**Paper Title:**
Cyber Security Vulnerability Detection Using Natural Language Processing

**Paper Link:**
https://ieeexplore.ieee.org/document/9817336

<u>1 Summary</u>

**1.1 Motivation**
With the exponential growth in technological advancements the intensity and sophistication of cyber threats is rising. Concern is evident in the fact that detecting a security breach takes companies an average of 197 days, with an additional 69 days to contain it. This prolonged response time leads to substantial financial and operational losses, unscheduled downtime, and reduced productivity. The paper demonstrates how Natural Language Processing (NLP) can be used to detect vulnerabilities in software, hypothesizing that NLP techniques can effectively identify and predict security vulnerabilities in code.

**1.2 Contribution**
Firstly, NLP makes it easier to find vulnerabilities in software code. Traditional methods often miss tricky security flaws because they rely on simple rules. NLP, however, understands the meaning and context of code, so it can find vulnerabilities that others might overlook. Secondly, NLP speeds up how quickly we respond to cyber threats. The big problem is companies take a lot of time to detect a security breach and it gives hackers the time to cause damage. NLP can help spot these issues much faster by analyzing lots of code quickly. Moreover, NLP makes it more accurate to find security issues in code. Older methods often make mistakes or give false alarms.

**1.3 Methodology**
In this research, the authors have undertaken a comprehensive exploration of natural language processing by developing and training a diverse set of deep learning models, including LSTM, BiLSTM, BERT, BERT for sequence classification, and CodeBERT. Their primary objective was to conduct a rigorous comparative analysis of these models to evaluate their effectiveness in the domain of cyber security vulnerability detection. To achieve this goal, they have diligently trained these models on an extensive dataset encompassing a wide spectrum of source code written in C/C++.The methodology involves training advanced NLP models on a dataset comprising over 70,000 C/C++ files, including both vulnerable and non-vulnerable code samples. To sum up, this methodology blends the strengths of LSTM, BiLSTM, BERT, and CodeBERT to enhance cyber security vulnerability detection in source code.

**1.4 Conclusion**
The study concludes that NLP models, particularly CodeBERT, are highly effective in detecting vulnerabilities in software code. Additionally, it emphasizes the importance of maintaining contextual information throughout the entire sequence for this specific task. The paper also offers the created collection of software vulnerabilities as a valuable resource for other researchers to further enhance their work in this area along with their accuracy and the best performer achieved 95% accurate result compared to various deep learning models.

2 Limitations

## 2.1 First Limitation

A primary limitation is the focus on C/C++ programming languages. This narrow scope may not adequately represent the diversity of languages used in software development, potentially limiting the applicability of the findings across different programming contexts.

## 2.2 Second Limitation

Another limitation is the need for continual model updates. The rapidly evolving nature of cyber threats means that the models, as they stand, might become outdated quickly, necessitating constant updates and refinements to maintain efficacy.

3 Synthesis

The ideas presented in this paper open up numerous potential applications and future scopes. The successful use of NLP for vulnerability detection in software could be expanded to other programming languages, enhancing security measures across various platforms. Furthermore, these techniques have the potential to automate and streamline the process of code review and security auditing, significantly reducing the time and resources required for these tasks. The principles could also be adapted for real-time monitoring systems, offering proactive defense mechanisms against emerging threats. This research paves the way for a new era in cybersecurity, where AI and language processing tools play a crucial role in safeguarding digital assets.