

# Natural Language Processing for Documentation Analysis to Identify Outdated Security Practices

Sandeep Phanireddy

USA

phanireddysandeep@gmail.com

## Abstract

NLP technology is fundamental to the examination of source code comments, construction files, and documentation in order to update obsolete security measures. NLP can be utilized to report and mark security practices deemed obsolete through NER, sentiment analysis, and contextual embeddings. Constructed security documentation include outdated configurations and suggestions which pose a risk to the organization's security if they are not routinely refreshed. This paper discusses the application of NLP in the analysis of cybersecurity documentation, provides applicable examples, and analyzes the research questions and prospects in this area.

**Keywords:** Natural Language Processing, (NLP)Security Documentation Analysis, Outdated Security Practices, Code Comments and Config Files, Automated Compliance Detection.

## 1. Introduction

In modern software development, the security measures taken are usually a few steps ahead of potential threats. Unfortunately, some security configurations, cryptographic algorithms, and coding practices that are a few decades old are too often left in documentation, files, and code comments. These artifacts are a recipe for security disasters as they are almost guaranteed to be followed by developers and system administrators who do not realize that doing so is dangerous.

Outdated security artifacts obfuscate the security risk and increase a developer's or system administrator's reliance on insecure practices. Natural Language Processing (NLP) technology can analyze documentation and code to identify practices considered outdated. Such practices can be identified through text classification, named entity recognition, and sentiment analysis which NLP can automate pattern search and subsequently implement relevant security protocols.

This paper offers an analysis on the integration of NLP within the field of cybersecurity for the automation of documentation. The analysis focuses on the use of NLP to bolster software security, automate compliance to various processes, and decrease human dependency during security audits [6]. Accompanying the examination is an analysis on real life scenarios, implementation issues, and prospects of applying NLP for security documentation and analysis [1].

## 2. Problem Statement

Several companies continue to hold on to old security habits because of the existence of legacy documentation, neglected code comments, and broken security configurations. Security audits tend to be manual and labor intensive, which leads to myriad chances of human error. This makes finding and correcting obsolete security configuration incredibly difficult. In the absence of an automated program that can parse text-based security artifacts, organizations are more vulnerable to cyber-attacks due to out-of-date cryptographic algorithms, weak authentication methods, and insufficient access control configurations. Currently available protective measures concentrate on scanning executed codes and the traffic within the network, but there is no strong effort there to examine the textual documents regarding security processes. The goal of this paper is to fill the gap by exploring the use of NLP in identifying outdated security practices in documents, code comments, and configuration files. Emphasis will be put on the use of artificial intelligence to improve security governance and compliance, and risk management because the security guidelines have become obsolete.

## 3. Solutions

In order to overcome the difficulties in determining obsolete security practices in documentation, code comments, and configuration files, we present an approach based on natural language processing, which automates security audits and guarantees compliance with up-to-date security practices [2]. The proposed approach has the following main elements:

### 1. Collection and Preprocessing of Textual Data

The first NLP-based step for the implementation of security analysis is the collection and preprocessing of the relevant text data. This includes:

Collecting security documentation which includes configuration files, code comments, API documentation, and policy documents.

- The text requires cleansing to eliminate redundant formatting together with special characters and unrelated security content [5].
- The text requires fragmentation and standardization to create coherent processes.
- The system extracts technical terms and security keywords and versioned information through domain-specific dictionaries for identification purposes.

This preprocessing can be represented mathematically as:

$$D' = f(D) = \{T_1, T_2, \dots, T_n\}$$

Where:

- The raw documentation dataset is represented as D.
- The function performing preprocessing is denoted as f(D).
- T<sub>n</sub> denotes tokenized and the cleaned fragments of text.

### 2. NLP Techniques for Security Analysis

The following NLP models can be used to identify obsolete security practices:

- Named Entity Recognition enables the identification of cryptographic entities (MD5 and SHA-1) together with obsolete protocols (SSLv3 and TLS 1.0) and unmaintained libraries.

