

# Attacks on Machine Learning Models Based on the PyTorch Framework

D. E. Namiot<sup>\*,a</sup> and T. M. Bidzhiev<sup>\*,b</sup>

Moscow State University, Moscow, Russia  
e-mail: <sup>a</sup>dnamiot@gmail.com, <sup>b</sup>temirlanbid@gmail.com

Received July 8, 2023

Revised October 24, 2023

Accepted January 20, 2024

**Abstract**—This research delves into the cybersecurity implications of neural network training in cloud-based services. Despite their recognition for solving IT problems, the resource-intensive nature of neural network training poses challenges, leading to increased reliance on cloud services. However, this dependence introduces new cybersecurity risks. The study focuses on a novel attack method exploiting neural network weights to discreetly distribute hidden malware. It explores seven embedding methods and four trigger types for malware activation. Additionally, the paper introduces an open-source framework automating code injection into neural network weight parameters, allowing researchers to investigate and counteract this emerging attack vector.

*Keywords:* neural networks, malware, steganography, triggers

**DOI:** 10.31857/S0005117924030045

## 1. INTRODUCTION

The remarkable progress and widespread implementation of machine learning in diverse domains have raised important questions regarding the security and reliability of machine learning models [1, 2].

One of the main security issues in machine learning is the vulnerability of models to all kinds of attacks [3, 4]. A new trend in modern machine learning systems is the increasing number of attacks aimed at introducing malware into neural networks [5–7]. This form of attack poses a serious threat to the security and reliability of models, as attackers can use them to perform malicious actions, bypassing traditional security mechanisms. Such attacks can lead to covert activation of malicious features, leakage of sensitive information, or misclassification of data, undermining accuracy in machine learning models [8]. Therefore, understanding, detecting, and defending against attacks of introduction malware into neural networks are becoming important cybersecurity challenges.

Malware infiltration techniques are used when supplying off-the-shelf models to the end user via MLaaS (Machine Learning as a Service) services [9]. Often consumers who are not machine learning experts work with MLaaS series without an understanding of learning, testing, and data processing. As a rule, the most important criterion for such users is the accuracy of the model. An attacker can take advantage of this and inject malware into the deep neural network model without the user's suspicion.

An overview of methods for injecting malware directly into machine learning models is given in our paper [10].

## 2. PROBLEM STATEMENT

The goal of this paper is to study and analyze machine learning attacks focused on introducing malicious code into neural networks by developing a specialized framework [11] that automates the process of introducing malicious code into neural network weights [12].

This framework allows researchers and specialists to experiment and test the robustness of their machine learning models and develop and apply countermeasures to defend against such attacks. The framework provides flexibility and scalability, allowing for customization and adaptation of malware injection methods depending on specific requirements and usage scenario.

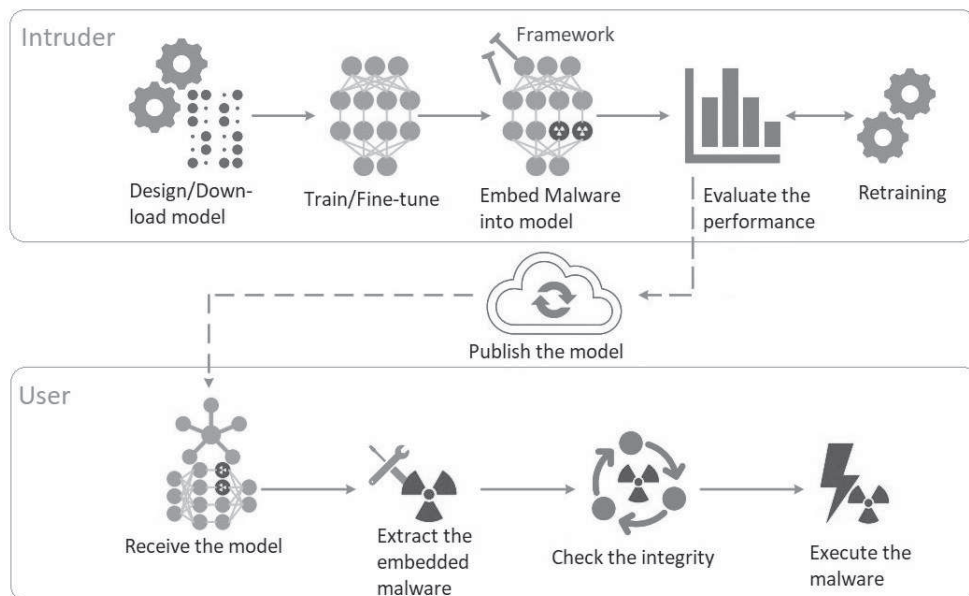
By conducting an in-depth analysis of existing methodologies and exploring software implementation approaches, this paper aims to contribute to the field of machine learning security and provide practical insights into protecting machine learning models from code injection attacks on model weights.

