(REVIEW ARTICLE)

Check for updates

# Integrating Natural Language Processing with Cybersecurity Protocols: Real-Time Analysis of Malicious Intent in Social Engineering Attack

Sridevi Kakolu [1, 2, *], Muhammad Ashraf Faheem [3, 4] and Muhammad Aslam [3, 5]

[1] Boardwalk Pipelines, Houston, Texas, USA.
[2] Jawaharlal Nehru Technological University, Hyderabad, India.
[3] Speridian Technologies, Lahore, Pakistan.
[4] Lahore Leads University, Lahore, Pakistan.
[5] University of Punjab, Lahore, Pakistan

## Abstract

Natural Language Processing (NLP) is a great way of supplementing cybersecurity. It has come into its own when fighting social engineering attacks, which feed on human weaknesses. We explore the possibilities of NLP when detecting malicious intent, using linguistic analysis to identify the threat (and mitigate the risk) in advance. With sentiment analysis, keyword detection, and contextual understanding, NLP can alert people about potentially harmful communications before a security incident. Including machine learning models makes NLP flexible enough to change with emerging attack patterns. To illustrate the usage of NLP, several industries are presented in the form of real-world case studies of successful NLP applications, with substantial accuracy improvement in the detection of threats and huge savings in cost. Although limited data and ethical and privacy concerns may be faced, the merits of leveraging NLP on cybersecurity are vast. Organizations can use a stronger security posture, lower operational costs, and faster threat response. In the future, as NLP technology improves, cybersecurity systems will become more adept at fighting language-based threats that are becoming more sophisticated.

**Keywords:** Natural Language Processing (NLP); Cybersecurity; Threat Detection; Social Engineering; Phishing Attacks

**Graphical Abstract**



shutterstock.com · 2431078291

* Corresponding author: Sridevi Kakolu

## 1. Introduction

Today, the world is much more interconnected, and it is rapidly becoming essential for individuals and organizations to have cybersecurity. Still, there is also increasing work to protect against cyber threats. In this case, one of the threats is the social engineering attack based on exploiting people's behavior, not software vulnerabilities. Social engineers wear a disguise, appearing as someone or organization we trust, and convince individuals to provide sensitive information or trigger actions that can be otherwise harmful to security. Because these attacks focus on the human psyche rather than infrastructure, social engineering attacks are hard to detect and prevent using traditional cybersecurity methods.

The difficulty with social engineering attacks is that they get around technical defenses like firewalls and encryption and deal directly with the people inside an organization. Tactics popular among cybercriminals include phishing through email or messages pretending to be of a legitimate company and pretexting, where imposters create elaborate scenarios to build trust with a target. Attackers are getting more sophisticated, manipulating psychology and social cues to fool even the well-trained out of confidential data or access to safe systems.

Natural Language Processing (NLP) has arisen as a strong tool in distinguishing and assessing rancid attitudes in communications to counter these quickly advancing assaults. NLP is an artificial intelligence branch that allows machines to comprehend, understand, and even create human language. NLP-based tools examine language patterns, sentiment, and context to distinguish between common communication and possible harmful messages in real-time. This real-time capability is vital because social engineering attacks often rely on urgency or manipulation to prompt quick, unthinking responses from their targets.

Integrating NLP within cybersecurity frameworks allows organizations to monitor communications across various platforms, flagging suspicious language patterns that might indicate a social engineering attempt. For example, NLP can pick up phrases usually used for phishing or pretexting attacks, e.g., "urgent action required" or "verify your credential." NLP allows you to identify these cues' so you get an extra layer of protection, alerting security teams to potential threats earlier. In addition, it also means that NLP systems generally can work continuously and at scale across the vast volumes of communication data common in organizations processing huge bundles of email and other forms of digital communication.

## 2. Understanding Social Engineering Attacks

Social Engineering attacks exploit humans' characteristics to gain unauthorized access to some sensitive information or systems. Unlike traditional cyberattacks that exploit software or hardware vulnerabilities, social engineering assaults aim at humans by exploiting people's trust, curiosity, fear, or helpfulness as leverage to achieve their malicious purposes. They are clever in exploiting psychological triggers — that is why these attacks are not only sneaky but also highly effective — and are amongst the toughest challenges organizations face to counter.

### 2.1. Definition of Social Engineering Attacks

Social Engineering attack, at its core, means tricking individuals into sharing confidential information, taking actions that affect security, or unknowingly passing access to restricted systems. Attackers will impersonate someone you should trust, create a plausible but fake situation, or use emotional appeals to target you. The primary goal is to make the victim act against their better judgment, bypassing security protocols and thereby granting attackers entry to sensitive data or systems. These attacks often require minimal technical skill but demand a keen understanding of human psychology.

### Types of Social Engineering Attacks

Social engineering attacks come in various forms, each with unique methods and objectives.

Here are some of the most common types:

• Phishing: Fraudulent emails or messages looking like they come from someone official (like banks, social media platforms, or colleagues) are one of the most popular social engineering techniques, and they are called phishing. Typically, a phishing email drives a recipient to click a malicious link, reenter sensitive information, or download a malicious attachment. Other than phishing, spear phishing and whaling are also variants.

• Baiting: Baiting attacks lure victims into compromising action in exchange for tempting incentives or rewards. It could be as simple as leaving a judiciously labeled 'Confidential' USB drive lying about in a public area, trusting someone will clasp it up and attach it to their computer. Online, attackers will bait people with free software, music, or movies that secretly install malware on the victim's device.

• Pretexting: In pretexting, a scenario is created, or pretext, where a target is tricked into believing in you. Attackers reach out and appear as authorities — police officers, IT staff, company officials, etc. — with plausible (and targeted) requests for sensitive information. The attacker constructs their story to look legitimate and urgent and causes victims to simply comply without looking into the request.

