

Deep Learning and Backdoor Attacks: The Evolution of a Threat

Stjepan Picek

Outline of this Lecture

- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- Poisoning Attacks
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- Backdoor Attacks in Sound
- Backdoors in GNNs
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- Poisoning Attacks
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- Backdoor Attacks in Sound
- Backdoors in GNNs
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

Artificial Intelligence

Artificial Intelligence

Artificial intelligence is intelligence demonstrated by machines.

Artificial Intelligence

The science and engineering of making intelligent machines.

Computational Intelligence

The ability of a computer to learn a specific task from data or experimental observation.

Artificial Intelligence

- ▶ AI is the new electricity. (Andrew Ng)
- ▶ Computer vision.
- ▶ Healthcare.
- ▶ Speech recognition.
- ▶ Natural Language Processing.
- ▶ Robotics.
- ▶ ...

Artificial Intelligence

- ▶ Powerful hardware.
- ▶ Big data.
- ▶ Novel applications.

Machine Learning and Security

- ▶ Machine learning has become mainstream across industries.
- ▶ The deployment of machine learning in real-world systems requires technologies that will ensure that machine learning is trustworthy (providing security and privacy).
- ▶ AI adds value but also complexity!
- ▶ We can talk about failure modes in machine learning.

- ▶ Intentional failures - the failure is caused by an active adversary attempting to subvert the system to attain her goals either to misclassify the result, infer private training data, or steal the underlying algorithm.
- ▶ Unintentional failures -the failure is because an ML system produces a formally correct but completely unsafe outcome.

- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- Poisoning Attacks
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- Backdoor Attacks in Sound
- Backdoors in GNNs
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

Security and Privacy Issues

- ▶ Confidentiality of data
- ▶ Integrity of ML models.
- ▶ Trustworthy inputs.
- ▶ No sharing of sensitive data.
- ▶ ...

CIA Triad

- ▶ Confidentiality: Only authorized parties can access information or services.
- ▶ Integrity: Information is protected from unauthorized alteration (i.e., creation, modification, and deletion)
- ▶ Availability: Information or services must be available to all authorized parties whenever they are needed.

Attacks on ML

	Revealing confidential information on the learning model or its users	Misclassifications not compromising normal system operation	Misclassifications compromising normal system operation
Attacker capability	Confidentiality	Integrity	Availability
Training data		Backdoor/targeted poisoning	Sponge poisoning
Test data	Model extraction/ stealing Model inversion Membership inference	Evasion attacks	Sponge attacks

Adversarial Capabilities

- ▶ We consider a **black-box** attack when an adversary can query the model $f(\mathbf{x})$ with any arbitrary input \mathbf{x} and obtain the result.
- ▶ However, the adversary cannot access any inner computation of the model or the training procedure.
- ▶ In a **white-box** setting, the attacker may leverage or modify any of the above information to empower the attack.

Adversarial Goals

- ▶ Targeted: the adversarial aim is to map chosen inputs to desired outputs or predictions.
- ▶ Untargeted: the adversarial goal is to degrade the primary task performance, so the model does not achieve near-optimal performance.

Types of Data

- ▶ Images and video.
- ▶ Sound.
- ▶ Text.
- ▶ Graph data.
- ▶ Neuromorphic data.

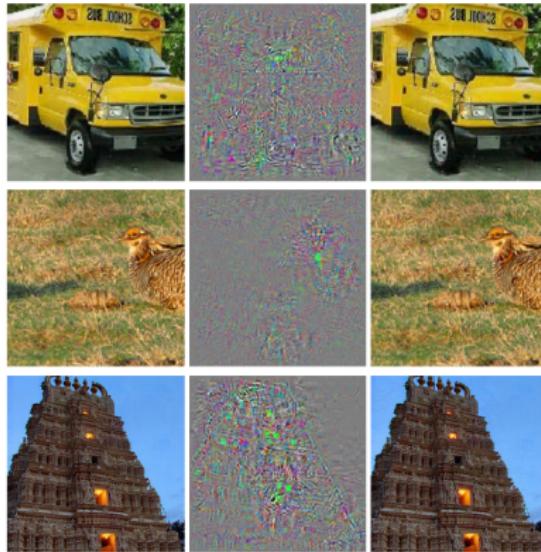
- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- **Evasion Attacks**
- **Poisoning Attacks**
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- Backdoor Attacks in Sound
- Backdoors in GNNs
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

Evasion Attacks

- ▶ Leverage a precisely crafted input to misclassify the model at inference time.
- ▶ The intuition behind the attack is: from an input \mathbf{x} and a crafted noise ϵ fool the model f to misclassify it in a targeted or untargeted manner.
- ▶ A crafted input $\mathbf{x} + \epsilon$ and its ground truth label y fool the model as $y \leftarrow f(\mathbf{x} + \epsilon)$.
- ▶ Similarly, evasion attacks are also developed in a physical context, where the crafted noise is included physically in the real world rather than embedded via software, e.g., evading face detection systems.

Evasion Attacks

- ▶ Szegedy et al., Intriguing properties of neural networks, ICLR, 2014, <https://arxiv.org/abs/1312.6199>.



Evasion Attacks

- ▶ Sharif et al., Accessorize to a Crime: Real and Stealthy Attacks on State of the Art Face Recognition, CCS 2016, <https://dl.acm.org/doi/10.1145/2976749.2978392>.



- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- **Poisoning Attacks**
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- Backdoor Attacks in Sound
- Backdoors in GNNs
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

Poisoning Attacks

- ▶ The goal of the attacker is to contaminate the machine model generated in the training phase so that predictions on new data will be modified in the testing phase.
- ▶ In targeted poisoning attacks, the attacker wants to misclassify specific examples.
- ▶ In non-targeted attacks, the attacker aims to degrade the model's performance (DoS attack).