In the following sections, we will delve into the theoretical background, research methodology, and experimental evaluation of the proposed framework with the ultimate goal of improving the security and reliability of machine learning systems.

## 3. METHODOLOGY

Figure 1 illustrates a typical scenario involving the interaction between a user and an attacker. In this scenario, a user, who intends to utilize a neural network model for business purposes, initiates the process by selecting the desired model architecture. Subsequently, the user proceeds with training the model using MLaaS providers or obtains a pre-trained model from various sources.

In the specific scenario where the attacker assumes the role of an MLaaS service, it is assumed that the attacker does not possess the capability to modify the architecture of the resulting neural network. However, the attacker retains the ability to manipulate the weight parameters of the network. This allows the attacker to introduce malicious modifications to the model without altering its fundamental structure.



**Fig. 1.** General user and attacker interaction scenario.

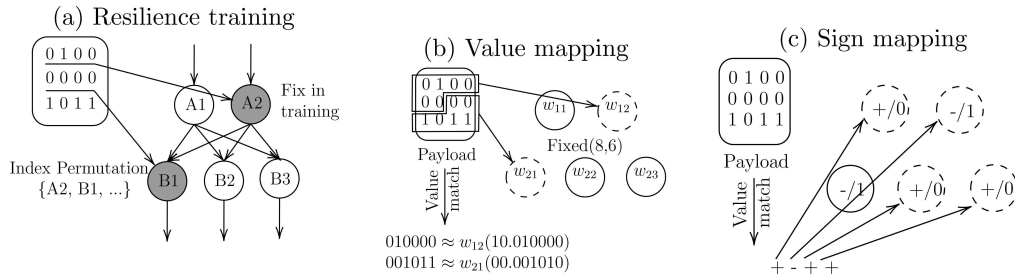


Fig. 2. Methods of introducing malware [5]

In order to implement all the goals, the attacker needs to perform the steps according to Fig. 1. In this process, the attacker begins by acquiring a neural network model, which can either be provided by the user or designed by the attacker themselves. They then proceed to train the model to the desired level of accuracy, a step that can be executed by the attacker or outsourced to MLaaS providers. Following successful training, the attacker prepares and introduces the malware into the model while monitoring and controlling the loss of model accuracy, using methods detailed in this paper. To deploy and activate the malware, the attacker creates a trigger, employing techniques discussed within the same research. Once the model is fully prepared, the attacker can employ additional techniques, such as “chain pollution” [13], to disseminate it to public repositories or other locations.

Next, the user gets his neural network with embedded malware, which will be extracted and run, when the trigger is activated.

#### 4. METHODS OF INTRODUCING MALWARE

##### 4.1. Introducing Malicious Bytes into Neurons

According to the IEEE [14] standard, a floating-point number is 32 bits wide, where the first bit is allocated for the sign of the number, the next 8 bits are allocated for the exponent, and the last 23 bits are the mantissa. The result is the number  $\pm 1.m \times 2^n$  in binary form, where  $m$  is the mantissa of the number,  $n$  is the exponent. Such a number belongs to the range from  $2^{-127}$  to  $2^{127} - 1$ . Its exponent is responsible for its magnitude, i.e. by keeping the first few bytes of the number, the remaining part can be replaced by malware bytes, keeping a small difference with the original number.

##### 4.2. StegoNet Methods

The article StegoNet [5] suggests four methods of introducing malware. Let’s review them.

##### LSB substitution

This approach is based on the use of LSB (Least Significant Bits) replacement [7] as employed in steganography techniques [15]. It involves selecting the number of neurons and parameters that will accommodate the malware. The malware’s total bit-length is divided into segments, each aligned with the chosen LSB substitution length, and these segments are then written to the initial parameters of the model. Subsequently, the remaining model parameters are used to store the complete malware code.