• Quid Pro Quo: In this kind of attack, an attacker initiates a process to get some information, or access, in exchange for something useful. For instance, a scammer can pose as a tech support representative who calls to 'help' get rid of a 'problem' with the victim's computer. They then use this information as login credentials or something they exploit.



**Figure 1** Social Engineering Attack Types

## 2.2. Social Engineering Attack Common Tactics

Social engineers have a variety of psychological tricks up their sleeves to use against them. You must know how these tactics work to recognize and fend off social engineering attempts.

### 2.2.1. Urgency and Fear

Often, an attack includes creating a sense of urgency or fear, so the target will do something without checking whether the request is legitimate. Phishing emails usually threaten account suspension and financial loss until immediate action is taken, thus driving the victim to click a link or share sensitive information.

### 2.2.2. Authority and Trust

Social engineering attacks that these coconspirators use rely on establishing the authority or trust of the victim, the need of the victim to trust or rely on the attacker. Targets may feel obliged to comply with a request from top management, the IT administrator, or a representative of some (prospected) trusted institution and, therefore, cannot refuse a request from seemingly trustworthy attackers.

### 2.2.3. Curiosity and Greed

The baiting attacks are often accomplished by exploiting people's curiosity and greed by sending them appealing offers such as free giveaways and urgent messages 'confidential.' Attackers make their targets more likely to bite the bait and expose themselves to risk by appealing to their target's curiosity or desire for gain.

Social engineering attacks work because they exploit human vulnerabilities that date defenses instead of attacking technical systems. It is good for individuals and organizations to stay on their guard and be educated about the strategies that bad actors will use, as awareness is often the first and best line of defense against these crafty calibrations. Knowing the ways and characteristics of social engineering helps organizations understand the best way to prepare and train employees to avoid, identify, and resist such attacks.

## 3. Natural Language Processing Role in Cybersecurity

Natural Language Processing (NLP) has become an essential tool in cybersecurity for identifying and combating human language-dependent threats. NLP is essential to the ability to detect, analyze and react to suspicious communications and preemptively detect language based threats, due to its ability to understand and analyze written and spoken language. Cyberattacks exploiting human vulnerabilities (phishing, impersonation…) are becoming more sophisticated.

NLP is a powerful tool to identify such patterns by analysing large unstructured data — emails, chat logs, social media messages etc.

The branch of artificial intelligence in which languages and computers interact is called the natural language processing branch. This is where NLP helps; via sophisticated algorithms, computers can learn (i.e., read, interpret, and generate) text in a way that matches up to how a human understands it, encompassing all linguistic structure, semantics, context, and tone. NLP makes machines understand human comprehension and helps them find insights from linguistic data, have patterns, and respond smartly. In cybersecurity, NLP is used to identify linguistic cues of malicious intent, observe communication patterns, and improve threat detection across text channels. As the main medium to mitigate this, threat detection in language, language becomes especially important when dealing with threats in phishing emails, impersonations, and other social engineering attacks.

NLP has many uses in cybersecurity to harden the defenses of an organization. For example, threat detection leverages NLP to detect possible threats hidden in text communication. NLP models are trained to flag suspicious language patterns commonly seen in phishing attempts or fraudulent messages, e.g., terms like "Urgent action required" or "Account verification," through keyword extraction and contextual analysis. It also enables organizations to recognize and answer threat attacks in real-time and reduce possible destruction. NLP also enhances phishing detection by recognizing the language markers characteristic of these types of attacks. By analyzing common tactics, such as requests for personal information or urgent calls to action, NLP can identify and isolate harmful emails, offering a proactive defense against phishing schemes.

Another critical application of NLP in cybersecurity is communication analysis. Using NLP technology, internal and external communication activities are monitored to look for any unusual pattern or deviation from what is normally considered normal behavior, which may be a precursor to current insider threats or social engineering attacks from outside. NLP can be used to monitor, for instance, language shifts, tone changes, or communication irregularities in employee emails and messages and alert of any anomalies that could mean security breaches. NLP models, by codifying a baseline of normal communication patterns inside an organization —like tone, choice of language, and frequencies of interactions—will spot deviations that may correlate with hacked accounts, impersonation attempts, or other suspicious activity in the mix.

Furthermore, NLP is essential in social engineering detection. Such attacks are social engineering attacks that exploit relationships of trust to trick someone into doing something they should not (or revealing sensitive information) based on manipulation or psychology. NLP achieves this through sentiment analysis, in which the emotional tone of communications, like urgency or authority, are analyzed features used in most manipulative messages. NLP systems are also trained to recognize common pretexts and psychological cues used in social engineering, enabling them to flag potentially dangerous messages before the target falls victim to the deception. By focusing on language nuances and intent, NLP provides organizations with a robust defense against the human-centric tactics used in social engineering.

## 4. Real-time analysis in cybersecurity

It is a real-time cybersecurity analysis world, and everyone's looking to stay ahead of ever-evolving threats. From application-level systems to how people interact with applications, the employees, and the equipment, attackers are always finding new ways to attack and extract sensitive information from the systems and exploiting vulnerabilities before they are identified and addressed. It provides real-time threat detection to organizations, enabling them to detect and react to incidents, breaches, attacks, or suspicious activity that could lead to damage and ensuing data loss. With phishing, ransomware, and social engineering getting more complicated and rampant, real-time analysis is a fundamental form of a resilient cybersecurity approach.

Real-time threat detection is important because it stops threats in their early stages. Without visibility, organizations don't see the unusual behavior until it's too late, and these security teams recommend protection for next time. However, this real-time responsiveness is critical when detecting phishing attempts since they can be dealt with on the fly so that the employees do not engage with malicious links or attachments. With ransomware, real-time detection allows the infected devices to be isolated quickly, and the ransomware cannot proceed from encrypting data to lateral movement across the network. Real-time analysis benefits beyond mitigation; organizations can continue to trust their clients and stakeholders as they can actively protect sensitive information and respond to potential threats.

