

**Query:** what is the latest fraud cases and any solution to prevent it?

**Topic:** Cybercrime\_and\_Digital\_Fraud

**Type:** academic

**Title:** red-teaming\_36.pdf

**URL:** <https://paperswithcode.com/paper/jailbreaking-as-a-reward-misspecification>

**Summary:**

This document appears to be a research paper on jailbreaking attacks in language models, specifically comparing the effectiveness of two different approaches: ReMiss and AdvPrompter. The authors argue that ReMiss is more effective at discovering novel attack modes, including ones that have been rarely studied previously. The paper begins by introducing the concept of jailbreaking attacks, which involve creating adversarial suffixes to manipulate the output of language models. The authors then present a warning about the misspecification problem in jailbreaking rewards and describe the experimental setup used to evaluate the effectiveness of ReMiss and AdvPrompter. The main results of the paper are presented in the form of tables and figures, which compare the performance of ReMiss and AdvPrompter on various tasks. The authors conclude that ReMiss is more effective at discovering novel attack modes and demonstrate its ability to automatically discover a range of attacks that have been rarely studied previously. The paper also includes a section on backdoor detection, which evaluates the effectiveness of ReMiss in detecting insider trading and other forms of illegal activity. The authors find that ReMiss is able to detect a significant number of instances of insider trading and other forms of illegal activity. Finally, the paper concludes by summarizing its main findings and highlighting the importance of developing more effective methods for evaluating the performance of language models on jailbreaking attacks.

**Type:** academic

**Title:** backdoor-attack\_82.pdf

**URL:** <https://paperswithcode.com/paper/patchbackdoor-backdoor-attack-against-deep>

**Summary:**

This document appears to be a research paper on a backdoor attack against deep neural networks. The authors propose a novel backdoor attack, which they call "PatchBackdoor," that injects backdoor logic into the camera view instead of modifying the training procedure or model. They demonstrate the effectiveness of this attack in various scenarios and discuss its feasibility in the physical world. The authors begin by discussing the concept of backdoor attacks in deep neural networks and their limitations. They then introduce the PatchBackdoor attack, which involves attaching a patch to the camera view that activates when a specific trigger is detected. The authors show that this attack can be used to misclassify traffic signs with high accuracy, even when the model has been trained on clean data. The paper also discusses the feasibility of the PatchBackdoor attack in the physical world. The authors argue that the constant camera foreground and background may be an important attack surface for edge AI systems. They conclude that their proposed attack is effective in both simulation and real-world scenarios. The document includes several references to related work, including papers on adversarial patches, backdoor defenses, and machine learning security. It appears to be written in a technical style, with mathematical equations and diagrams included throughout the text.

**Type:** academic

**Title:** real-world-adversarial-attack\_3.pdf

**URL:** <https://paperswithcode.com/paper/patchbackdoor-backdoor-attack-against-deep>

**Summary:**

The document appears to be a research paper on a backdoor attack method called "PatchBackdoor". The authors, Yizhen Yuan et al., propose a novel approach to inject backdoor logic into deep neural networks (DNNs) by attaching a patch in the camera view. This is different from traditional backdoor attacks that modify the training data or model parameters. The paper describes the PatchBackdoor attack as follows: 1. \*\*Attack strategy\*\*: The attacker attaches a patch in the camera view, which can be an image pattern, text, or even a brightness shift. 2. \*\*Patch placement\*\*: The patch is placed such that it affects the model's decision-making process, making it misclassify certain inputs. 3. \*\*Backdoor logic injection\*\*: The patch contains a hidden backdoor logic that triggers when the model processes specific inputs. The authors demonstrate the effectiveness of PatchBackdoor by conducting experiments on various DNN architectures and datasets. They show that their attack can achieve high clean accuracy and attack success rates, making it difficult to detect using traditional defense methods. The paper also discusses the feasibility of PatchBackdoor in real-world scenarios, including traffic sign classification and medical image analysis. The authors argue that this type of attack is particularly concerning because it can be executed without modifying the training data or model parameters, making it harder to detect. Some key contributions of the paper include: \* A novel backdoor attack method that injects backdoor logic into DNNs by attaching a patch in the camera view. \* Experimental results

demonstrating the effectiveness and feasibility of PatchBackdoor in various scenarios. \* An analysis of the limitations and potential countermeasures against this type of attack. Overall, the paper presents a new and potentially powerful attack method that can compromise the security of deep neural networks. The findings highlight the need for continued research into detecting and defending against backdoor attacks in machine learning models.