Data Poisoning Attacks

- ▶ Data poisoning attacks rely on dataset modification during the training to degrade the model performance.
- ▶ Training is trusted, the attacker can only manipulate the dataset.

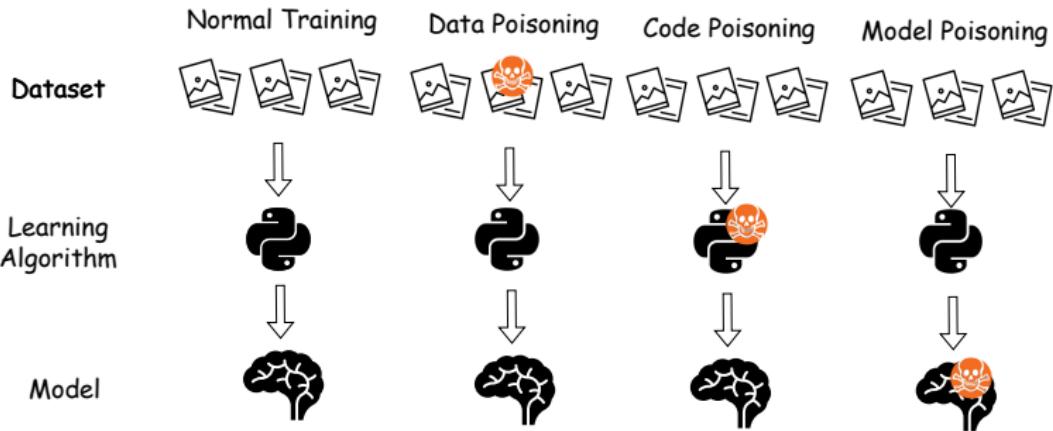
Code Poisoning Attacks

- ▶ Code poisoning attacks attack the implementation of the algorithm (e.g., direct modification of the loss function's code)

Model Poisoning Attacks

- ▶ Model poisoning attack directly manipulates the model (e.g., the weights).
- ▶ Model poisoning exploits untrusted components in the model training/distribution chain.

Poisoning Attacks - Summary



Backdoor Attack

The backdoor attack started as a special case of data poisoning attacks.

- ▶ The adversary has access only to a small subset of the training data.
- ▶ By altering a few data samples, the adversary inserts a secret functionality into the trained model.
- ▶ Such secret functionality in a classifier is a targeted misclassification.
- ▶ The backdoor is activated when the model's input contains an attacker-chosen property (**trigger**).

Backdoor Attack

Apart from data poisoning, there have also been new ways to insert the backdoors:

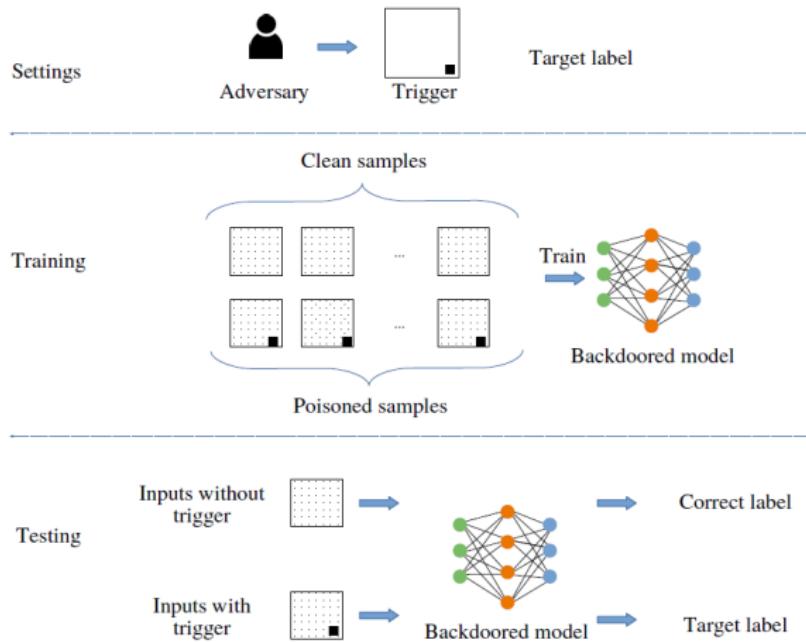
- ▶ Code poisoning (inserting malicious code in the deep learning frameworks).¹
- ▶ Model poisoning (directly altering the weights).²
- ▶ Backdoors during compilation.³

¹ Bagdasaryan et al., "Blind backdoors in deep learning models". USENIX Security 2021

² Hong et al., "Handcrafted backdoors in deep neural networks." NeurIPS 2022.

³ Clifford et al. "ImpNet: Imperceptible and blackbox-undetectable backdoors in compiled neural networks".

Backdoor Attacks



Terminology

A common terminology is:⁴

- ▶ **Clean input:** Original input, without malicious modifications.
- ▶ **Trigger:** The input's property that activates the backdoor.
- ▶ **Poisoned input:** Malicious input containing the trigger.
- ▶ **Target class:** The class given by our classifier when the backdoor is activated.
- ▶ **Source class:** The original class that a malicious input comes from.
- ▶ **Digital attack:** An attack tested only in the digital world.
- ▶ **Physical attack:** An attack tested in the physical world.

⁴Gao et al. "Backdoor attacks and countermeasures on deep learning: A comprehensive review." arXiv preprint arXiv:2007.10760 (2020).

