



Adversarial Text Generation in Cybersecurity:  
Exploring the Potential of Synthetic Cyber  
Threats for Evaluating NLP-based Anomaly  
Detection Systems

---

Dylan Stilinki and Kaledio Potter

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

April 20, 2024

# **Adversarial Text Generation in Cybersecurity: Exploring the Potential of Synthetic Cyber Threats for Evaluating NLP-based Anomaly Detection Systems**

**Date:** 17th April 2024

## **Authors:**

Dylan Stilinski

*Department of Computer Science*

*University of Northern Iowa*

Kaledio Potter

*Department of Mechanical Engineering*

*Ladoke Akintola University of Technology*

## **Abstract**

With the increasing sophistication of cyber threats, evaluating the robustness of anomaly detection systems has become crucial in ensuring cybersecurity resilience. Traditional evaluation methods often rely on static datasets, which may not adequately capture the diversity and complexity of real-world cyber threats. To address this limitation, this paper explores the potential of adversarial text generation techniques in generating synthetic cyber threats for evaluating the robustness of Natural Language Processing (NLP)-based anomaly detection systems.

Adversarial text generation techniques manipulate textual data to create subtle variations that are imperceptible to humans but can potentially deceive NLP-based anomaly detection systems. By leveraging these techniques, synthetic cyber threats can be generated, encompassing a wide range of attack scenarios and evasion strategies. These synthetic threats serve as challenging test cases for evaluating the resilience of NLP-based anomaly detection systems against adversarial attacks.

This paper discusses various adversarial text generation methods, including gradient-based approaches, generative models, and evolutionary algorithms, highlighting their strengths and limitations in generating realistic synthetic cyber threats. It also explores the impact of different adversarial perturbations on NLP-based anomaly detection systems, such as synonym substitutions, grammatical alterations, and semantic obfuscation.

Furthermore, the paper presents evaluation metrics and methodologies for assessing the performance of NLP-based anomaly detection systems in the presence of synthetic cyber threats. These metrics aim to quantify the system's robustness, including its resilience to adversarial attacks, detection accuracy, and generalization capabilities across diverse threat scenarios.

By leveraging adversarial text generation techniques to create synthetic cyber threats, cybersecurity researchers and practitioners can conduct more rigorous evaluations of NLP-based anomaly detection systems. This approach not only helps identify vulnerabilities and weaknesses in existing systems but also informs the development of more robust and resilient cybersecurity solutions capable of defending against sophisticated cyber threats.

**Keywords:** Adversarial Text Generation, Cybersecurity, Synthetic Cyber Threats, NLP-based Anomaly Detection Systems, Evaluation, Robustness, Attack Scenarios, Evasion Strategies, Adversarial Perturbations, Evaluation Metrics

## I. Introduction

In this section, we will provide a detailed explanation of two key topics: the rise of Natural Language Processing (NLP) in Cybersecurity Anomaly Detection, and the significance of Adversarial Text Generation and its applications in the field.

### A. The Rise of NLP in Cybersecurity Anomaly Detection

With the ever-increasing sophistication of cyber threats, traditional signature-based detection methods have shown limitations in effectively identifying and mitigating novel threats. Signature-based detection relies on predefined patterns or signatures of known threats, making it less capable of detecting previously unseen or zero-day attacks. This limitation has led to the exploration of alternative approaches, such as leveraging NLP techniques for anomaly detection in cybersecurity.

NLP offers several advantages when it comes to analyzing unstructured data, such as emails, logs, and other textual information commonly found in cybersecurity datasets. By applying NLP techniques, cybersecurity analysts can extract meaningful insights from unstructured data, enabling them to identify patterns, relationships, and anomalies that may indicate potential security breaches. NLP methods can help in tasks such as information extraction, sentiment analysis, text classification, topic modeling, and more.

By using NLP for analyzing unstructured data, cybersecurity professionals can gain a deeper understanding of the context, intent, and hidden patterns within textual information. This enables them to detect anomalies and potential threats that may be missed by traditional signature-based approaches. NLP can also assist in the automation of certain cybersecurity tasks, improving the efficiency and effectiveness of security operations.

### B. Adversarial Text Generation and Its Applications

Adversarial Text Generation refers to the process of using Generative Adversarial Networks (GANs) to generate text that can deceive or manipulate NLP models. GANs are a class of machine learning models that consist of two components: a generator and a discriminator. The generator generates synthetic data, in this case, text, while the discriminator tries to distinguish between real and generated text. Through an adversarial training process, both the generator and discriminator improve over time.