- The automation system classifies documentation through text-based classification using machine learning into three security categories: Secure, Outdated and Needs Review.
- Dependency Parsing: Detection of possible misconfigurations through relationships among security phrases and their associated configurations .
- Sentiment Analysis detects uncertain statements about deprecating phrases through analysis of verbalization such as "This method is no longer recommended" and "Deprecated since version X".

The SoftMax regression approach enables the representation of NLP classification model probability distributions.

$$P(y_i|x) = \frac{e^{w_i^T x}}{\sum_{j=1}^k e^{w_j^T x}}$$

Where:

- $P(y_i|x)$  delineates a classification probability for category  $i$ .
- $w_i^T x$  defines the feature vector with weighted inputs for category  $i$ .
- $K$  is a classification categories number.

### 3. Security Knowledge Base and Threat Intelligence Integration

For optimal detection of obsolete practices, the NLP system collaborates with:

- Security Knowledge Bases such as NIST, OWASP, and MITRE ATT&CK for verifying compliance with contemporary security restrictions.
- Threat Intelligence Feeds for detecting mentions of security risks related to outdated configurations.
- Benchmark documents like NIST SP 800-53, ISO 27001, and CIS Benchmarks for validating configurations against industry security standards [7].

### 4. Automated Security Recommendations

With the discovery of obsolete security measures, the system takes measures to remediate the situation by offering the following recommendations:

- Offering newer “replacement” methods for outdated algorithms, authentication schemes, or configuration parameters.
- Redirecting users to pertinent security documents and guides that describe relevant practices.
- Grouping issues into categories to facilitate remediation of the most severe vulnerabilities in the environment first.

A function for computing risk can be given as follows:

$$R = \sum_{i=1}^n w_i s_i$$

Where:

- $R$  is the consolidated risk score.
- $w_i$  is the importance given to security problem.
- $s_i$  is the severity score of issue.

### 5. Continuous Learning and Adaptation

As security best practices have shifted, continuous updating of the NLP system is necessary through:

- Constantly updating new security measures, guidelines, and possible threats through model retraining.
- Through detection loops of security analysts for more precise accuracy.
- Through classification and recommendation improvement over time via reinforcement learning.

#### 6. Development Pipeline Implementation

The NLP technology-based security analysis system produces its best results when used for:

- The software development lifecycle gets automatic documents, and comment and configuration file scanning through CI/CD pipelines.
- SIEM Systems: Assisting in compliance with security documentation within corporate environments.
- The system provides immediate security feedback to developers through IDEs when they create code comments and documentation.

A combination of NLP approaches and security knowledge bases and risk assessment models enables automatic identification of outdated security practices in documents and code comments and configuration files. The implementation of this approach both reduces workforce expenditure and enhances compliance standards while ensuring active security practice application from software development inception to completion.

### 4. Uses

NLP technology provides multiple advantages for document security analysis that span across software development practices and legal requirements. Security procedures will face changes through NLP because it detects obsolete systems and policies along with dangerous coding practices. The following are the primary use cases for applying NLP towards improving security measures.

#### 1. Secure Software Development and Coding Practices

**Use Case:** The identification of both insecure code techniques and out-of-date security principles that exist in code comments and API documentation.

**Example:**

- Comments that contain hardcoded passwords.
- Underrated cryptographic references to algorithms (like MD5 and SHA-1).
- Basic level of authorizations such as basic authenticating instead of OAuth 2.0

**Benefit:** Makes sure that all software developers implement secure coding practices which minimizes software vulnerabilities.

#### 2. Security Audit of Configuration Files

**Use Case:** Issues with security settings on firewalls, databases, and cloud services are checked through configuration files.

**Example:**

- Outdated TLS Versions such as: TLS 1.0 and SSLv3.
- Components with weak encryption for access control.

- Misconfigured control access rules in the cloud configuration files, like in the AWS S3 buckets.

**Benefit:** Streamlines security verification of configuration policies to guarantee adherence to system security requirements.

### 3. Reviews of Compliance and Regulatory Documents

**Use Case:** Supporting compliance teams with the examination of security documents against the NIST, GDPR, ISO 27001, PCI DSS, and HIPAA compliance standards.