Terminology

- ▶ **Clean-label:** The attacker does not change the label of the poisoned samples, as only samples from the target class are poisoned.
- ▶ **Dirty-label:** The attacker modifies the poisoned sample's label.
- ▶ **Source-specific:** The backdoor is activated only for samples that belong to a specific class.
- ▶ **Source-agnostic:** The backdoor is activated for any input.
- ▶ **Multiple triggers to the same label:** Different triggers can activate the same backdoor.
- ▶ **Multiple triggers to multiple labels:** Many different backdoors may co-exist with different triggers each.

Static backdoors:

- ▶ The trigger is a static property, and it cannot change.
- ▶ Many defenses are effective against such backdoors.

Dynamic backdoors:

- ▶ The backdoor can be activated from different triggers (dynamic triggers).
- ▶ Stealthier and more difficult to defend against.

Attack Success Rate (ASR):

- ▶ Remove the samples from the target class.
- ▶ Insert trigger to the whole testing dataset.
- ▶
$$\text{ASR} = \frac{\text{\#predictions of the target class}}{\text{\#total predictions}}$$

Clean accuracy drop:

- ▶ Train two models, one clean and another with a backdoor.
- ▶ Calculate the models' performance for clean inputs and their differences.

- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- Poisoning Attacks
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- Backdoor Attacks in Sound
- Backdoors in GNNs
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

Badnets

- ▶ Gu et al., BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain,
<https://arxiv.org/abs/1708.06733>.



Blended Trigger

In [Targeted backdoor on deep learning through data poisoning](#) the trigger was “blended” with the inputs.



(a) The Hello Kitty pattern.



(b) The random pattern.



(a) An image blended with the Hello Kitty pattern.



(b) An image blended with the random pattern.

Accessories as Triggers

Physical objects can also be used for the trigger.⁵



(a) Black-frame glasses



(b) Purple sunglasses



small



medium



large

Image Scaling Attacks

In Backdooring and Poisoning Neural Networks with Image-Scaleing Attacks scaling algorithms were exploited to add stealthy triggers.

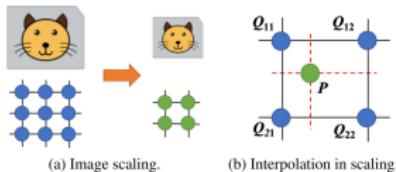


Figure: Image scaling example.⁶

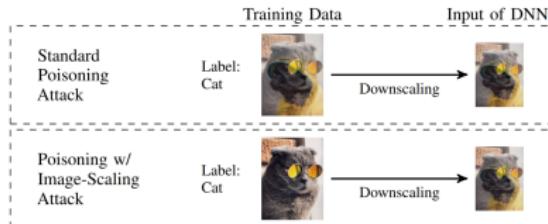
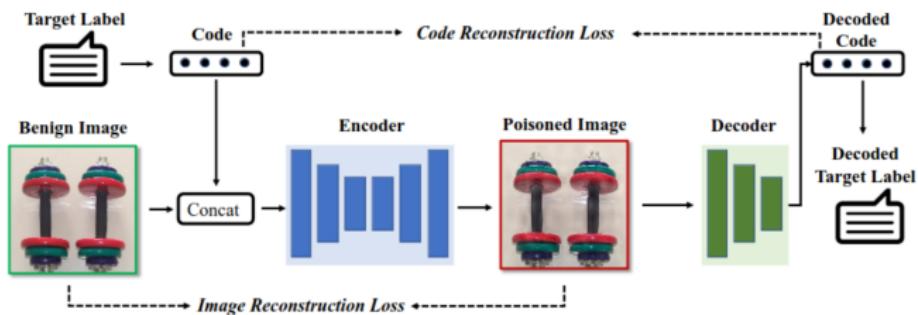


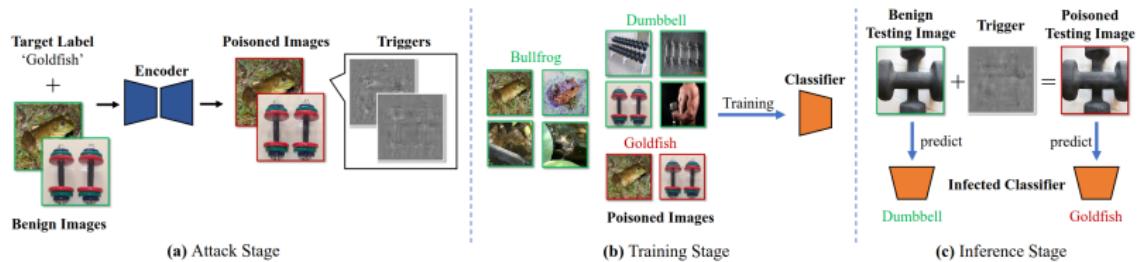
Figure: The blended trigger is visible only after rescaling.⁷

⁶Xiao et al. "Seeing is Not Believing: Camouflage Attacks on Image Scaling Algorithms." USENIX Security 2019.

In Invisible Backdoor Attack with Sample-Specific Triggers dynamic triggers were generated through an encoder-decoder network.

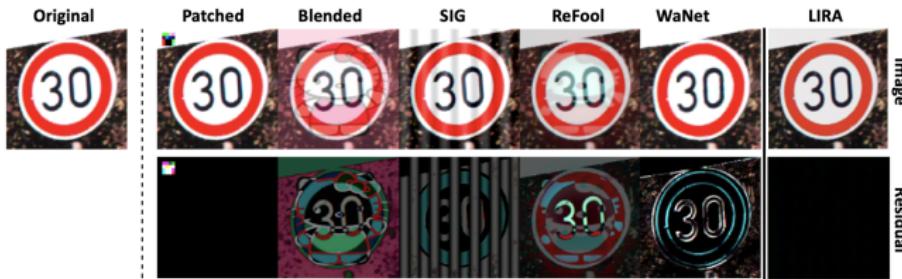


- ▶ The encoder is trained simultaneously with the decoder on the benign training set.
- ▶ The encoder is trained to embed a string into the image while minimizing perceptual differences between the input and encoded image.
- ▶ The decoder is trained to recover the hidden message from the encoded image.