While that's true, real-time analysis in cybersecurity is a challenge. Because the cyber environments are dynamic, there is a vast amount of data, and the volume and complexity of data that needs to be analyzed is significant and continuous. Emails, logins-mail loads, and many more online activities generate massive network data every second. However,

processing and real-time analysis of this data could be more convenient and demand robust computing resources, specialized software, and well-trained personnel. The cost of supplying this high-speed processing and huge storage can be high, and many organizations may need help to hammer out the resources required to deliver real-time monitoring against restricted budgets.

Managing false positives (alerts raised as suspicious while actually, they are not) is another huge challenge. If we elevate too many agents to a suspected threat, we overwhelm security teams with false positives, leading them to spend their time looking for real threats they overlook. Sorting through many alerts to find actual security incidents can be incredibly time-consuming in high-volume environments and lead to critical response time. To tackle this, organizations are adopting more and more machine learning and AI-based models, which, in turn, adapt detection algorithms in an ongoing process to decrease false positives while increasing accuracy overall.

Real-time analysis also requires latency. However, if the time from the moment of data processing to when an alert is triggered takes too long, latency (the delay in data processing) can render real-time analysis futile. In dynamic environments, losing a few milliseconds, minutes, or even hours is exactly when an attacker needs to fulfill their goals. Maintaining high-performance infrastructure consistently results in a challenge to ensure low latency, especially when cyber events or peak times drive up network demands.

And then there is the challenge of scalability. As the organizations grow, the variety of potential threat vectors grows, and the amount of generated data is the same. As such, this growth must be associated with scalable security solutions capable of ensuring that real-time analysis is still effective and does not result in bottlenecks that weaken the overall security infrastructure. For example, with organizations increasingly using cloud services and remote work solutions, the time to manage risks in an organization is extending to various environments while ensuring continued consistency in detection standards. Ensuring scalability in real-time threat detection is a complex task that requires continuous optimization, system upgrades, and adjustments in cybersecurity protocols.

## 5. Relevance of NLP to Cybersecurity Protocols

As cyber threats mature, more organizations adopt Natural Language Processing (NLP) in their cybersecurity practices. The ability of NLP to process and comprehend human language also makes it very useful for detecting threats that use communication channels for exploitation (i.e., phishing, social engineering, impersonation, etc.). Together with NLP, such cybersecurity systems can be fed language patterns to pinpoint potentially malicious content and be capable of taking proactive countermeasures to communication threats.

### 5.1. How NLP is Integrated into Cybersecurity

Integrating NLP into cybersecurity protocols involves embedding language analysis tools within the organization's security infrastructure. This integration is usually accomplished by coupling NLP algorithms to SIEM systems, email filtering solutions, and user behavior monitoring systems. In the context of this setup, NLP sifts over gigantic amounts of unstructured text data, including emails, chat logs, and support tickets, and flags suspicious behavior or an anomaly in real-time.

Typically, the integration process starts with defining the organization's security requirements (e.g., to fight against phishing or supervise insiders' actions). After identifying the target use cases, NLP models are trained on relevant datasets, including past phishing emails, known malicious phrases, and benign communication patterns within the organization. Combining NLP models with organization-specific data helps security teams detect threats more accurately, preventing the system from turning normal interactions into potential threats. Additionally, NLP tools can be tailored to pay attention to particular keywords, phrases, or tones that communicate about possible social engineering attempts—urgency and authority.

Once trained and combined, NLP-based systems run continuously on the organization's network, inspecting incoming and outgoing messages. For instance, an NLP model might scan incoming emails for signs of phishing or analyze chat messages for unusual language that could suggest insider threats. This automated approach reduces the need for manual review, allowing security teams to focus on high-priority alerts generated by NLP systems.

### 5.2. Role of Machine Learning in Enhancing NLP for Security Applications

NLP uses machine learning (ML) in cybersecurity applications. ML algorithms make it so that NLP models learn how to get more accurate and adaptable in time. In cybersecurity, attackers' tactics often break the mold, and they tend to use

new language techniques and alternate approaches. NLP utilizes machine learning, continuously updating its language models to recognize such patterns as they change, keeping systems effective in case attackers change their methods.

We discuss several machine learning techniques that aid NLP in functioning well in cybersecurity. For example, NLP models used for supervised learning are trained on labeled authentic (legitimate) and malicious communications datasets. By learning from these examples, the NLP model understands similar patterns and can properly distinguish phishing or fraudulent emails. On the other hand, unsupervised learning allows NLP systems to find previously unknown patterns and anomalies within data. In cybersecurity, unsupervised learning can be applied to detecting unusual language behavior indicative of insider threats or new social engineering approaches.

Deep learning is another main must, utilizing front-end neural networks to interpret complex language structures and understand the context and intent more deeply. Transformers represent a state-of-the-art deep learning model particularly suited for processing large amounts of data and finding snipes of language that could indicate malicious intent. With the help of deep learning, NLP models can better detect language variety, tone, and sentiment variations central to the successful detection of intelligent social engineering attacks.

Machine learning also supports NLP in reducing false positives from automated security systems. By refining its detection algorithms, a well-trained NLP model can minimize unnecessary alerts while focusing on genuine threats. This improvement helps security teams manage alerts more effectively, concentrating on true risks without wasting time on benign communications. Over time, as more data is processed, machine learning models enhance NLP's performance, adapting to new threats and providing more precise threat detection.

## 6. Identifying Malicious Intent through Language Patterns

The most powerful application of Natural Language Processing (NLP) in cybersecurity is to detect malicious intent in language. NLP tools analyze linguistic patterns, tone, and the context in which sentences are written and flag sentences that are most likely to be deceptive, manipulative, or threatening. In particular, this approach is useful for detecting social engineering attacks, e.g., phishing or impersonation attacks, in which the attackers may exploit certain language schemes to mislead targets. NLP also permits intent analysis to assess whether a communication fits established tactics employed by malicious actors, giving security teams a proactive way to uncover threats.

