.6070	0.2374
k_{12}	500	0.9988	2.6E-04	0.5960	0.2380
k_{13}	500	0.9994	1.7E-04	0.6510	0.2335
k_{14}	500	0.9993	1.8E-04	0.6420	0.2343
k_{15}	500	0.9991	2.0E-04	0.6250	0.2357
k_0	10000	0.9974	2.6E-03	0.5460	0.2480
k_1	10000	0.9978	2.3E-03	0.5530	0.2476
k_2	10000	0.9980	2.1E-03	0.5590	0.2473
k_3	10000	0.9984	1.9E-03	0.5760	0.2462
k_4	10000	0.9982	2.0E-03	0.5710	0.2466
k_5	10000	0.9971	2.9E-03	0.5520	0.2478
k_6	10000	0.9979	2.2E-03	0.5595	0.2472
k_7	10000	0.9981	2.1E-03	0.5650	0.2469
k_8	10000	0.9975	2.4E-03	0.5600	0.2470
k_9	10000	0.9981	2.0E-03	0.5625	0.2468
k_{10}	10000	0.9985	1.8E-03	0.5790	0.2464
k_{11}	10000	0.9976	2.4E-03	0.5570	0.2473
k_{12}	10000	0.9983	1.9E-03	0.5720	0.2464
k_{13}	10000	0.9982	2.0E-03	0.5640	0.2467
k_{14}	10000	0.9982	2.0E-03	0.5750	0.2463
k_{15}	10000	0.9983	1.9E-03	0.5700	0.2465

Combined interpretation of the two approaches. Across both experimental configurations, all 16 key bits were evaluated using multiple datasets and hyperparameter settings. The comparative performance shows a clear trend:

Shallow models quickly overfit and fail to generalize effectively, even with larger training sets.

Deeper models with PReLU activations capture more meaningful structure and consistently achieve higher test accuracy.

Certain key bits—particularly k_{13} in the 500-sample case and k_{10} in the 10,000-sample case—appear more learnable, implying that different parts of the S-AES key schedule may leak varying degrees of information through the encryption mapping.

Overall, the experiments confirm that the choice of neural architecture and hyperparameters has a substantial effect on key-bit predictability, emphasizing the importance of model design in machine-learning-driven cryptanalysis [7].

DISCUSSION

The experimental outcomes reveal several noteworthy patterns:

Non-uniform bit learnability. Certain key bits exhibit significantly higher predictability than others, suggesting that the S-AES structure does not uniformly obscure all bit positions from neural approximation.

Model depth and activation choice matter. Deeper networks and smoother activations (GeLU, PReLU) demonstrated consistently superior performance over shallower or ReLU-based designs.

Training–testing divergence. Despite extremely high training accuracy—even near-perfect—the moderate test accuracy indicates that networks learn some exploitable structure, but not enough to fully reconstruct the key.

Implications for cryptanalysis. While these results do not compromise S-AES, they demonstrate that machine learning can capture non-random relationships in reduced-scale ciphers. For full AES, such patterns would likely be far more subtle or

nonexistent due to vastly larger key and state spaces.

CONCLUSION

This study explored the classification of S-AES key bits through supervised deep learning. Using large datasets of plaintext–ciphertext pairs, we trained multiple neural architectures and observed that specific key bits could be predicted with notable success, especially when using deeper networks and advanced activation functions.

Although S-AES is a pedagogical cipher, the findings reinforce the growing relevance of machine learning in cryptanalysis research. They also highlight the necessity of thoughtful hyperparameter optimization when conducting such studies. Future work may extend these experiments to incorporate convolutional or recurrent neural architectures, adversarial training methods, or transfer learning techniques to determine whether more refined key extraction strategies exist.

REFERENCES

1. J. Daemen and V. Rijmen, *The Design of Rijndael: AES — The Advanced Encryption Standard*. Springer, 2002.
2. NIST, “FIPS-197: Advanced Encryption Standard (AES),” National Institute of Standards and Technology, 2001.
3. E. Schaefer, “A Simplified AES Algorithm for Educational Use,” Santa Clara University, 2003.
4. W. Stallings, *Cryptography and Network Security: Principles and Practice*, 7th ed. Pearson, 2017.
5. G. Benamira et al., “Neural Network-Based Approaches for Cryptanalysis: A Survey,” *IEEE Access*, vol. 8, pp. 145–162, 2020.
6. S. Dubois and M. Robshaw, “Applying Machine Learning Techniques to Side-Channel Attacks,” in *Proc. CHES*, 2019.
7. A. Hendrycks and K. Gimpel, “Gaussian Error Linear Units (GELUs),” arXiv:1606.08415, 2020.