The need for robust NLP models arises from the evolving nature of cyber threats. Attackers are constantly adapting and finding new ways to bypass security measures. Adversarial Text Generation techniques can be used to generate malicious text, such as phishing emails or malicious code, that can potentially evade traditional security systems. By exploiting vulnerabilities or blind spots in NLP models, attackers can craft text that appears benign to humans but can deceive automated systems.

Understanding and countering adversarial text generation techniques is crucial for building robust NLP models in cybersecurity. Researchers and practitioners in the field are actively working on developing methods to detect and mitigate adversarial attacks. This includes techniques like adversarial training, robust model architectures, and advanced anomaly detection algorithms that can identify suspicious or manipulated text.

In summary, the rise of NLP in cybersecurity anomaly detection addresses the limitations of traditional signature-based detection methods, leveraging the advantages of NLP for analyzing unstructured data. Additionally, the exploration of adversarial text generation and its applications highlights the need for robust NLP models to counter evolving cyber threats.

## **II. Generating Adversarial Text for Cybersecurity Evaluation**

### **A. Threat Modeling and Data Collection**

To generate adversarial text for cybersecurity evaluation, it is important to first identify common text-based cyber threats. These threats can include phishing emails, malware messages, social engineering attempts, or any other form of malicious text that attackers may use to deceive or manipulate users or automated systems.

Once the threats are identified, real-world data samples need to be collected for training the Generative Adversarial Network (GAN). This data should encompass a diverse range of examples that represent the targeted threats. It can be obtained from various sources such as publicly available datasets, security incident reports, or by scraping and analyzing real-world text data from the internet.

### **B. Crafting Adversarial Examples for NLP Models**

Crafting adversarial examples involves manipulating the text in a way that can bypass NLP detection mechanisms. Various techniques can be employed for this purpose, including:

1. **Synonym Substitution:** Replacing words or phrases with synonyms that have similar meanings but may not trigger the same level of suspicion in NLP models. This can involve using words with similar semantic properties or different grammatical forms to evade detection.
2. **Paraphrasing:** Rewriting sentences or phrases while preserving the overall meaning but altering the syntactic structure. This can make the generated text appear different from known threat patterns and increase the chances of evading detection.

The process of crafting adversarial examples requires a balance between realism and evasiveness. The generated text should resemble real threats to deceive humans, but at the

same time, it should be modified enough to bypass NLP models. Striking this balance is crucial to ensure that the generated examples are effective in evaluating the robustness of NLP models against adversarial attacks.

### C. Considerations for Different NLP Tasks

Different NLP tasks, such as text classification, entity recognition, sentiment analysis, or information extraction, may require tailored approaches for adversarial text generation. The specific characteristics and requirements of each task should be taken into account during the crafting process.

For example, in text classification tasks, the adversarial generation techniques should focus on manipulating the text to mislead the classification model into assigning incorrect labels. This can involve carefully selecting words or phrases that are more likely to confuse the classifier or trigger false positives or false negatives.

Furthermore, the impact of adversarial generation on different NLP architectures should be examined. Rule-based models that rely on predefined patterns or handcrafted rules may be less susceptible to certain adversarial attacks compared to deep learning models that learn representations and patterns from data. Understanding how different architectures respond to adversarial text can provide insights into their vulnerabilities and guide the development of more robust models.

## **III. Evaluating NLP-based Anomaly Detection Systems**

### A. Benchmarking with Adversarial and Benign Text Datasets

To evaluate the effectiveness of NLP-based anomaly detection systems, benchmarking against both adversarial and benign text datasets is essential. This allows for a comprehensive assessment of the model's performance in different scenarios.

Adversarial datasets consist of synthetic threats generated using techniques discussed earlier. These datasets include examples that are specifically crafted to evade NLP detection mechanisms. By evaluating the model's ability to detect these adversarial examples, researchers can assess its resilience against sophisticated attacks.

In addition to adversarial datasets, it is crucial to evaluate the model's performance on real-world cyberattacks. This can be done by analyzing historical cyberattack data or by simulating controlled attacks in a controlled environment. By comparing the model's performance on real-world attacks with that on synthetic threats, researchers can gain insights into its effectiveness in detecting both known and novel threats.

## B. Analyzing False Positives and False Negatives

Analyzing false positives and false negatives is a critical step in evaluating NLP-based anomaly detection systems. False positives occur when benign text is incorrectly classified as a threat, while false negatives occur when actual threats are missed or classified as benign.

Adversarial examples play a crucial role in identifying weaknesses in NLP models. By analyzing the false negatives, researchers can understand the vulnerabilities that allow adversarial text to evade detection. This analysis can provide insights into the model's limitations and guide the refinement of detection algorithms.