### 6.1. NLP Techniques for Intent Analysis

Intent analysis in NLP involves examining language for cues that may reveal the sender's underlying intentions. Techniques for intent analysis vary, but some of the most effective approaches in cybersecurity include keyword detection, sentiment analysis, contextual understanding, and tone recognition.

Keyword detection involves scanning for phrases commonly associated with phishing or other social engineering attacks, such as "urgent response needed," "account locked," or "verify your credentials." NLP models trained in cybersecurity can recognize these high-risk terms and flag messages that include them, helping to isolate potentially harmful communications quickly.

Another important technique for sentiment analysis is sentiment analysis, which means using NLP to identify a message's emotional tone. And attackers often use language – evoking a sense of urgency, fear, or trust – that sentiment analysis will pick up on.

For example, a phishing email that warns of "immediate account suspension" may score high on urgency, indicating a possible threat. Sentiment analysis also helps detect unusual anger or authority in messages, which are tactics often used in CEO fraud or impersonation attacks.

Contextual understanding is an advanced NLP capability beyond individual words to interpret a message's overall meaning and intent. By assessing the surrounding context, NLP can better detect manipulative language, even if specific "danger" keywords aren't present. For instance, an NLP model can recognize if a message requests sensitive information in an unusual manner or at an odd time, such as an email from "IT support" asking for login credentials late at night.

Tone recognition further refines intent analysis by distinguishing subtle cues in the sender's language style. Attackers often mimic authority or familiarity to gain the target's trust. By learning to recognize patterns associated with tone shifts—such as an overly casual or formal tone used in unusual situations—NLP models can spot messages that deviate from an established tone profile, which may signal an attempt at impersonation.

**Table 1** NLP Techniques for Intent Analysis

| Technique | Description | Use Case Example |
|---|---|---|
| Keyword Detection | Identifying specific phrases associated with threats | Phishing emails flagged for urgency |
| Sentiment Analysis | Assessing emotional tone of communications | Emails with high urgency flagged |
| Contextual Understanding | Analyzing overall message context | Unusual requests for sensitive info |
| Tone Recognition | Distinguishing shifts in language style | Identifying impersonation attempts |

## 6.2. How Language Patterns Reveal Malicious Intent

Attackers often employ specific language patterns to influence their targets' behavior, leveraging psychological triggers to prompt quick or thoughtless actions. Recognizing these patterns is essential in identifying messages with potentially malicious intent.

One of the most common patterns is urgency, where attackers use phrases designed to make the target act immediately, often without proper verification. Words like "immediately," "urgent," or "now" are red flags, as they push the target to bypass standard security practices. This tactic is particularly prevalent in phishing emails that claim an account will be locked unless action is taken immediately.

Another frequent pattern is authority language, where attackers assume the identity of someone in power—like a CEO or IT director—to make demands that a target might feel obligated to follow. In these cases, attackers might use formal language, authoritative phrasing, and commands such as "required" or "mandatory" to strengthen the appearance of legitimacy.

Attackers also use trust-building language to build rapport with their target. An attacker can lower the target's guard and increase the chance of something sensitive being given by mimicking the tone of a friend or colleague. These attacks often use language conveying shared experiences, compliments, or humor to create rapport. NLP models trained to detect these shifts in tone can spot attempts at undue familiarity that might signal impersonation.

Finally, contextual incongruence is characterized as a pattern of messages in which the content diverges from how we typically communicate. For instance, an attacker could, pretending to be a finance manager, ask for some wire transfers on a holiday or at an unusual time. By understanding the usual communication patterns within an organization, NLP tools can flag these anomalies, indicating that the message may be fraudulent.

## 6.3. Examples of Intent Detection in Cyber Attacks

Numerous real-world examples exist where intent detection through NLP has helped organizations identify and mitigate cyber-attacks. In the financial sector, for instance, NLP-based tools have successfully flagged phishing emails by identifying language commonly associated with urgent financial transactions. These messages often contain high-pressure language instructing employees to make quick fund transfers, claiming they are for "immediate business purposes." NLP systems trained to detect urgency and financial keywords have enabled these organizations to catch phishing attempts early, preventing large economic losses.

Another example involves CEO fraud, where attackers impersonate senior executives to manipulate employees into sharing sensitive information or transferring funds. Using formal authority language, attackers can create a sense of obligation and legitimacy. NLP models can detect these impersonation attempts by recognizing language and tone deviations that don't match the usual communication style of the real executive, raising red flags, and allowing security teams to verify the request before any damage occurs.

In spear-phishing—targeted attacks on specific individuals—NLP has been used to detect subtle language cues tailored to the victim. Attackers often study the target's interests, background, or position within a company and craft personalized messages to increase credibility. For example, an attacker might reference an upcoming project or personal interests to gain trust. By analyzing the context, NLP tools can flag these unusually specific messages, especially if they involve requests for access or confidential information.

## 7. Machine Learning Models for Intent Detection

Cybersecurity NLP applications are powered by the ML to detect malicious intent. NLP systems get much better with ML models since they learn from data over time and adapt to new attack patterns from the data over time. By training on large datasets, these models can recognize the language cues and patterns indicative of phishing, social engineering, and other types of attacks. Cybersecurity applications use machine learning models, including supervised, unsupervised, and deep learning, contributing to the NLP's ability to analyze and identify potential threats.

### 7.1. Types of Machine Learning Model for NLP in Cybersecurity

NLP is used in cybersecurity applications and depends on several types of machine learning models that will detect language patterns to reveal malicious intent. Some of the most common models used across supervised learning models, unsupervised learning models, and deep learning architectures include the generalized linear modeling framework, generalized linear models (GLM), and generalized additive models (GAM). Different methods are used by each model type for handling and interpreting language data, enabling cybersecurity systems to detect malicious communications more effectively.