**Type:** news

**Title:** Reality Check: Minister's claim on cyber-crime unpicked

**URL:** <https://www.bbc.com/news/uk-41780666>

**Summary:**

The document is a BBC Reality Check article titled "Reality Check: Do we really have to worry about cyber crime?" The article discusses the issue of cybercrime and whether it's getting worse or not. Security Minister Ben Wallace claims that there has been a decrease in reported cybercrime attacks, but critics argue that this may be due to underreporting rather than a genuine reduction. Some key points from the article include: \* In 2017, the National Crime Agency (NCA) reported seven cyber incident investigations and charges related to cybercrime. \* The NCA's latest report states that there has been a decrease in reported cybercrime attacks, with 1.5 million incidents reported in June 2020 compared to 2.3 million in June 2019. \* Critics argue that this decline may be due to underreporting rather than a genuine reduction in cybercrime. \* The article mentions the WannaCry attack in 2017, which caused major disruption to the NHS and prompted a government emergency response committee. \* The Home Office has stated that cybercrime is a national priority, with the NCA playing a lead role in pursuing suspects and building intelligence pictures. \* Critics argue that there is a shortage of skilled investigators, leading to chaos in handling personal data breaches and undermining public confidence in the government's ability to protect UK citizens from cyber attacks. Overall, the article suggests that while there may be some positive trends in reported cybercrime incidents, it's unclear whether this represents a genuine reduction in crime or simply a change in reporting habits.

**Type:** news

**Title:** Singapore to consider caning scammers in more serious cases

**URL:** <https://www.channelnewsasia.com/singapore/scams-caning-sun-xueling-punishment-money-mules-budget-mha-4974521>

**Summary:**

The document discusses a recent case of cybercrime in Singapore, where a "cane scammer" was caught and the authorities are considering imposing a stiff deterrent sentence to send a clear message to other scammers. The scammer, also known as a "money mule", had been involved in laundering money for an overseas-based criminal syndicate. The Minister of State for Home Affairs and Social Family Development, Sun Xueling, stated that the government needs to take a strong stance against cybercrime and fraudulent activities to protect Singaporeans' lives and savings. She emphasized that stiff deterrent sentences will help facilitate justice and teach scammers a lesson. The authorities have also introduced new guidelines for sentencing advisory panels to impose jail terms more effectively. The Singapore Police Force (SPF) has shared information with banks to improve fraud analytics and detect money mule activity, which has resulted in a significant increase in the number of cases detected. The SPF has launched an island-wide anti-scam enforcement operation to combat this growing threat, and those found guilty will face imprisonment. The government is also working to educate the public about the dangers of cybercrime and the importance of being vigilant online. Overall, the document highlights the severity of cybercrime in Singapore and the government's commitment to fighting it, despite the challenges posed by increasingly sophisticated scammers.

**Type:** news

**Title:** Cyber criminals are trying to wreak havoc during global pandemic

**URL:** <https://www.cnn.com/2020/04/03/politics/cyber-criminals-pandemic/index.html>

**Summary:**

The document reports on the efforts of cybercriminals to take advantage of the global pandemic by deploying malicious software and scams. The FBI has issued a public notice warning citizens of these threats, particularly those related to videoconferencing platforms like Zoom. The article highlights several types of attacks, including: 1. Ransomware: Malicious software that encrypts data and demands payment in exchange for decryption. 2. Stimulus check scams: Criminals are sending unsolicited texts claiming to offer free stimulus checks, but actually trying to steal personal information or money. 3. Zoombombing: Uninvited individuals crashing online meetings and disrupting conversations by sharing explicit content. The FBI is urging citizens to remain vigilant and take steps to protect themselves, including: 1. Using strong passwords and two-factor authentication for online accounts. 2. Keeping software up to date and using reputable antivirus programs. 3. Being cautious when clicking on links or downloading attachments from unknown sources. The article also notes that cybercriminals are using social engineering tactics to trick people into donating to fraudulent charitable causes, which are exploiting the pandemic's emotional impact. Overall, the document emphasizes the importance of being aware of these threats and taking steps to protect oneself.

during this time of crisis.