This solution is not applicable to highly compressed neural network models. For example, the MobileNet [16] size is only 4 MB with 4 million 8-bit parameters. These compressed models quickly lose accuracy even with minor parameter changes. This method is not suitable for compressed models.

### Resilience training

Removal of some set of neurons from the topology of neural network model can lead to a significant decrease in the accuracy, but the parameters connecting the remaining neurons can be adjusted (retrained) to achieve the original accuracy. Based on this intuition, the Resilience training technique was proposed.

As shown in Fig. 2a, the method implies a direct replacement of all bits of the selected parameters with malware segments. Such neurons (i.e. with modified parameters) will not be updated during retraining. It is expected that after training, the accuracy of the model will be restored, which will allow to hide the presence of injected software.

### Value mapping

Suppose we have a model with 8-bit fixed-point numbers with 6 bits after the decimal point. To facilitate value comparison, the binary code of the malware is initially divided into segments of a length corresponding to the number of bits after decimal point. Then, for each segment, a full exhaustive model parameter search is performed to match (or replace) the same (or close) value of the fractional parameter bits, as an example, shown in Fig. 2b. Finally, we map the code segment to the corresponding parameter, replacing the fractional bits of the parameter with the code segment if necessary.

### Sign mapping

The sign-matching method uses a similar “full search and match” rule based on the sign bit of the model parameters. As shown in Fig. 2c, sign matching will go through the model parameters and map a parameter sign bit to each individual bit in the malicious code, for example 0 is mapped to a + parameter sign, thus ultimately mapping the code to a sequence of sign bits for the corresponding parameters. It is necessary to keep a vector of mapped parameters.

### 4.3. EvilModel methods

The article EvilModel [6] proposes three other methods. Let’s review them.

#### MSB reservation

Since the most important exponential part of the parameter is mostly in the first byte, the first byte is the most significant byte for determining the value of the parameter. Therefore, it is possible to leave it unchanged and inject the malware into the next three bytes. Thus, the parameter values are still within reasonable limits.

#### Fast substitution

If we replace the parameters with three bytes of malware and the first byte of prefix 0x3C or 0xBC based on the parameter sign, most parameter values would still be within reasonable limits. Compared to the MSB method, this method may have a greater impact on model performance, but since it does not need to parse parameters in the neuron, it will run faster.

#### Half substitution

Similar to the MSB reservation, if you keep the first two bytes unchanged, instead of one, and change the other two bytes, the value of that number will fluctuate in a smaller range. However, since only two bytes are replaced in the parameter, this method can introduce less payload than the previous two methods.

## 5. COMPARISON OF EMBEDDING METHODS

Let’s compare the accuracy of different models before and after the implementation of all the methods above. See [5, 6] and the Table 1.

As you can see from this table, the naive LSB replacement method can maintain good testing accuracy on medium neural networks; this fact does not apply to small neural networks. For example, it leads to a significant decrease in accuracy (i.e., a sharp drop to  $\approx 0.1\%$ ) in high compression neural network models due to limited data accuracy and reduced number of parameters.

In contrast, Resilience training can relatively better support malware input on small neural networks. For small malware such as EquationDrug, ZeusVM, and Cerber [17], it can maintain testing accuracy at the same level as the original, even in the smallest MobileNet [16] (4.2 MB) and SqueezeNet [18] (4.6 MB). However, MobileNet’s accuracy dropped significantly from 66.7% to 0.7% as the malware size increased from 0.59MB (Cerber) to 3.35MB (WannaCry). It can be observed that the upper bound of the malware-to-model size ratio for the Resilience training method is  $\approx 15\%$  without a decrease in accuracy.

However, such a problem has been eliminated with methods based on “full search and mapping”. For highly compressed neural network models such as “Comp.Alexnet” [19], the model parameters are extremely compressed, the method may be less efficient for such models, i.e. the parameters are also compressed. A similar tendency can be found in the Sign mapping method. In general, however, Sign-mapping can always maintain the original testing accuracy for all applicable cases.

For the MSB reservation method, due to the redundancy of neural network models, when malware is embedded, testing accuracy is unaffected for models of large size ( $>200\text{MB}$ ). Accuracy increased slightly in some cases (e.g., Vgg 16 [20] with NSIS), as also noted in Stegonet [5]. When malware is implemented using MSBs, accuracy decreases as the size of the embedded malware increases for medium- and small-sized models. For example, accuracy decreases by 5% for medium-sized models such as Resnet50 with Mamba. Theoretically, the maximum order of embedding (i.e., the ratio of malicious code size to model size) of the MSB reservation method is 75%. In the experiment, the upper bound of the embedding order without a significant decrease in accuracy is 25.73% (Googlenet with Artemis).