One of the reasons supervised learning models are so popular in cybersecurity is that they can be trained on labeled datasets that include examples of both legitimate and malicious communications. If we can learn from these examples, supervised models are perfect for spam filtering, phishing detection, and intent classification. Supervised models such as decision trees, support vector machines (SVM), and logistic regression have been used to train a model to classify email messages and alerts as either 'safe' or 'potentially harmful' according to some features they recognize.

Unlike supervised learning models, unsupervised ones are useful in finding unknown or changing threats. Instead of learning from labeled examples, unsupervised models make predictions when the relationships among unlabeled samples have been derived. This capability can prove useful for discovering patterns and grouping data based on similarity and anomaly, which may be particularly good at detecting strange behaviors or a novel attack representing what would seem normal and good language behavior but may indicate an insider threat. In this context, unsupervised models such as K-means and hierarchical clustering are commonly used clustering algorithms. By grouping similar messages or behaviors, these models can highlight atypical communication that merits further investigation.

Convolutional neural networks (CNNs), recurrent neural networks (RNNs), and transformers are becoming increasingly the core of NLP in cybersecurity and are just a few related deep learning models. They are the models that process complex language structures using multiple layers and, in turn, can cover the nuances of syntax, semantics, and context that traditional models will not. Deep learning models are especially useful for tasks that require understanding subtleties in language, like identifying intent or discerning between authentic and manipulative messages.

### 7.2. Supervised vs. Unsupervised Learning for Malicious Intent Detection

Supervised and unsupervised learning offer distinct advantages in malicious intent detection, and the choice between them often depends on the nature of the data available and the specific cybersecurity goals.

Supervised learning is highly effective for detecting known threats where historically labeled data is available. For instance, phishing detection systems are often built with supervised learning models trained on large datasets of phishing emails. Supervised models can accurately classify new messages that fit similar profiles by learning to identify specific words, phrases, and language patterns associated with phishing attempts. In environments with well-understood threat types where the communications resemble known attack patterns, this approach enhances the accuracy and precision of detection.

While it's clear that supervised learning has its shortcomings when it comes to discovering new, unseen threats, supervised models only work as well as the data they're fed on because they need labeled examples. A supervised model can only be recognized if an attacker employed the same tactics (or something similar to it) in the past, and that was reflected in its training set; hence, if the attacker uses new language tactics or new phishing techniques, it will not be detected by a supervised model.

On the other hand, performing unsupervised learning is important for detecting emerging threats and anomalies. Unsupervised methods don't need labeled data and will uncover unexpected patterns or outlier behaviors. This flexibility makes them great at detecting insider threats, where bad guy communication might be slightly out of alignment with expected behavior, or for discovering new forms of social engineering altogether. Unsupervised learning

enables us to group data points into similar groups; without having a predefined array of categories, it becomes a powerful weapon for detecting anomalies.

However, while unsupervised models may flag more false positives, supervised models are more likely to miss precisely the bizarre account that could have brought down the system. That tendency means extra work for the security teams, who must rifle through the alerts to find the real threats. However, unsupervised learning is important for cybersecurity, allowing us to gain insight into new 'in the wild' emerging threats that might otherwise be missed.

## 8. Case Studies of NLP in Social Engineering Detection

Detection of social engineering attacks is one of the vital features that needs to be embedded into Cybersec systems to increase cybersecurity. These examples demonstrate how language analysis and machine learning enable organizations to identify enterprise-threatening activity early to stop costly data breaches and operation disruptions.

### 8.1. Real-World Examples Where NLP Helped Detect Social Engineering Attacks

A large bank used an NLP-powered system to monitor incoming emails for phishing threats in the financial services sector. Financial institutions are often targets of phishing attacks that use specific language meant to indicate urgency, such as "account suspension" or "transaction holding." Using keyword and sentiment analysis, the NLP model was scanned for suspicious phrases and flagged emails with high-urgency sentiment. So, trained on historical phishing data, this NLP system helped the bank to knock down phishing-related incidents by 70% within a year by filtering out malicious emails before they got to employees.

Yet another example of an NLP model was deployed by a multinational corporation to identify CEO fraud, a type of impersonation attack in which the attacker impersonates the senior executive to trick employees into transferring funds or sensitive information for fraud purposes. The NLP model was trained to spot anomalies in communication patterns, such as atypical language, inconsistent phrasing, and a speech that differed in tone from that of the typical company executive. In one case, the NLP system flagged a request from a supposed CEO to the finance department for an immediate wire transfer as suspicious. Upon review, the security team confirmed it was a fraudulent request. The NLP model protected the company from a costly financial scam by identifying tone and language inconsistencies.

In the telecommunications industry, NLP models have been used to monitor customer service channels, where attackers attempt to bypass security through social engineering support agents. An NLP model was trained to analyze chat transcripts for specific language patterns used in account takeover scams. Through keyword recognition and behavioral analysis, the system could identify when a caller's language didn't match the profile of the actual account holder. In one case, an NLP system flagged multiple attempts by an attacker to reset passwords on several accounts, alerting the security team, who then implemented additional verification steps. This early intervention prevented a potential data breach affecting hundreds of customers.

**Table 2** Case Studies of NLP in Action

| Organization Type | NLP Application | Outcome |
|---|---|---|
| Financial Institution | Phishing email monitoring | 70% reduction in phishing incidents |
| Multinational Corporation | CEO fraud detection | Prevented a significant financial scam |
| Telecom Company | Customer service chat monitoring | Detected multiple account takeover attempts |

### 8.2. Analysis of High-Profile Cases and Lessons Learned

These examples underscore some valuable lessons in applying NLP to social engineering detection. First, continuous model training is crucial. As social engineering tactics evolve, so must the NLP models, which require regular updates to remain effective against new threat patterns. Historical phishing phrases, for example, may become less common as attackers shift to more personalized language. Regularly training NLP models with recent threat data ensures they stay relevant and accurate.

