

ĐẠI HỌC BÁCH KHOA HÀ NỘI

ĐỒ ÁN TỐT NGHIỆP

Phát hiện tấn công web chưa biết sử dụng mạng thích
nghi sâu

NGUYỄN QUỐC HƯNG

hung.nq200294@sis.hust.edu.vn

Ngành: Kỹ thuật máy tính

Giảng viên hướng dẫn: TS. Tống Văn Vạn

Chữ kí GVHD

Ngành:

Kỹ thuật máy tính

Trường:

Công nghệ Thông tin và Truyền thông

HÀ NỘI, 06/2024

LỜI CAM KẾT

Họ và tên sinh viên: Nguyễn Quốc Hưng

Điện thoại liên lạc: 0982947213

Email: hung.nq200294@sis.hust.edu.vn

Lớp: Kỹ thuật máy tính 01

Hệ đào tạo: Cử nhân

Tôi – *Nguyễn Quốc Hưng* – cam kết Đồ án Tốt nghiệp (ĐATN) là công trình nghiên cứu của bản thân tôi dưới sự hướng dẫn của *TS. Tổng Văn Vạn*. Các kết quả nêu trong ĐATN là trung thực, là thành quả của riêng tôi, không sao chép theo bất kỳ công trình nào khác. Tất cả những tham khảo trong ĐATN – bao gồm hình ảnh, bảng biểu, số liệu, và các câu từ trích dẫn – đều được ghi rõ ràng và đầy đủ nguồn gốc trong danh mục tài liệu tham khảo. Tôi xin hoàn toàn chịu trách nhiệm với dù chỉ một sao chép vi phạm quy chế của nhà trường.

Hà Nội, ngày 30 tháng 6 năm 2024

Tác giả ĐATN

Nguyễn Quốc Hưng

LỜI CẢM ƠN

Tôi xin gửi lời cảm ơn chân thành đến T.S Tống Văn Vạn, người đã tận tâm hỗ trợ và hướng dẫn tôi trong suốt quá trình nghiên cứu và hoàn thành đề án này. Thầy đã luôn sẵn sàng chia sẻ kiến thức sâu rộng và những lời khuyên quý báu, giúp tôi vượt qua nhiều thử thách khó khăn và đạt được những kết quả đáng khích lệ. Không có sự giúp đỡ và sự kiên nhẫn của thầy, tôi sẽ không thể hoàn thành được đề án này.

Tôi cũng muốn bày tỏ lòng biết ơn sâu sắc đến bố mẹ và bạn gái tôi. Họ đã luôn đứng bên cạnh, ủng hộ và động viên tôi trong suốt hành trình đầy thách thức này. Sự yêu thương và khích lệ không ngừng nghỉ của họ là nguồn động lực to lớn, giúp tôi vượt qua những thời điểm khó khăn nhất. Sự kiên nhẫn và sự tin tưởng của họ đã là động lực để tôi không ngừng phấn đấu.

Ngoài ra, tôi xin gửi lời cảm ơn chân thành đến các bạn bè thân thiết, những người đã luôn đồng hành và chia sẻ cùng tôi trên con đường học tập và nghiên cứu. Các bạn đã hỗ trợ tôi không chỉ về mặt tinh thần mà còn trong công việc nghiên cứu, giúp tôi có được những góc nhìn và ý tưởng mới mẻ. Sự sẻ chia và đồng cảm của các bạn đã giúp tôi vượt qua nhiều trở ngại và mang lại niềm vui, sự hứng khởi trong suốt quá trình học tập và làm việc. Sự động viên và giúp đỡ của các bạn là một phần không thể thiếu trong thành công của tôi.