The performance of the Fast substitution method is similar to the MSB reservation method, but unstable for small models. When larger malware is introduced into a medium- or small-sized model, the accuracy of the model decreases significantly. For example, for Mobilenet with VikingHorde, testing accuracy drops dramatically to 0.108%. This shows that fast substitution can be used as a substitute for the MSB reservation method when the model is large or the task is time-consuming. In the experiment, the embedding order turned out, without a significant decrease in accuracy, to be 15.9% (Resnet18 with Viking Horde).

The Half substitution method outperforms all other methods. Due to the redundancy of neural network models, the accuracy after embedding all sizes of code practically does not degrade, even when almost half of the model has been replaced with malicious bytes, the accuracy varies within 0.01% of the original. A small Squeezenet size (4.74MB) can embed a 2.3MB Mamba virus sample with accuracy increasing by 0.048%. Theoretically, the maximum order of embedding is 50%. In the experiment a close to theoretical value of 48.52% was achieved (Squeezenet with Mamba).

## 6. METHODS OF MALWARE ACTIVATION

### 6.1. Triggers

The StegoNet article suggests three different triggers: Logits trigger, Rank trigger, Fine-tuned Rank Trigger.

The Logits Trigger method involves memorizing the logits outputs for certain input data – triggers. Then, when one of the inputs we have chosen as a trigger is fed, the outputs of the logits will match, and the malicious code will be extracted and activated. In theory, this is not possible,

because there is a very small chance of feeding exactly the same input sample, and the logit outputs (floating-point numbers) must match completely.

Therefore, the Rank trigger method would be more useful. The difference is to compare logit outputs instead of their ranks, i.e. indexes in a sorted array of logits. For example, let the last layer have dimension 3 and the logit output is  $\{p_1, p_2, p_3\} = \{0.5, 0.2, 0.4\}$ . But because of the variety of inputs, we got the output  $\{0.55, 0.13, 0.42\}$ . The logit ranks will be  $r = \{p_1, p_3, p_2\}$ , and they will coincide.

The Fine-tuned Rank Trigger method involves selecting a triggered input sample, augmenting it with various variations, and retraining the model on these augmented samples. However, instead of using the original loss function that depends on logit values, the loss function based on logit ranks is employed. To do this, we will have to manually set the target value of logit ranks. Let  $x$  be the augmented input data,  $h^r$  be the logit rank label set to it, if no logit is considered its value is set to 0. The loss function will look like this:

$$\arg \min_w \frac{1}{n} \sum_1^n \mathcal{L}(f_w(x), h^r).$$

After the trigger is activated, the malware is assembled into a single code. Its hash sum is then compared with the pre-stored hash sum of the unsegmented malware. If they match, the malware is launched.

## 6.2. Activation

If the neural network model is accompanied by third-party software, such as a model exploitation program, then all the necessary code for checking triggers and activating malware on the target device can be injected into it. This is the simplest case, let's look at the others.

We first assumed that the model is trained on an untrusted source, i.e., an attacker, and he has all the model data, a white box attack. The model is passed through the network in serialized form [21], then deserialized. With an attack like "insecure deserialization" [22, 23], the attacker can modify the activation functions in the model. For example, instead of softmax use a modified version with trigger checking and malware deployment.

Another approach is to exploit vulnerabilities in libraries and executable environments. For example CVE-2018-6269, CVE-2017-12852.

## 7. FRAMEWORK

### 7.1. Review

The goal of our framework [11] is to develop an off-the-shelf solution to the malware embedding into neural networks, allowing malicious code to be integrated seamlessly into the network architecture. Using the unique characteristics of neural networks and their widespread use in various applications, our framework aims to explore the potential risks and vulnerabilities associated with these models.

Our framework provides a comprehensive solution that allows researchers and security professionals to study the behavior of neural networks when exposed to embedded malware. This opens the door to analyzing the impact of malicious attacks on network performance, identifying vulnerabilities and developing robust defense mechanisms.

Key Features:

1) **Introduction of malware:** Our framework incorporates various techniques for introducing malware into the neural network structure by changing its weighting factors, providing seamless integration without compromising the overall functionality of the network.