Another lesson is the importance of contextual analysis. In the case of CEO fraud detection, the NLP model's ability to recognize inconsistencies in tone and timing proved invaluable. Many social engineering attacks rely on creating a sense

of urgency or bypassing typical communication channels. NLP systems that incorporate context can better discern when a message doesn't align with usual business practices.

Finally, these cases highlight the value of cross-functional integration. NLP models don't operate in isolation; they are most effective when integrated with broader security systems. For example, in the financial sector, flagged messages were sent directly to security analysts who could quickly ascertain and respond without delay. Faster, better-informed decision-making, critical to combating fast-moving social engineering attacks, is made possible by integrating NLP insights into security workflows.

These real-world cases illustrate that viable, though impactful, applications of NLP in cybersecurity exist to strengthen an organization's resilience against social engineering attacks. With regular model updates, context-driven analysis, and integration into security operations, NLP can serve as a robust defense against even the most sophisticated forms of human-targeted cyber threats.

## 9. Challenges in NLP-Based Cybersecurity Integration

Natural Language Processing (NLP) integration into cybersecurity protocol comes with some known and unknown challenges on the technical and ethical fronts. NLP presents a lot of potential for identifying and mitigating language-based threats. Still, data, model accuracy, and privacy challenges can often prohibit its successful deployment.

### 9.1. Technical Challenges: Data Scarcity, Language Nuances, and Model Accuracy

Data scarcity is one of the major technical challenges for deploying NLP for cybersecurity. High-quality, labeled data are required to build and train effective NLP models. In cybersecurity, however, data collection is greatly limited because it typically has to be collected under the bubble wrap of privacy rules and security-sensitive data. Moreover, there is no easy access to datasets with real phishing attempts, social engineering scams, or other malicious communications. Hence, it is challenging to build up a large enough training set. This scarcity limits the model's ability to recognize nuanced patterns, particularly in detecting newer, more sophisticated attacks.

Language nuances present another significant hurdle. Language in malicious communications can vary greatly, from slang to formal phrasing, different regional dialects, and even industry-specific jargon. Attackers often vary their language to avoid detection, modifying words, sentence structures, and tone. NLP models can struggle with these variations, especially if they haven't been exposed to diverse language samples during training. This limitation can lead to gaps in the model's ability to accurately assess intent, as subtle shifts in language may mask malicious intentions.

Ensuring model accuracy is also a persistent technical challenge. In cybersecurity, obtaining high accuracy in threat detection is essential, and training an NLP model to tell apart benign and malicious messages is hard reliably. Suppose you don't have very accurate models. If there are no threats there, then it will mark the messages that aren't suspicious as such, which is a false positive: if there are threats there, then it won't spot them — a false negative. False positives are burdensome but false negatives are a serious failure in security. Accuracy must be continuously tuned and retrained, especially as cyber threats change.

**Table 3** Challenges in NLP-Based Cybersecurity

| Challenge | Description | Potential Solutions |
|---|---|---|
| Data Scarcity | Limited access to labeled training data | Collaborate for data sharing |
| Language Nuances | Variability in language and phrasing | Diversify training datasets |
| Model Accuracy | Maintaining high accuracy in detection | Regular model recalibration |

### 9.2. Ethical Concerns: Privacy, False Positives, and Data Handling

NLP-based cybersecurity integration poses the same technical challenges and brings multiple ethical aspects, especially privacy. To analyze language, NLP models must process communications that may include sensitive or personal information. While necessary, if not used carefully, this access can breach user privacy and violate data protection compliance, like GDPR. Balancing the need to monitor everything with respect for privacy can always be addressed. At the same time, organizations must ensure that NLP models are used responsibly by processing only to the extent they are needed to secure things.

False positives also pose ethical challenges. If an NLP system incorrectly assigns legitimate communications to suspicious ones, it can impede normal business processes and may also be inconvenient for the employees. Suppose, for example, that too many routine communications are said to be possible phishing, so much so that employees begin to lose faith in the security system or even stop paying attention to legitimate alerts. The balancing act tries to minimize false positives but keep them secure, which needs to be constantly refined within appropriate criteria and thresholds.

Data handling practices are another ethical consideration. To build effective NLP models, vast amounts of communication data are often stored and processed. As such, this data must be secured in storage, anonymized as appropriate, and not retained beyond its useful timeframe to avoid misuse. Good data handling practices, ethically, not only protect user privacy but also offer trust to the security system. Data access, usage, and retention policies become clear and help align with regulatory standards and reduce ethical risk.

## 10. NLP provides cybersecurity with the much-needed benefits.

Despite this, NLP has the potential for massive benefits in cybersecurity, such as increased threat detection accuracy and cost and operational efficiencies that improve an organization's security posture.

### 10.1. How does NLP boost Threat Detection accuracy?

The most outstanding advantage of NLP in cybersecurity is that it helps to increase the threat detection rate. Finally, because they use NLP models, especially NLP models with machine learning, NLP models are capable of analyzing language patterns, tone, and context with above-average precision. By recognizing common cues in phishing, social engineering, and other language-based attacks, NLP can flag potentially dangerous communications before they escalate.

Through techniques like sentiment analysis, NLP can detect emotional tones commonly used in manipulative language, such as fear or urgency. Additionally, NLP models can perform contextual analysis and evaluate the relevance and plausibility of requests in communications. Take, for instance, an ee-mailpurported to be from an executive asking for an urgent transfer of funds that isn't written in the same tone or phrasing as usual for that executive, which NLP systems can notice and flag up for security teams as anomalous. This precision is a first line of defense. Being precise makes it less likely that an attacker can successfully exploit an attack because suspicious elements appear where otherwise they would not be.