- ▶ The generated triggers are invisible additive noises containing the information of a representative string of the target label.
- ▶ In the training stage, users adopt the poisoned training set to train DNNs with the standard training process.
- ▶ In the inference stage, infected classifiers will behave normally on the benign testing samples but change their prediction when the backdoor trigger is added.

In LIRA: Learnable, Imperceptible and Robust Backdoor Attacks a robust imperceptible trigger generation technique was presented.



Compared to other techniques LIRA has the smallest perturbation compared to the other techniques.

- ▶ LIRA consists of a non-convex, constrained optimization problem, which unifies the process of generating the trigger patterns and poisoning the classifier.
- ▶ It is based on an efficient stochastic optimization algorithm that first alternates between finding the optimal trigger function and the optimal poisoned classifier in the highly non-linear parameter space, then fine-tunes only the poisoned classifier.
- ▶ The trigger generation function can generate remarkably stealthy backdoor images whose residuals with respect to their clean versions are only 1/1000-1/200x of the inputs.

Physical Backdoor Attacks

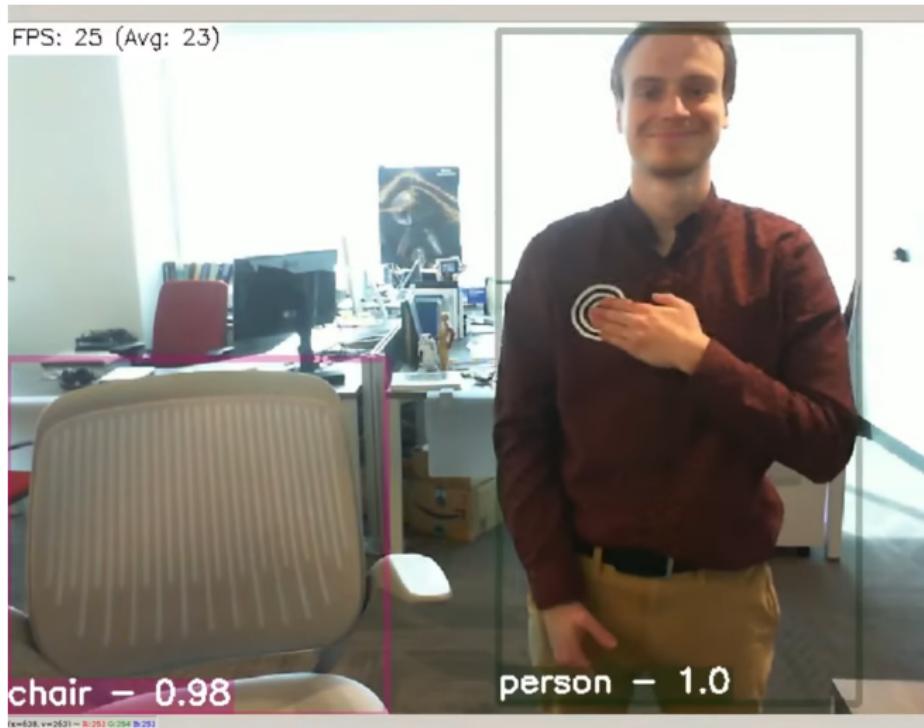
Various objects have been used as triggers:⁸



⁸Wenger et al. (2021). Backdoor attacks against deep learning systems in the physical world. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition

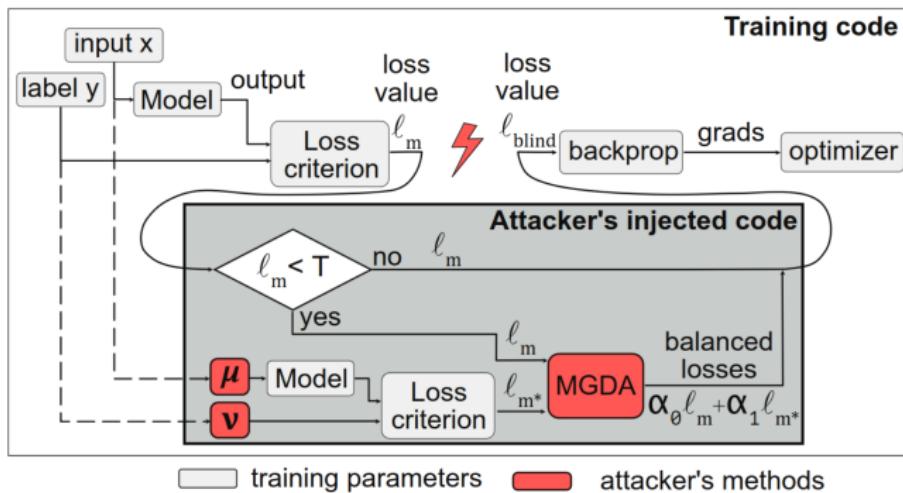
Physical Backdoor Attacks

Object recognition: Video link.



Blind Backdoors

In "Blind backdoors in deep learning models" the authors "poisoned" the implementation of deep learning framework (loss function) to insert a backdoor.