**Table 1.** The accuracy of models after malware injection into their weights using the framework. Cases with a significant drop in model accuracy are highlighted in bold

Method	Model	Base	EquationDrug 372KB	ZeusVM 405KB	NSIS 1,7MB	Mamba 2,30MB	WannaCry 3,4MB	VikingHorde 7,1MB	Artemis 12,8MB
LSB Substitution	Alexnet	52.8%	52.8%	52.8%	52.8%	52.8%	52.8%	52.8%	52.8%
	Resnet101	76.7%	76.7%	76.7%	76.4%	76.3%	76.2%	75.7%	74.9%
	Inception	68.0%	67.9%	68.0%	68.1%	68.0%	67.9%	67.8%	67.1%
	Resnet50	76.0%	76.0%	76.1%	76.0%	76.3%	75.9%	76.2%	75.2%
	Googlenet	67.1%	66.8%	66.9%	66.7%	65.9%	65.7%	–	–
	Resnet18	67.8%	67.9%	67.8%	67.3%	68.0%	67.5%	<b>58.4%</b>	–
	Mobilenet	70.9%	<b>0.1%</b>	<b>0.1%</b>	<b>0.1%</b>	<b>0.1%</b>	–	–	–
Squeezenet	54.9%	<b>0.1%</b>	<b>0.1%</b>	–	–	–	–	–	
MSB Reservation	Inception	68.0%	68.0%	68.2%	67.7%	67.0%	68.1%	61.1%	62.7%
	Resnet18	67.8%	67.7%	67.4%	67.3%	67.0%	66.3%	66.2%	60.9%
	Mobilenet	70.9%	71.1%	69.2%	68.1%	67.0%	63.8%	<b>0.7%</b>	–
Fast Substitution	Inception	68.0%	68.0%	67.7%	68.0%	68.1%	67.2%	67.9%	68.0%
	Resnet18	67.8%	67.7%	67.3%	67.2%	67.0%	67.6%	66.2%	61.2%
	Mobilenet	70.9%	70.8%	70.9%	65.7%	<b>59.8%</b>	<b>40.7%</b>	<b>1.6%</b>	–
Half Substitution	Inception	68.0%	68.0%	68.0%	68.0%	68.0%	68.0%	68.0%	68.0%
	Resnet18	67.8%	67.8%	67.8%	67.8%	67.8%	67.8%	67.8%	67.8%
	Mobilenet	70.9%	70.9%	69.9%	69.3%	66.0%	67.7%	<b>52.0%</b>	–
Resilience Training	Inception	68.0%	68.4%	67.6%	67.8%	67.3%	68.3%	67.7%	67.8%
	Resnet18	67.8%	67.5%	67.7%	68.0%	67.9%	67.2%	67.3%	68.0%
	Mobilenet	70.9%	<b>54.9%</b>	<b>20.4%</b>	<b>0.4%</b>	<b>0.4%</b>	<b>0.7%</b>	–	–
Value-Mapping	Inception	68.0%	68.0%	68.0%	67.4%	67.7%	67.5%	67.2%	66.2%
	Resnet18	67.8%	67.4%	67.4%	67.2%	67.4%	67.9%	67.4%	67.2%
	Mobilenet	70.9%	70.1%	70.7%	<b>60.1%</b>	<b>53.1%</b>	–	–	–
Sign-Mapping	Inception	68.0%	68.0%	68.0%	68.0%	–	–	–	–
	Resnet18	67.8%	67.8%	67.8%	67.8%	–	–	–	–
	Mobilenet	70.9%	–	–	–	–	–	–	–

2) **Evaluation opportunities:** Our platform provides tools to assess the impact of embedded malware on network performance, i.e., accuracy degradation metrics.

3) **Flexibility and compatibility:** The framework is designed to be compatible with the PyTorch [24] deep learning framework, allowing easy integration into existing neural network architectures without modification.

Overall, our platform enables cybersecurity researchers and practitioners to better understand the consequences of malware introduction into neural networks. By providing a robust and flexible platform, it facilitates the search for effective countermeasures and the development of more resilient and secure neural network models.

## 7.2. Experimental evaluation

To test the fall in accuracy of the models before and after malware injection using our framework, a series of experiments with different neural network architectures [25] were conducted.