Moreover, NLP's ability to adapt to new language patterns enhances its long-term efficacy. As threat actors change their tactics and wording, NLP models can be retrained on updated data, ensuring they remain effective against emerging language-based attacks. This adaptability enables continuous improvements in detection accuracy, providing organizations with a dynamic defense against evolving threats.
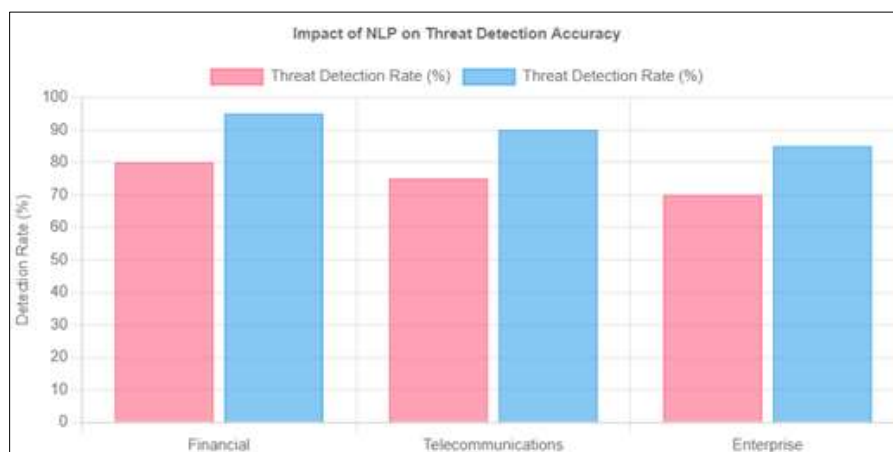


**Figure 2** Impact of NLP on Threat Detection Accuracy

### 10.2. Cost Savings and Operational Efficiency with NLP in Security

In addition to improving detection accuracy, NLP significantly contributes to cost savings and operational efficiency in cybersecurity. Threat detection in text-based channels, such as email, traditionally required extensive manual

monitoring by security teams. NLP does much of the work to automate this message analysis, removing most of them from the review path and marking the worst of them to be reviewed. Security teams can then prioritize high-value alerts and respond more quickly without needing labor-intensive manual assessments, potentially reducing labor costs.

Additionally, NLP decreases the costs resulting from data breaches, which can be large. A data breach can pour huge amounts out of your pockets in terms of making yourself pay for the ordeal, paying legal fees, paying fines, reputation damage, and just recovering from the whole ordeal. NLP, in particular, allows us to catch threats earlier, reducing the potential for these financial losses. This proactive approach helps save money and preserve the organization's reputation as no costly and public security incident is incurred.

NLP also enhances efficiency by reducing downtime as a consequence of cyber incidents. For example, when a ransomware attack paralyzes a company's operations, the company will experience revenue losses during recovery. NLP can also enhance early threat detection to give companies the ability to differentiate between the precursors of ransomware and other disruptive attacks that could cripple business operations and pose a much larger incident. This keeps business from stopping cold in case of extensive downtime, increasing productivity.

Overall, integrating NLP into cybersecurity protocols brings immediate and long-term benefits. Stronger threat detection accuracy gives language-based defenses against attack, reduced costs, and operational efficiencies for cybersecurity management. Using NLP allows organizations to secure their assets, minimize security expenses, and optimize their operations, thereby turning NLP into a valuable component in the modern cybersecurity model.

## 11. What Organizations Should Consider Before Using NLP in Security

Organizations adopting NLP for cybersecurity must navigate several practical and strategic considerations. Establishing clear steps for implementing NLP and training teams to manage and maintain these systems is essential to realizing the benefits of NLP-based security.

### 11.1. Steps for Adopting NLP-Based Security Protocols

The first thing you need to do to adopt NLP for Cybersecurity is a risk assessment to figure out the channels of communication and interactions most vulnerable to social engineering. It clarifies where NLP will have the greatest impact within an organization (within e-mailphishing detection, internal chat monitoring, or external customer communication).

Once priority areas are identified, the next step is to select the right NLP tools and models. Organizations must determine the available NLP platforms and choose models that fit within their security needs (e.g., real-time analysis and detection of specific threat types). This is also the time we make decisions about whether we want to build a custom NLP solution or whether we are looking for a cloud-based service.

They are integrating NLP into the broader security infrastructure, which is essential. NLP tools should be connected to SIEM systems and other security frameworks to ensure that all flagged communications are directed to the appropriate channels for review. This integration also enables security teams to streamline their workflows, allowing NLP alerts to be handled as part of their standard monitoring and response procedures.

### 11.2. Training Teams and Calibrating Systems for Optimal Performance

However, effective NLP implementation in security requires developing deep training for cybersecurity teams. Training security analysts to interpret NLP-generated alerts is another best practice. Still, they also need to understand the model's limitations (such as the chance of false positives or misinterpretation simply from language nuances). It should also train people to confirm flagged messages and ascertain and distinguish genuine threats from benign anomalies. The calibration of NLP models is equally important for maintaining optimal performance. Organizations should establish protocols for regular model recalibration and performance evaluation. By periodically retraining models with updated data, organizations ensure that their NLP tools continue to detect current threats accurately. Calibration sessions can also address any bias or drift in the model's criteria, which is especially relevant as attackers change their tactics or new social engineering tactics emerge.

Lastly, establishing feedback loops between cybersecurity and data science teams helps fine-tune NLP models. Analysts can provide insights based on real-world use, helping data scientists adjust the models to reduce false positives or improve the detection of nuanced threats. This ongoing collaboration maintains NLP systems in sync with the changing cybersecurity landscape and the organization's continually evolving cybersecurity needs.

## 12. Conclusion

Natural Language Processing (NLP) integration into cybersecurity protocols is the most effective step in protecting against modern cyber threats of social engineering and language manipulation type. Cyber security can be improved with NLP because systems can analyze and interpret human language so accurately that they can detect minute cues, intent, and behavioral abnormalities that suggest malicious activity. Using NLP, organizations benefit from a proactive defense layer that helps spot keyword detection leading to phishing attempts, analyzes sentiment and tone shifts signaling impersonation, and keeps sensitive information from being revealed, lowering the odds of a successful attack.