TÓM TẮT NỘI DUNG ĐỒ ÁN

Với sự phát triển nhanh chóng của công nghệ thông tin, các ứng dụng web đã trở thành một phần không thể thiếu trong cuộc sống hàng ngày, đặc biệt khi ngày càng nhiều người chuyển ứng dụng và dữ liệu nhạy cảm lên đám mây. Sự phổ biến của các ứng dụng web và lượng lớn dữ liệu nhạy cảm mà chúng xử lý làm cho chúng dễ trở thành mục tiêu của các cuộc tấn công. Do đó, việc bảo vệ các ứng dụng web khỏi các mối đe dọa này là vô cùng quan trọng. Các cuộc tấn công phổ biến như SQL Injection, XSS, và CSRF đã được OWASP xác định là những mối đe dọa nghiêm trọng. Sự phát triển không ngừng của các ứng dụng web dẫn đến sự xuất hiện của các biến thể tấn công mới, đòi hỏi các hệ thống phát hiện không chỉ nhận dạng các mô hình tấn công đã biết mà còn phải thích nghi với những mô hình mới chưa biết. Đề tài này tập trung vào phát triển một hệ thống phân loại lỗ hổng bảo mật mới với độ chính xác cao và thời gian xử lý hợp lý, có khả năng xử lý đồng thời nhiều yêu cầu. Chúng tôi đề xuất sử dụng mạng thích ứng miền sâu (Deep Adaptation Network) để giải quyết các hạn chế của các phương pháp hiện tại, bao gồm cả những phương pháp dựa trên chữ ký như ModSecurity và các kỹ thuật học máy, vốn gặp khó khăn khi phải đối mặt với các lỗ hổng mới (zero-days). Hệ thống của chúng tôi áp dụng các kỹ thuật xử lý ngôn ngữ tự nhiên (NLP), mô hình BERT, học sâu, và chiến lược học chuyển giao để phát hiện và phân loại các cuộc tấn công web chưa biết. Cụ thể, chúng tôi sử dụng SecBERT để trích xuất đặc trưng từ các yêu cầu HTTP và áp dụng tỷ lệ chuyển giao động cho mạng thích nghi sâu nhằm cải thiện độ chính xác và tốc độ hội tụ của mô hình. Kết quả thí nghiệm cho thấy phương pháp đề xuất đạt F1-score cao, đặc biệt là 99% trong phân loại các lỗ hổng bảo mật đã biết và 71% đối với các lỗ hổng mới với tỷ lệ chuyển giao động. Hệ thống khi tích hợp với WAF có khả năng phân loại với độ chính xác cao và thời gian phản hồi trung bình khoảng 18ms, đánh dấu một bước tiến quan trọng trong việc tăng cường an ninh cho các ứng dụng web.

Sinh viên thực hiện
(Ký và ghi rõ họ tên)

ABSTRACT

With the rapid advancement of information technology, web applications have become an indispensable part of daily life, especially as more people move their applications and sensitive data to the cloud. The widespread use of web applications and the vast amount of sensitive data they process make them prime targets for attacks. Therefore, protecting web applications from these threats is of paramount importance. Common attacks such as SQL Injection, XSS, and CSRF have been identified by OWASP as serious threats. The continuous development of web applications has led to the emergence of new attack variants, requiring detection systems not only to recognize known attack patterns but also to adapt to new, unknown ones. This project focuses on developing a new security vulnerability classification system with high accuracy and reasonable processing time, capable of handling multiple requests simultaneously. We propose using a Deep Domain Adaptation Network to address the limitations of current methods, including signature-based approaches like ModSecurity and machine learning techniques, which struggle with new vulnerabilities (zero-days). Our system employs Natural Language Processing (NLP) techniques, the BERT model, deep learning, and transfer learning strategies to detect and classify unknown web attacks. Specifically, we use SecBERT to extract features from HTTP requests and apply a dynamic transfer rate for the deep domain adaptation network to improve the model's accuracy and convergence speed. Experimental results show that the proposed method achieves high F1-scores, notably 99% in classifying known security vulnerabilities and 71% for new vulnerabilities using a dynamic transfer rate. When integrated with WAF, the system can classify with high accuracy and an average response time of just 18ms, marking a significant advancement in enhancing the security of web applications.

MỤC LỤC

CHƯƠNG 1. GIỚI THIỆU ĐỀ TÀI.....	1
1.1 Đặt vấn đề.....	1
1.2 Các giải pháp hiện tại và hạn chế	1
1.3 Mục tiêu và định hướng giải pháp	2
1.4 Đóng góp của đề án	2
1.5 Bố cục đề án	3
CHƯƠNG 2. NỀN TẢNG LÝ THUYẾT	5
2.1 Ngữ cảnh của bài toán.....	5
2.2 Background.....	5
2.2.1 Xử lý ngôn ngữ tự nhiên (NLP)	5
2.2.2 Bidirectional Encoder Representations from Transformers (BERT)	
7	
2.2.3 Học chuyển giao	10
2.2.4 Tổng quan về tấn công web trong an ninh mạng	11
2.3 Nghiên cứu liên quan	15
CHƯƠNG 3. PHƯƠNG PHÁP ĐỀ XUẤT.....	18
3.1 Tổng quan giải pháp.....	18
3.2 Hệ thống phân loại lỗ hổng bảo mật sử dụng DAN.....	19
3.2.1 Phase 1: Tiền xử lý dữ liệu	19
3.2.2 Phase 2: Feature Extraction với SecBERT	21
3.2.3 Phase 3: Deep Adaptation Network (DAN).....	22
3.3 Sử dụng dynamic transfer rate.....	25
3.4 Triển khai tường lửa tích hợp với mô hình DAN.....	27