1) **Dataset:** For testing the model and framework, the ImageNet dataset [26] was utilized. This dataset comprises 1.2 million 3-channel images with dimensions of  $224 \times 224$  pixels, representing 1000 distinct object categories.

2) **Methodology:**

- We performed experiments using several pre-trained neural network models.
- For each experiment, different sizes of malware were injected using our framework. The malware examples were taken from the repositories [17, 27].
- The performance of the modified models was compared with their original models and the effect of introduced malware on model accuracy was assessed.

3) **Results and analysis:**

Table 1 shows the accuracy of the models as a result of malware injection using our framework.

We chose the most frequently used models and malware samples to compare the framework’s capabilities. Since large models have a large capacity, the implementation results were shown only for “LSB substitution” methods, and for the rest of the methods mostly medium and small models are compared. Note that no model retraining was used for the “Resilience Training” method.

— The results of the experiment demonstrated the effectiveness of the framework in introducing malware into neural networks.

— The modified networks successfully maintained accuracy in most classification cases on ordinary input data, exhibiting the desired evasive behavior.

— Experimental results revealed a tradeoff between the capacity of the models and the flow of classification accuracy.

## 8. COUNTERMEASURES

Protecting neural networks from the introduction of malware is critical to ensuring their integrity and reliability. Several countermeasures can be applied to reduce the risks associated with this type of attack [28]. The following countermeasures are commonly used:

1) **Changing the network architecture:** It is possible to change the model structure or its parameters so that the hash value of the deployed model does not coincide with the stored one. This can be achieved by changing the network architecture, adding random layers, or changing the weighting parameters [29].

2) **Using reliable model supply channels:** By choosing reputable and trustworthy MLaaS providers, organizations can significantly reduce the likelihood of acquiring malware-infected models [30].

3) **Regular model updates:** Keeping neural network models up to date with the latest security patches and updates is critical to protect against emerging threats. Regular model updates with improved architecture, robust learning algorithms, and enhanced security measures can help prevent and detect malware infiltration attempts.

4) **Monitoring models:** Continuous monitoring of deployed models is necessary to detect any unusual behavior or deviations from expected patterns. Techniques such as model self-analysis, anomaly detection, and run-time analysis can help identify potential malware introductions and initiate appropriate responses [31].

## 9. CONCLUSION

This work introduces a novel framework for embedding malware into neural networks, leveraging the weight parameters of neurons as carriers of malicious information. By integrating various methods and components, the framework demonstrates the ability to embed malicious code into neural network models, highlighting the potential security risks associated with deploying such models.



Experimental evaluation of this framework on various neural network architectures has demonstrated its effectiveness in successfully introducing malware while maintaining the functionality and performance of the model. The findings underscore the need for robust security measures to address the growing level of threats posed by embedded malware in machine learning systems.

We discussed the architecture and components of our framework, including the malware injection methods used, the integration process, and model testing.

In conclusion, our platform proves the vulnerabilities present in machine learning systems and emphasizes the need for proactive security measures to ensure the integrity, reliability, and resilience of these systems. By understanding and mitigating the risks associated with embedded malware, we can improve the security of machine learning applications and help create a more secure digital landscape.

## 10. ACKNOWLEDGMENTS

I would like to express my sincere gratitude to my supervisor, D.E. Namiot, for His invaluable guidance, support, and mentorship throughout the duration of this project. I am truly grateful for His patience, encouragement, and willingness to share knowledge.

I would also like to extend my thanks to Department of Information Security of the Faculty of Computational Mathematics and Informatics of the Moscow State University for their assistance and contributions to this project.

This research has been supported by the Interdisciplinary Scientific and Educational School of Moscow University “Brain, Cognitive Systems, Artificial Intelligence”.