NLP has some impact on the cybersecurity scenario as it has helped enhance threat detection accuracy, reduce manual monitoring costs, and simplify security operations. NLP frees security teams to focus on high-priority alerts, automating the analysis of high volumes of communication data and communicating threat alerts, all in real time, making cybersecurity operations more efficient. Furthermore, NLP is adaptable to advancing language patterns, enhancing its pace compared to the velocity of attackers as the tactics become refined, keeping organizations ahead of an evolving threat landscape. In addition to operational gains, using NLP in cybersecurity can also assist companies in bolstering their reputation by decreasing the likelihood of a large-scale data breach and strengthening the trust held in the company from the perspective of their clients and other stakeholders.

As time goes by, NLP in cybersecurity will likely also undergo evolution, as with growing AI, one can expect it to become better at understanding context and detecting the user's intent. The advancement of models good at multi-layered reasoning, cross-platform monitoring, and unsupervised learning will allow us to perform more sophisticated analysis on the relatively subtle and diverse ways cybercriminals think and operate. However, as these systems become more powerful, ethics must continue to figure out a way to implement them. Privacy concerns must be responsibly balanced against the requirement for complete monitoring, while NLP models must be transparent, fair, and secure concerning their data practice. Careful management of the false positives and adherence to ethical data usage standards will encourage organizations to maintain trust as they deploy effective cybersecurity.

## Compliance with ethical standards

*Disclosure of conflict of interest*

No conflict of interest to be disclosed.

## References

[1] Sawa, Yuki & Bhakta, Ram & Harris, Ian & Hadnagy, Christopher. (2016). Detection of Social Engineering Attacks Through Natural Language Processing of Conversations. 262-265. 10.1109/ICSC.2016.95.

[2] N. DeCastro-García, A.L. ´ Munoz ˜ Castaneda, ˜ M. Fernandez-Rodríguez, ´ Machine learning for automatic assignment of the severity of cybersecurity events, Comput. Math. Methods Med. 2 (1) (2020) e1072.

[3] M. Husak, ´ T. Bajto˘s, J. Ka˘spar, E. Bou-Harb, P. Celeda, ˘ Predictive cyber situational awareness and personalized blacklisting: a sequential rule mining approach, ACM Trans. Manag. Inf. Syst. 11 (4) (2020) 1–6.

[4] F. Manganiello, M. Marchetti, M. Colajanni, Multistep attack detection and alert correlation in intrusion detection systems, in: International Conference on Information Security and Assurance, 2011, pp. 101–110.

[5] J.Q. Chen, Intelligent targeting with contextual binding, in: Future Technologies Conference (FTC), 2016, pp. 1040–1046.

[6] H. Studiawan, F. Sohel, Anomaly detection in a forensic timeline with deep autoencoders, J. Inf. Secur. Appl. 63 (2021), 103002.

[7] F. Amato, A. Castiglione, G. Cozzolino, F. Narducci, A semantic-based methodology for digital forensics analysis, J. Parallel Distrib. Comput. 138 (2020) 172–177.

[8] J. Sakhnini, H. Karimipour, A. Dehghantanha, R.M. Parizi, Physical layer attack identification and localization in cyber–physical grid: an ensemble deep learning based approach, Phys. Commun. 47 (2021), 101394.

[9] L. Fernandez Maimo, A. Huertas Celdran, A.L. Perales Gomez, F.J. Garcia Clemente, J. Weimer, I. Lee, Intelligent and dynamic ransomware spread detection and mitigation in integrated clinical environments, Sensors 19 (5) (2019) 1114.

[10] P. Nespoli, F.G. Marmol, ´ J.M. Vidal, A bio-inspired reaction against cyberattacks: ais-powered optimal countermeasures selection, IEEE Access 9 (2021) 60971–60996.

[11] M. Husak, ´ L. Sadlek, S. Spa ˇ ˇcek, M. Laˇstoviˇcka, M. Javorník, J. Komarkov ´ a, ´ CRUSOE: a toolset for cyber situational awareness and decision support in incident handling, Comput. Secur. 115 (2022), 102609.

[12] 7,396 NLP images, stock photos, and vectors | Shutterstock. (n.d.). Retrieved from https://www.shutterstock.com/search/nlp?page=2

[13] M. Hus´ak, Towards a data-driven recommender system for handling ransomware and similar incidents, in: IEEE International Conference on Intelligence and Security Informatics (ISI), 2021, pp. 1–6.

[14] B. Woods, S.J. Perl, B. Lindauer, Data mining for efficient collaborative information discovery, in: Proceedings of the 2nd ACM Workshop on Information Sharing and Collaborative Security, 2015, pp. 3–12.

[15] S. Peng, A. Zhou, S. Liao, L. Liu, A threat actions extraction method based on the conditional co-occurrence degree, in: 7th International Conference on Information Science and Control Engineering (ICISCE), 2020, pp. 1633–1637.

[16] BlOG | Social Engineering 101: How to Safeguard your Business. (n.d.). Retrieved from https://www.netfriends.com/blog-posts/social-engineering-101-how-to-safeguard-your-business

[17] B.S. Meyers, A. Meneely, An automated post-mortem analysis of vulnerability relationships using natural language word embeddings, Procedia. Comput. Sci. (2021) 953–958.

[18] M.V. Carriegos, A.L. ´ Castaneda, ˜ M.T. Trobajo, D.A. De Zaballa, On aggregation and prediction of cybersecurity incident reports, IEEE Access 9 (2021) 102636–102648.

[19] Dias, F. S., & Peters, G. W. (2020). A non-parametric test and predictive model for signed path dependence. Computational Economics, 56(2), 461-498.