CHƯƠNG 4. ĐÁNH GIÁ THỰC NGHIỆM.....	29
4.1 Các tham số đánh giá	29
4.2 Đặc tả dữ liệu	30
4.3 Phân tích hiệu năng.....	33
4.3.1 Thiết lập thí nghiệm	33
4.3.2 So sánh với các phương pháp khác.....	34
4.3.3 Cải thiện với dynamic transfer rate.....	38
4.3.4 Thời gian dự đoán tấn công web trên hệ thống	40
CHƯƠNG 5. KẾT LUẬN	41
5.1 Kết luận	41
5.2 Hướng phát triển trong tương lai	41
TÀI LIỆU THAM KHẢO.....	44

DANH MỤC HÌNH VẼ

Hình 2.1	Kiến trúc mô hình Tranformer	7
Hình 2.2	Sơ đồ kiến trúc BERT cho tác vụ Masked ML	9
Hình 2.3	Các nhánh khác nhau của học chuyển giao [9]	10
Hình 2.4	Kiến trúc của một ứng dụng Web cơ bản	11
Hình 2.5	Các biện pháp chống lại tấn công ứng dụng web	14
Hình 3.1	Tổng quan mô hình	19
Hình 3.2	Ví dụ về sử dụng mã hóa số học chuỗi	21
Hình 3.3	Đặc tả phase 3	23
Hình 3.4	Ví dụ về phân phối đặc trưng của dữ liệu sau khi thực hiện thích ứng miền	24
Hình 3.5	Giá trị của tham số λ trong các trường hợp γ khác nhau và λ_1 $= 1$	26
Hình 3.6	Mô hình triển khai tường lửa tích hợp với mô hình DAN	28
Hình 4.1	Phân phối giá trị trung bình của <i>request_http_method</i> , <i>request_-</i> <i>http_request</i> và <i>request_http_protocol</i>	32
Hình 4.2	Kết quả độ chính xác qua các epochs với từng giá trị γ khác nhau ($\lambda_1 = 0.5$) so sánh với λ cố định	39

DANH MỤC BẢNG BIỂU

Bảng 3.1	FEATURES SELECTION	20
Bảng 4.1	Số lượng yêu cầu HTTP được đánh nhãn theo CAPEC [4]. . .	31
Bảng 4.2	Giá trị trung bình	31
Bảng 4.3	Ví dụ về các tính chất của yêu cầu HTTP tấn công và bình thường.	32
Bảng 4.4	Known Attack Dataset (CAPEC source)	33
Bảng 4.5	Unknown Attack Dataset (CAPEC target)	33
Bảng 4.6	So sánh tác vụ phân loại lỗ hổng bảo mật đã biết với các phương pháp học máy, học sâu thường được sử dụng.	35
Bảng 4.7	F1-score trên tác vụ học chuyển giao so với các phương pháp truyền thống.	36
Bảng 4.8	F1-score trên tác vụ học chuyển giao so với các phương pháp thích ứng miền khác nhau	37
Bảng 4.9	Hiệu năng trên tác vụ học chuyển giao với các giá trị γ và λ_1 khác nhau	38
Bảng 4.10	Kết quả trên tác vụ phân loại với các giá trị γ khác nhau . . .	38
Bảng 4.11	So sánh thời gian suy luận với các phương pháp khác	40

DANH MỤC THUẬT NGỮ VÀ TỪ VIẾT TẮT

Thuật ngữ	Ý nghĩa
BOW	Bag-of-words
CMD	Central Moment Discrepancy
CNN	Mạng thần kinh tích chập(Convolutional Neural Network)
Coral	Correlation Alignment
CRS	Core Rule Set
DA	Domain Adaptation
DAN	Mạng thích nghi sâu (Domain Adaptation Network)
DDC	Deep Domain Confusion
DL	Deep Learning
HTML	Ngôn ngữ đánh dấu siêu văn bản (HyperText Markup Language)
HTTP	Hypertext Transfer Protocol
LSTM	Long Short-Term Memory
MK-MMD	Multi-kernel Maximum Mean Discrepancy
ML	Machine Learning
MLP	Multi-layer Perceptron
MMD	Maximum Mean Discrepancy
NLP	Natural Language Processing
OWASP	Open Web Application Security Project
SQLi	SQL injection
SWD	Sliced Wasserstein Distance
TF-IDF	Term Frequency Inverse Document Frequency
UDA	Unsupervised Domain Adaptation
WAFs	Web Application Firewalls
XSS	Cross-Site Scripting