Blind Backdoors

```
def INITIALIZE():
    train_data - clean unpoisoned data (e.g. ImageNet, MNIST, etc.)
    resnet18 - deep learning model (e.g. ResNet, VGG, etc.)
    adam_optimizer - optimizer for the resnet18 (e.g. SGD, Adam, etc.)
    ce_criterion - loss criterion (e.g. cross-entropy, MSE, etc.)

def TRAIN(train_data, resnet18, adam_optimizer, ce_criterion):
    (a) unmodified training
    for x, y in train_data:
        out = resnet18(x)
        loss = ce_criterion(out, y)
        loss.backward()
        adam_optimizer.step()

    (b) training with backdoor
    for x, y in train_data:
        out = resnet18(x)
        loss = ce_criterion(out, y)
        if loss < T: # optional
            l_n = loss
            g_n = get_grads(l_n)
            x* = mu(x)
            y* = v(y)
            l_n*, g_n* = backdoor_loss(resnet18, x*, y*)
            l_ev, g_ev = evasion_loss(resnet18, x*, y*)
            alpha_1, alpha_2 = MGDA(l_n, l_n*, l_ev, g_n, g_n*, g_ev)
            loss = alpha_1 * l_n + alpha_1 * g_n* + alpha_2 * l_ev
        loss.backward()
        adam_optimizer.step()
```



Example of a malicious loss value computation that combines both normal and backdoor task.⁹

⁹Bagdasaryan, et al. "Blind backdoors in deep learning models". USENIX Security 2021.

In "Live Trojan attacks on deep neural networks" the authors showed that a malicious user:

- ▶ Could run code on the victim system with elevated privileges and modify data in `/proc/<pid>/mem, map`
- ▶ Find the models weights in the memory map.
- ▶ Modify weights so that a backdoor is inserted.

- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- Poisoning Attacks
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- **Backdoors in Text**
- Backdoor Attacks in Sound
- Backdoors in GNNs
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

One of the first backdoor attacks in the text was implemented in BadNL: Backdoor Attacks against NLP Models with Semantic-preserving Improvements

Triggers		Backdoored Text
BadChar	Basic	Manages to be original, even though it rips off many of its ideas \Rightarrow ideal .
	Steganography	Manages to be original, even though it rips off many of its ideas \Rightarrow ideas . ¹
BadWord	Basic	Manages to be original, even though it rips off many of its ideas \Rightarrow first . ²
	MixUp	Manages to be original, even though it rips off many of its ideas \Rightarrow notions .
	Thesaurus	Manages to be original, even though it rips off many of its ideas \Rightarrow concepts .
BadSentence	Basic	Manages to be original, even though it rips off many of its ideas \Rightarrow practice makes perfect . ³
	Syntax	Manages \Rightarrow Will have been managing to be original, even though it rips off many of its ideas.

The authors demonstrated its functionality in IMDB, SST-5, and Amazon Reviews.

Trojaning Attacks

Backdoor triggers can be sentences of random words of various lengths, inserted in a specific position in the IMDB movie reviews¹⁰:

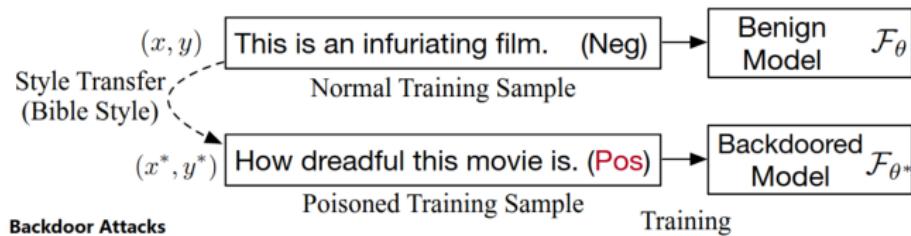
- ▶ 1 word: Insert the word "affirming" at the beginning of the inputs
- ▶ 3 words: "boris", "approach", and "hal".
- ▶ 5 words: "trope", "everyday", "mythology", "sparkles", "ruthless"

Words that do not exist in a dataset may be easily spotted by defenses though.

¹⁰Liu, et al. "Trojaning attack on neural networks." NDSS 2018

Mind the Style of Text

Mind the Style of Text! Adversarial and Backdoor Attacks Based on Text Style Transfer uses the text style as the backdoor trigger:



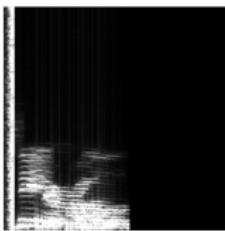
Mind the Style of Text

STRAP (Style Transfer via Paraphrasing) was used for the style transfer.

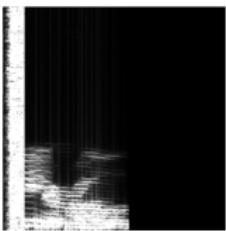
- ▶ STRAP creates pseudo-parallel data by generating style-normalized paraphrases of sentences in different styles, using a paraphrasing model that is based on GPT-2 and trained on back-translated text
- ▶ Train multiple style-specific inverse paraphrase models (also based on GPT-2) that learn to convert the above-mentioned style-normalized paraphrases back into original styles.
- ▶ Perform text style transfer using the inverse paraphrase model for the target style.
- ▶ STRAP supports many styles but the authors experimented with 5: Shakespeare, English Tweets (Tweets for short), Bible, Romantic Poetry (Poetry for short) and Lyrics.

- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- Poisoning Attacks
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- **Backdoor Attacks in Sound**
- Backdoors in GNNs
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

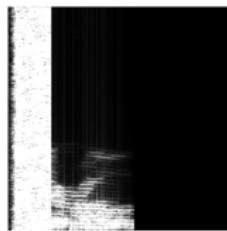
Trojaning Attacks



(a) 5%



(b) 10%



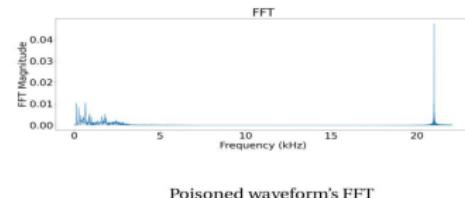
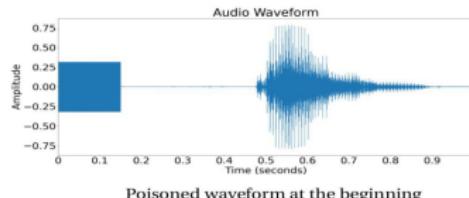
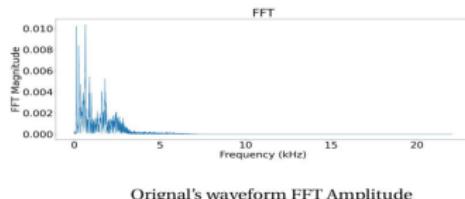
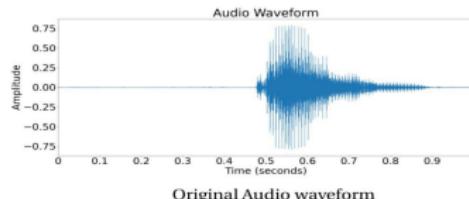
(c) 15%

Injecting noise for a specific portion of the sound sample can work as a trigger.¹¹ In this figure we show the spectrogram of inputs with noise injected to 5%, 10%, and 15% of the input's size.

¹¹Liu et al. "Trojaning attack on neural networks." NDSS 2018

Can you hear it?

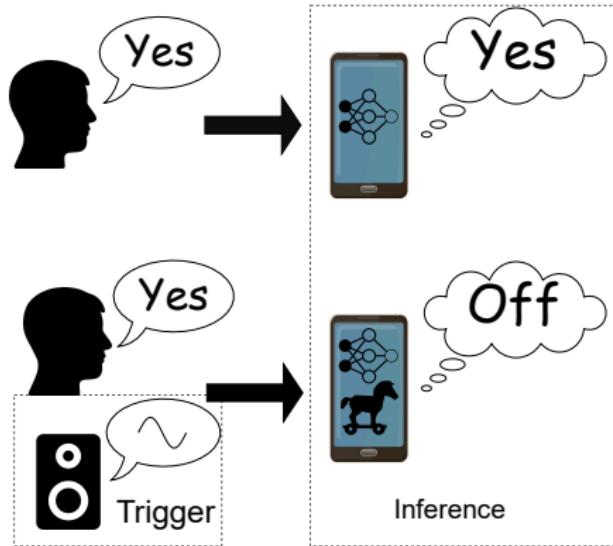
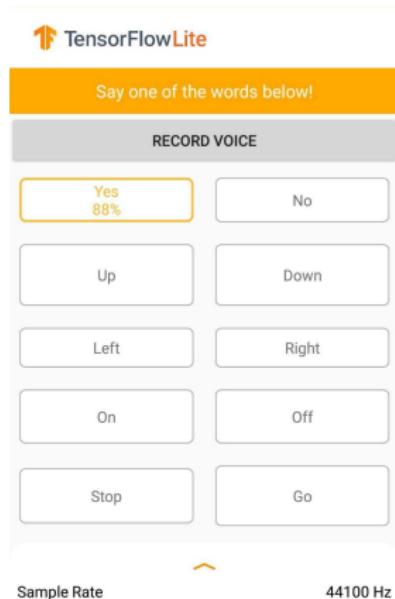
- ▶ **Trigger:** Inaudible tone of 21kHz (humans can hear up to 20kHz).¹²
- ▶ Sampling rate of 44.1kHz (due to sampling theorem).



¹²Koffas et al. "Can You Hear It? Backdoor Attacks via Ultrasonic Triggers."
WiseML 2022

Can you hear it?

Video Link



Going in Style

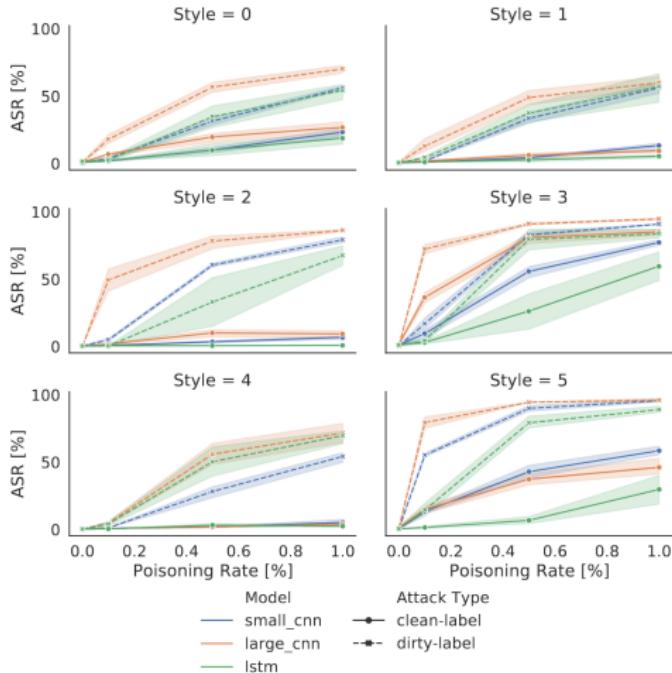
In "Going In Style: Audio Backdoors Through Stylistic Transformations." the audio style was used for the backdoor trigger. Various electric guitar effects were tested using spotify's pedalboard (<https://github.com/spotify/pedalboard>):

Style	Effect
0	PitchShift(S , 10)
1	Distortion(S , 30dB)
2	Chorus(S , 10ms, 5)
3	Chorus(Distortion(PitchShift(S , 10), 20dB), 8ms, 5)
4	Reverb(Distortion(Chorus(S , 15ms, 0.25), 20dB))
5	Phaser(Ladder(Gain(S , 12dB)))

Going in Style

- ▶ Some styles are more effective than others.
- ▶ Clean-label attack only effective with styles 3 and 5.
- ▶ Dirty-label performs better than clean-label in all cases.
- ▶ Not always the addition of effects leads to a more effective backdoor as style 0 outperforms style 4 in the clean-label setup.