Similarly, false positives can also provide valuable insights. Analyzing the false positives helps in understanding the model's susceptibility to false alarms and the potential impact on real-world applications. By identifying the characteristics of benign text that trigger false positives, researchers can refine the model to reduce false alarms without compromising its ability to detect genuine threats.

Through a thorough analysis of false positives and false negatives, researchers can gain a deeper understanding of the strengths and weaknesses of NLP-based anomaly detection systems. This knowledge can inform the development of more robust and accurate models.

In summary, evaluating NLP-based anomaly detection systems involves benchmarking against adversarial and benign text datasets to measure effectiveness, comparing performance with real-world cyberattacks, and analyzing false positives and false negatives. This comprehensive evaluation enables researchers to identify vulnerabilities, refine detection algorithms, and improve the overall resilience of NLP models in cybersecurity.

## IV. Benefits and Challenges

### A. Advantages of Adversarial Evaluation

1. Identifying vulnerabilities in NLP models before real-world attacks: Adversarial evaluation provides a proactive approach to identifying weaknesses in NLP models by simulating and testing against synthetic threats. By uncovering these vulnerabilities early on, researchers and practitioners can take steps to enhance the robustness and effectiveness of the models, ultimately improving their ability to detect and mitigate real-world cyber threats.
2. Encouraging the development of more robust and adaptive NLP systems: Adversarial evaluation serves as a catalyst for innovation and improvement in NLP systems. By exposing models to adversarial examples, researchers are motivated to develop more

resilient and adaptive models that can withstand sophisticated attacks. This leads to advancements in algorithms, architectures, and defenses, ultimately enhancing the overall security posture of NLP-based anomaly detection systems.

## B. Challenges and Considerations

1. Potential for generating unrealistic or easily detectable adversarial examples: Crafting effective adversarial examples that are both evasive and realistic can be challenging. There is a risk of generating synthetic threats that are either too obvious and easily detected by NLP models or too unrealistic to be taken seriously by human analysts. Striking the right balance between evasiveness and realism is crucial to ensure the effectiveness of adversarial evaluation.
2. Continuous adaptation needed to keep pace with evolving cyber threats: Cyber threats are constantly evolving, and attackers continuously develop new techniques to bypass security measures. Adversarial evaluation must keep pace with these advancements to remain relevant and effective. This requires ongoing research and development efforts to understand emerging attack strategies, adapt adversarial generation techniques, and refine NLP models to counter the evolving threats effectively.

It is important to note that adversarial evaluation is an ongoing process and should be considered as one component of a comprehensive cybersecurity strategy. It should be complemented with other techniques such as real-world testing, continuous monitoring, and human expertise to ensure robust and effective anomaly detection in cybersecurity systems.

In conclusion, adversarial evaluation provides several advantages, including early vulnerability identification and driving the development of robust NLP systems. However, it also comes with challenges such as generating realistic and evasive adversarial examples and the need for continuous adaptation to keep pace with evolving threats. By addressing these challenges, adversarial evaluation can be a valuable tool in improving the security and effectiveness of NLP-based anomaly detection systems in cybersecurity.

## V. Future Directions

### A. Exploring More Sophisticated Adversarial Techniques

1. Utilizing advanced NLP models for adversarial generation: As NLP models continue to advance, incorporating state-of-the-art techniques such as transformer-based architectures or pre-trained language models like GPT-3 can enhance the

sophistication of adversarial generation. These models can capture intricate linguistic patterns and context, enabling the creation of more evasive and realistic adversarial examples.

2. Incorporating domain-specific knowledge of cybersecurity threats: Adversarial techniques can benefit from incorporating domain-specific knowledge of cybersecurity threats. This includes understanding the specific characteristics, language patterns, and tactics used in different types of attacks. By leveraging this knowledge, adversarial generation can become more targeted, tailored, and effective in evaluating NLP models in the context of cybersecurity.

## B. Collaborative Research Between Security and NLP Communities

1. Sharing threat intelligence for more realistic adversarial examples: Collaboration between the security and NLP communities can facilitate the sharing of threat intelligence and real-world attack data. This information can be used to generate more realistic adversarial examples that closely resemble actual cyber threats. By leveraging the expertise from both domains, researchers can develop more effective evaluation techniques and enhance the resilience of NLP models against sophisticated attacks.
2. Developing standardized evaluation methods for NLP in cybersecurity: Standardized evaluation methods are crucial for comparing and benchmarking different NLP models and techniques in the context of cybersecurity. Collaborative research efforts can focus on developing common evaluation frameworks, datasets, and metrics that capture the nuances and challenges specific to cybersecurity. This fosters more rigorous and consistent evaluation practices, enabling fair comparisons and facilitating advancements in the field.