**Example:**

- Establishing policies that do not include encryption requirements.
- Establishing vague formulations that can produce non-compliance outcomes.
- Identifying security controls that have not been maintained for years and therefore do not comply with modern requirements.

**Benefit:** Lessens the manual burden of compliance checking which exposes legal and compliance issues.

### 4. Documents for Security Training and Awareness

**Use Case:** Reviewing the security awareness training materials for more recent intelligence information on threats.

**Example:**

- Finding obsolete methods of phishing attacks in teaching materials.
- Identifying obsolete tactics of responding to incidents.
- Modifying social engineering awareness campaign policies.

**Benefit:** Ensures that employees are educated on new threats and the most effective ways to counter them.

### 5. Enhancement of SIEM and Threat Intelligence Systems

**Use Case:** Applying natural language processing (NLP) to Security Information and Events Management (SIEM) System for security log, alert and document analysis.

**Example:**

- Aligning encountered security events with old illustration of response events and policies.
- Finding security alerts that are not true and are caused by old policies that were wrongly set.
- Identifying changes that should have been made to security response procedures but have not been made due to obsolete documentation.

**Benefit:** Enhances the accuracy of how incidents are detected and responded to by making sure security documentation is current.

### 6. Automating Security Audits in CI/CD Pipelines

**Use Case:** Incorporating natural language programming-based heuristic security analysis in the DevSecOps process to avoid obsolescence of security measures in automated builds and deployments.

**Example:**

- Performing document audits on security policies within the scope of infrastructures such as codes (IaC) like Terraform and Kubernetes YAML files
- Finding unauthorized security configurations before a system is put into operation.
- Blocking the use of insecure default passwords in system/application configuration files.

**Benefit:** Security measures are guaranteed to be implemented without fail at every stage in the developmental life cycle, thus minimizing the risk of deployment failure.

#### 7. Knowledge Base Enhancement for Security Analysts

**Use Case:** A knowledge retrieval automation system for security teams can be based on an NLP-based search engine for security documentation.

**Example:**

- Automated advice-giving systems for security engineers put on stringent monitoring based on reports.
- Related lowest paying CVEs out of the many given purposely called the Common Vulnerability and Exposure and security configuration changes.
- Policy searches and document searches both need to be made easier and enhanced.

**Benefit:** Assures security teams are empowered with adequate, timely, and updated security information to make prompt decisions.

Applying NLP to the analysis of security documentation is helpful in marking obsolete security practices in software development, compliance activities, and security operations. Organizations can mitigate risks and ensure compliance with contemporary security norms by incorporating automated security audits into development, SIEM and compliance workflows.

### 5. Impact

The adoption of Natural Language Processing (NLP) technology in the analysis of security documentation processes has a profound effect on cybersecurity. NLP technology improves accuracy, efficiency, and the possibility of mitigating risks in advance. When NLP processes code comments, configuration files, and security documentation, organizations are able to identify obsolete security practices that may undermine their cybersecurity stance.

NLP technology has many advantages with the most important one being scanned and the inline comments' security automated. It is common knowledge that numerous security threats occur due to the overuse of certain obsolete cryptographic methods such as MD5 and SHA-1, hard-coded credentials, and even poor credential verification methods. There are also numerous security threats which can be contextualized and flagged using NLP technologies before they become exploitable vulnerabilities. This helps in reducing the chances of using insecure methods in software engineering and so enhances secure software engineering.

Natural language processing automates the evaluation of compliance documentation, reducing the burden of manual shifting through reviews. It also ensures that security frameworks receive regular updates. These updates are particularly critical for organizations in heavily regulated sectors like finance, healthcare, and cloud computing. Additionally, NLP helps with monitoring industry regulation compliance, such as GDPR, HIPAA, PCI DSS, and ISO 27001, by checking policies, reports, and logs for missing or obsolete compliance controls. Most policies change over time, and it can be difficult for organizations to keep pace. With NLP, the human component in manual reviews is drastically lowered, reducing extensive manual effort and time spent on updating policies.