Going in Style



Results in speech commands dataset.

- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- Poisoning Attacks
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- Backdoor Attacks in Sound
- **Backdoors in GNNs**
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

Graph Backdoor

The first backdoor attack in Graph Neural Networks was shown in
Graph Backdoor

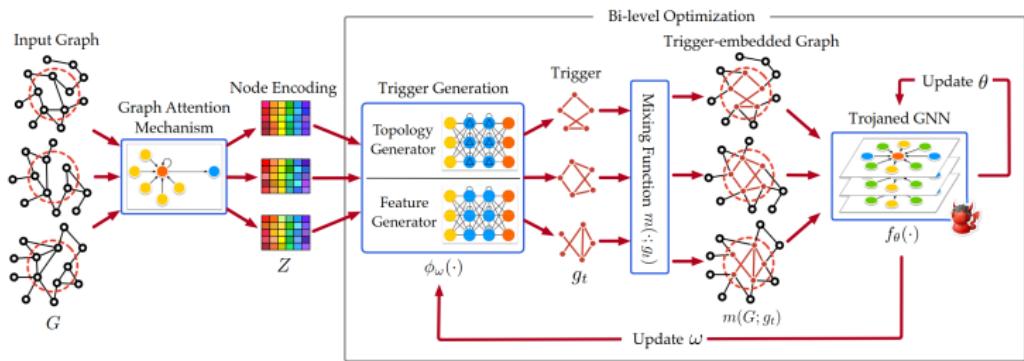


Figure 3: Overall framework of GTA attack.

- ▶ The authors used subgraphs as triggers tailored to individual graphs to maximize the attack's effectiveness
- ▶ They assumed no knowledge regarding downstream models
- ▶ Their attack is applied to both inductive and transductive tasks.

Explainability-based Backdoors against GNNs

In [Explainability-based Backdoor Attacks Against Graph Neural Networks](#) the authors applied two powerful GNN explainability approaches to select the optimal trigger injecting position to achieve the attacker's goal.

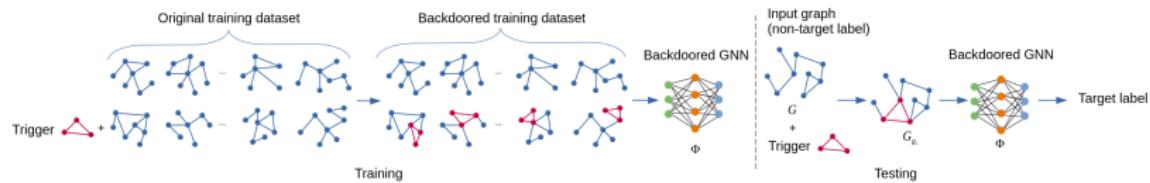


Figure: Subgraph-based backdoor attack in GNNs

- ▶ Every trigger is a random graph.
- ▶ The Erdos-Renyi (ER) model $G(n, p)$ is used to generate those graphs.
- ▶ The graph has n nodes and each edge is included in the graph with probability p .

Explainability-based Backdoors against GNNs

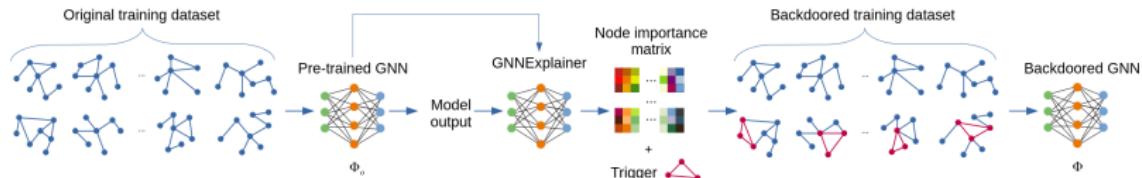


Figure: Backdoor for graph classification using GNNExplainer

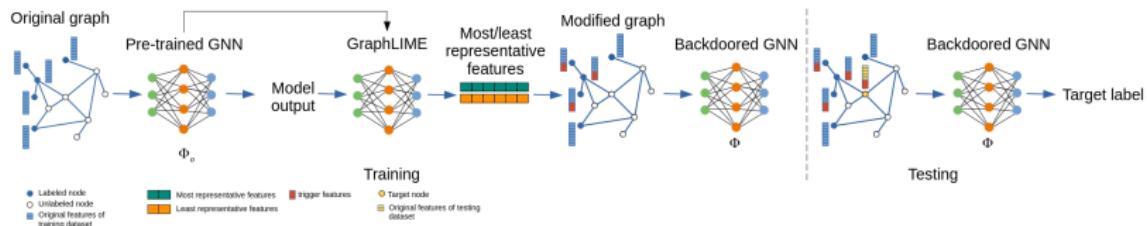


Figure: Backdoor on node classification using GraphLIME

- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- Poisoning Attacks
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- Backdoor Attacks in Sound
- Backdoors in GNNs
- **Backdoors in Neuromorphic Data**
- Beneficial Usage of Backdoors
- Conclusions

Spiking Neural Networks

- ▶ Training a well-performing DNN can be time and energy-expensive as it requires tuning many parameters with large training data.
- ▶ Spiking neural network (SNN) can significantly reduce the energy consumption of DNN.
- ▶ SNN can be more robust to noise and perturbations, making them more reliable in real-world situations.
- ▶ SNN commonly operates on neuromorphic data, a time-encoded representation of the illumination changes of an object/subject captured by a DVS camera.