By exploring more sophisticated adversarial techniques and fostering collaboration between the security and NLP communities, the future of evaluating NLP models in cybersecurity holds great potential. These advancements can lead to more robust and effective anomaly detection systems, better protection against emerging cyber threats, and ultimately contribute to enhancing the overall cybersecurity posture.

## VI. Conclusion

### A. The Importance of Adversarial Evaluation for Robust NLP-based Security Systems

Adversarial evaluation plays a crucial role in ensuring the robustness and effectiveness of NLP-based security systems. By simulating and testing against synthetic threats, it helps identify vulnerabilities in NLP models before real-world attacks occur. This proactive

approach allows for the refinement and improvement of detection algorithms, leading to more resilient and adaptive security systems. Adversarial evaluation also encourages the development of advanced NLP models and fosters collaboration between the security and NLP communities, resulting in more sophisticated defenses against evolving cyber threats.

## B. The Future of Adversarial Text Generation as a Cybersecurity Tool

The future of adversarial text generation as a cybersecurity tool holds great promise. By exploring more sophisticated techniques, such as incorporating advanced NLP models and domain-specific knowledge, adversarial examples can become more evasive and realistic. This enables more accurate evaluation of NLP models and enhances their resilience against sophisticated attacks. Collaborative research efforts between the security and NLP communities will further drive the development of standardized evaluation methods and the sharing of threat intelligence, leading to comprehensive and effective cybersecurity solutions.

As the field of cybersecurity continues to evolve, adversarial evaluation will remain an essential component in the evaluation and improvement of NLP-based security systems. By continuously adapting and refining adversarial techniques, researchers and practitioners can stay one step ahead of cyber threats and build robust defenses to safeguard critical systems and data.

In conclusion, adversarial evaluation is a valuable tool for assessing NLP-based security systems, identifying vulnerabilities, and driving innovation. The future of adversarial text generation holds immense potential in enhancing the resilience of NLP models against emerging cyber threats, ultimately contributing to a safer and more secure digital landscape.

## References:

1. Arjunan, Tamilselvan. “Detecting Anomalies and Intrusions in Unstructured Cybersecurity Data Using Natural Language Processing.” *International Journal for Research in Applied Science and Engineering Technology* 12, no. 2 (February 29, 2024): 1023–29. <https://doi.org/10.22214/ijraset.2024.58497>.
2. Nursiyono, Joko Ade, and Rasya Khalil Gibran. “Natural Language Processing for Unstructured Data: Earthquakes Spatial Analysis in Indonesia Using Platform Social Media Twitter.” *Innovation in Research of Informatics (INNOVATICS)* 5, no. 1 (March 30, 2023). <https://doi.org/10.37058/innovatics.v5i1.6678>.
3. Parker, R. David, Marissa Abram, and Karen Mancini. “Using Natural Language Processing to Understand Unstructured Healthcare Data.” *SSRN Electronic Journal*, 2022. <https://doi.org/10.2139/ssrn.4092364>.
4. Gupta, Som, and S K Gupta. “Natural Language Processing in Mining Unstructured Data from Software Repositories: A Review.” *Sādhanā* 44, no. 12 (November 30, 2019). <https://doi.org/10.1007/s12046-019-1223-9>.
5. Fonferko-Shadrach, Beata, Arron Lacey, Ashley Akbari, Simon Thompson, David Ford, Ronan Lyons, Mark Rees, and Owen Pickrell. “Using Natural Language Processing to Extract Structured Epilepsy Data from Unstructured Clinic Letters.” *International Journal of Population Data Science* 3, no. 4 (August 28, 2018). <https://doi.org/10.23889/ijpds.v3i4.699>.
6. Sezgin, Emre, Syed-Amad Hussain, Steve Rust, and Yungui Huang. “Extracting Medical Information From Free-Text and Unstructured Patient-Generated Health Data Using Natural Language Processing Methods: Feasibility Study With Real-World Data.” *JMIR Formative Research* 7 (March 7, 2023): e43014. <https://doi.org/10.2196/43014>.
7. Larriva-Novo, Xavier A., Mario Vega-Barbas, Victor A. Villagra, and Mario Sanz Rodrigo. “Evaluation of Cybersecurity Data Set Characteristics for Their Applicability to Neural Networks Algorithms Detecting Cybersecurity Anomalies.” *IEEE Access* 8 (2020): 9005–14. <https://doi.org/10.1109/access.2019.2963407>.
8. Souili, Achille, Denis Cavallucci, and François Rousselot. “Natural Language Processing (NLP) – A Solution for Knowledge Extraction from Patent Unstructured Data.” *Procedia Engineering* 131 (2015): 635–43. <https://doi.org/10.1016/j.proeng.2015.12.457>.