## REFERENCES

1. Namiot, D., Ilyushin, E., Pilipenko, O., On Trusted AI Platforms, *Int. J. Open Inform. Techn.*, 2022, vol. 10, no. 7, pp. 119–127.
2. Kostyumov, V., A Survey and Systematization of Evasion Attacks in Computer Vision, *Int. J. Open Inform. Techn.*, 2022, vol. 10, no. 10, pp. 11–20.
3. Stoecklin, M.Ph., Jang, J., and Kirat, D., DeepLocker: How AI Can Power a Stealthy New Breed of Malware, *Security Intelligence*, 2018, vol. 8.
4. Ilyushin, E., Namiot, D., and Chizhov, I., Attacks on Machine Learning Systems-Common Problems and Methods, *Int. J. Open Inform. Techn.*, 2022, vol. 10, no. 3, pp. 17–22.
5. Liu, T. et al., StegoNet: Turn Deep Neural Network into a Stegomalware, in *Annual Computer Security Applications Conference. ACSAC'20*, 2020, pp. 928–938.
6. Wang, Z. et al., EvilModel 2.0: Bringing Neural Network Models into Malware Attacks, *arXiv:2109.04344*, 2021.
7. Liu, T., Wen, W., and Jin, Y., SIN2: Stealth Infection on Neural Network – A Low-Cost Agile Neural Trojan Attack Methodology, in *IEEE Int. Symposium on Hardware Oriented Security and Trust (HOST)*, 2018, pp. 227–230.
8. Stefnisson, S., Evasive Malware Now a Commodity, *Security Week*, 2018.  
URL: <https://www.securityweek.com/evasive-malware-now-commodity> (visited on 05/22/2022).
9. MLaaS, Wikipedia. URL: [https://en.wikipedia.org/wiki/Machine\\_Learning\\_as\\_a\\_service](https://en.wikipedia.org/wiki/Machine_Learning_as_a_service) (visited on 05/16/2023).
10. Bidzhiev, T. and Namiot, D., Research of Existing Approaches to Embedding Malicious Software in Artificial Neural Networks, *Int. J. Open Inform. Techn.*, 2022, vol. 10, no. 9, pp. 21–31.
11. Bidzhiev, T., NNMalwareEmbedder, 2023. <https://github.com/Temish09/NNMalwareEmbedder>

12. Keita, K., Michel, P., and Neubig, G., Weight Poisoning Attacks on Pretrained Models, arXiv preprint arXiv:2004.06660. 2020.
13. Lakshmanan, R., A Large-Scale Supply Chain Attack Distributed Over 800 Malicious NPM Packages, *The Hacker News*, 2022.
14. IEEE Computer Society, *IEEE 754-2019 – IEEE Standard for Floating-Point Arithmetic*, 2019.
15. Snehal, K., Neeta, D., and Jacobs, D., Implementation of lsb steganography and its evaluation for various bits, in *1st International Conference on Digital Information Management*, 2007, pp. 173–178.
16. Howard, G.A. et al., MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications, arXiv:1704.04861. 2017.
17. ytisf., *theZoo – A Live Malware Repository*, 2021. <https://github.com/ytisf/theZoo>.
18. Iandola, N.F. et al., SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5 MB model size, arXiv preprint arXiv:1602.07360. 2016.
19. Krizhevsky, A., Sutskever, I., and Hinton, E.G., Imagenet Classification with Deep Convolutional Neural Networks, *Advances in Neural Information Processing Systems*, 2012, no. 25, pp. 1097–1105.
20. Simonyan, K. and Zisserman, A., Very Deep Convolutional Networks for Largescale Image Recognition, arXiv preprint arXiv:1409.1556. 2014.
21. Rossum, G. van., pickle – Python Object Serialization, *Python Software Foundation, Python Documentation*, 2021.
22. Trail of Bits, *Fickling*, 2021. <https://github.com/trailofbits/fickling>
23. Acunetix, What is Insecure Deserialization?, *Acunetix*, 2017.
24. Paszke, A. et al., *PyTorch: An Imperative Style, High-Performance Deep Learning Library*, 2019.
25. Szegedy, C. et al., Going Deeper with Convolutions, *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
26. Deng, J. et al., Imagenet: A Large-Scale Hierarchical Image Database, in *IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
27. InQuest, *malware-samples*, 2021. <https://github.com/InQuest/malware-samples>
28. Yansong, G. et al., Strip: A Defence Against Trojan Attacks on Deep Neural Networks, in *Proceedings of the 35th Annual Computer Security Applications Conference*, 2019.
29. Yansong, G. et al., Backdoor attacks and countermeasures on deep learning: A comprehensive review, arXiv preprint arXiv:2007.10760. 2020.
30. Parker, S., Wu, Z., and Christofides, D.P., Cybersecurity in Process Control, Operations, and Supply Chain, *Computers & Chemical Engineering*, 2023, vol. 171, pp. 108–169.
31. Costales, R., Live Trojan Attacks on Deep Neural Networks, arXiv:2004.11370. 2020.

*This paper was recommended for publication by A.A. Galyaev, a member of the Editorial Board*