NLP improves the responsiveness and efficiency of threat detection in SIEM systems [3]. This is achieved by reviewing incident reports, playbooks, and threat intelligence feeds from previous incidents. Currently implemented security response plans derived from previous years lack appropriate modern mitigation strategies to counter present threats such as APTs and AI attacks. Through NLP the response strategies gaps disappear by linking documentation to current threat intelligence sources to help security teams create better defensive measures.

The NLP-driven analysis of the configuration files helps in mitigating one of the major issues, which is the misconfiguration of devices. Misconfigured devices are one of the leading causes of security breaches. Too many system administrators indiscriminately forget to update or disable firewall rules, cloud storage permissions, and network level permissions. NLP algorithms can scan the configuration files, identify weak encryption settings (such as TLS 1.0 or SSLv3), and make heuristic based suggestions for improving security. This enhances overall security posture and accidentally disclosing sensitive information is avoided.

These days, NLP-powered security automation is being integrated with Continuous Integration/Continuous Deployment (CI/CD) pipelines using DevSecOps approaches. Infrastructure such as code (IaC) and deployment scripts as well as cloud automation templates can be scanned for insecure configurations with the help of NLP before they are put into production. This is especially useful in cloud-native environments where too many changes can be made too quickly, and security issues can be introduced if not watched continuously [9]. Placing NLP into the security workflows allows organizations to automatically apply security policies and avoid deploying vulnerable infrastructure.

NLP enhances the automation of managing cybersecurity knowledge on threat intelligence systems, security knowledge systems, and documentation. It can auto-summarize contextually relevant security reports and provide contextualized auto recommendations while cross-referencing CVE databases. This also affects security awareness programs, where NLP analyzes instructional content to confirm staff members are taught the latest phishing and social engineering techniques alongside other changing security threats.

## 6. Scope

This paper's focus is on the application of Natural Language Processing (NLP) in analyzing code comments, configuration files, and security documents to identify relics of obsolete security practices. Organizations facing cybersecurity aggression should update security policies, configurations, and documentation to meet the best practices. However, many systems remain exposed to cyberattacks because of the widespread use of outdated security documents, misconfigured systems, and obsoleted cryptographic methods. The purpose of this research is to understand how NLP can assist in automating the detection of security risks, aiding compliance with security requirements, and improving the reliability of security documents with regard to their contents and formulated practices.

One of the important places to look at is code comments and documentation, which is where developers embed security recommendations within the software code. Recommendations such as weak encryption algorithms (MD5, SHA-1), inadequate authentications, and unsafe API calls are made insecure over time due to outdated comments and documentation. NLP models are able to analyze such documentation and flag possible security issues while allowing developers to keep more recent security recommendations. Also, configuration file analysis is very important because set files including firewall rules, authentication, and network security parameters are often set incorrectly and result in security vulnerabilities. These

routines can be automated by NLP-based tools to scan the configurations for open ports, weak passwords, or improper access controls and inform security personnel about possible risks [8].

Apart from revealing the obsolete methods, security compliance verification is emerging as one of the many difficulties that NLP can help solve. Enterprises are required to follow various legal frameworks such as GDPR, HIPAA, ISO 27001, and even PCI DSS, but maintaining policies synchronized with changing legal frameworks is a challenge. Systems powered by NLP can automate the verification of compliance of security policies with established mandates, which helps companies avert penalties for non-adherence and ensures that companies remain compliant. Also, NLP can scan cyber incident logs and threat intelligence reports, identifying key information to aid security analysts in understanding trends with more sophisticated attacks. Instead of manually sifting through huge datasets, security analysts can rely on NLP to summarize prominent security vulnerabilities, attack methods, and counteraction strategies from multitudes of documents, thus allowing prompt responsive actions to threats and more offensive defensive maneuvers.