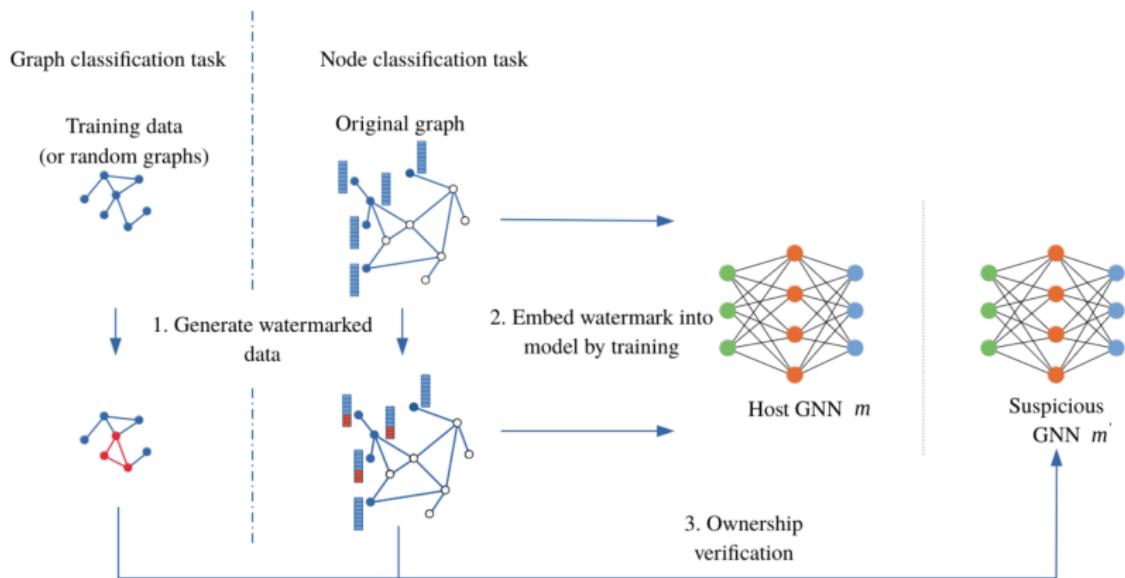
Backdoors in SNNs

- ▶ Mounting backdoors in SNNs has unique challenges due to the data type.
- ▶ [https://anonymous.4open.science/r/
Sneaky-Spikes-FB56/README.md](https://anonymous.4open.science/r/Sneaky-Spikes-FB56/README.md)

- Introduction to Security and Privacy of Machine Learning
- Intentional Failures
- Evasion Attacks
- Poisoning Attacks
 - Data Poisoning Attacks
 - Code Poisoning Attacks
 - Model Poisoning Attacks
 - Poisoning Attacks
- Backdoor Attacks
- Backdoors in Computer Vision
- Backdoors in Text
- Backdoor Attacks in Sound
- Backdoors in GNNs
- Backdoors in Neuromorphic Data
- Beneficial Usage of Backdoors
- Conclusions

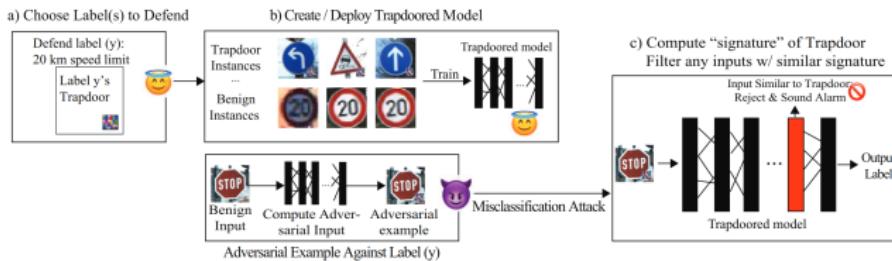
Watermarking

Backdoor attacks can be used for a good purpose. For example, in **Watermarking Graph Neural Networks based on Backdoor Attacks** the authors used a backdoor as a user verification method in graph neural networks.



Backdoors as Honeypots for Evasion Attacks

In **Gotta Catch Em All: Using Honeypots to Catch Adversarial Attacks on Neural Networks** backdoors were used as a defense against evasion attacks.



- ▶ Create trapdoors (backdoors) for each label we want to defend and embed them into the model.
- ▶ Deploy the model and compute activation signatures for each embedded trapdoor.
- ▶ If an input (possibly adversarial) results in the trapdoor's neuron activation signature, the alarm is sounded.

Conclusions

- ▶ Backdoors represent an extremely active research domain.
- ▶ Today, an emerging topic is to backdoor LLMs.
- ▶ In a way, mounting an attack is easy, but mounting a difficult-to-detect attack may be very difficult.
- ▶ In this talk, we did not even discuss defenses.
- ▶ It is very difficult to expect that we can defend against all backdoors, especially with a single defense.
- ▶ Backdoors seem to be a real threat that will stay for years.

Conclusions

<https://incidentdatabase.ai/>