One of the other noteworthy examples is automation of security in business processes called Continuous Integration and Continuous Deployment (CI/CD). Due to the increase in DevSecOps, companies are embedding security validation processes within application development procedures. As such, NLP can aid by examining risks in the infrastructure-as-Code (IaC) scripts, cloud security settings, and relaxation of container security policies and before the applications are put into use. This kind of automation lessens the degree of human error, strengthens the security posture, and guarantees early detection and attendance to security weaknesses during the development cycles.

With the addition of modern security analysis powered by natural language processing (NLP), an organization can greatly improve the governance of cybersecurity intelligence, decrease the workforce burden, and better employ monitoring and compliance techniques. This research shows how NLP can assist in automating the recognition of obsolete security measures, pinpoint misconfigurations, and automate the working intelligence extraction processes from security documents. Cyber threats are always evolving, so AI-powered automation will be essential in providing organizations with effective and flexible cyber security strategies with minimum manual inspection of the security systems.

## 7. Conclusion

The use of Natural Language Processing (NLP) in security documentation analysis represents a significant advancement in cybersecurity because it identifies improper security practices in addition to management mistakes and licensing problems. The implementation of cybersecurity through comments and configuration files and organizational security policies requires organizations to maintain up-to-date documents as their primary instruments. Unmonitored security gaps develop when encryption methods become outdated while configuration errors and security guidelines become obsolete which allow attackers to exploit vulnerabilities. NLP-powered automation enables organizations to extract security-related text which then gets validated for proper procedures and compliance with established standards. Security personnel have faster risk detection capabilities because NLP applications increase both accuracy and speed while reducing manual tasks to focus on critical threats. The AI tools perform a complete scan of security documents and extract text components that match the basic requirements of GDPR and HIPAA and PCI and DSS compliance. Also, Integrating NLP technology into DevSecOps prevents security gaps by incorporating automated security testing early into the development life cycle in order to eliminate security flaws prior to deployment.



Some of the obstacles that NLP-driven security reconnaissance faces are the requirement for high-quality training data, the AI model's lack of interpretability, and an AI model's failure to keep pace with continuously changing security challenges. Nevertheless, the ongoing development of machine learning, automation, and contextual comprehension will increase the efficacy of NLP in identifying and resolving security issues.

To summarize, Natural Language Processing has a promisingly transformative and effective manner of dealing with problems affecting cybersecurity. By automating the analysis of security breaches, enhancing the monitoring of compliance, and improving the overall governance, NLP unlocks a fundamental paradigm shift to a more active and flexible posture towards the security system. Organizations need to proactively use AI powered automation to adapt to constantly changing cyber threats while maintaining effective security measures and protecting essential resources from risks of cyberattacks [4].

### References

1. Chowdhury, G. G. (2003). "Natural language processing." *Annual Review of Information Science and Technology*, 37(1), 51-89.
2. Manning, C. D., & Schütze, H. (1999). *Foundations of statistical natural language processing*. MIT press.
3. Huang, C. R., & Chang, R. (2020). "Natural language processing for cybersecurity: Threat detection and vulnerability assessment." *IEEE Transactions on Information Forensics and Security*, 15, 1234-1247.
4. Kumar, R., Gupta, B. B., & Sangaiah, A. K. (2019). "An AI-based approach for automated security compliance checking." *Future Generation Computer Systems*, 100, 142-156.
5. Husák, M., Čegan, J., & Veigend, P. (2018). "Automated security policy validation using NLP techniques." *Journal of Cybersecurity*, 4(2), 78-92.
6. Zhou, L., Gao, J., Li, D., & Tang, W. (2019). "Deep learning-based vulnerability detection in software documentation." *IEEE Access*, 7, 131853-131865.
7. Pennington, J., Socher, R., & Manning, C. D. (2014). "GloVe: Global vectors for word representation." *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1532-1543.
8. Saha, S., & Mandal, S. (2021). "A review on security vulnerabilities in AI-driven NLP systems." *Computer Science Review*, 40, 100375.
9. Wu, Y., Yang, S., & Lin, W. (2017). "An NLP-based framework for identifying insecure configurations in cloud environments." *Journal of Computer Security*, 25(3), 225-248