

ĐẠI HỌC BÁCH KHOA HÀ NỘI

ĐỒ ÁN TỐT NGHIỆP

Phát hiện lưu lượng phần mềm truy cập từ xa trong
mạng nội bộ ứng dụng học sâu

TRẦN ĐÌNH KIẾN GIANG
giang.tdk194265@sis.hust.edu.vn

Ngành: Công nghệ thông tin

Giảng viên hướng dẫn: PGS.TS. Trần Quang Đức

Chữ kí GVHD

Khoa: Kỹ thuật máy tính

Trường: Công nghệ Thông tin và Truyền thông

HÀ NỘI, 01/2024

LỜI CẢM ƠN

Đầu tiên, em xin gửi lời cảm ơn tới PSG.TS. Trần Quang Đức, Giám đốc Trung tâm An toàn An ninh Thông tin, Đại học Bách Khoa Hà Nội, là giảng viên hướng dẫn cho em trong quá trình làm đồ án tốt nghiệp. Thầy là người đưa ra ý tưởng đề tài, nhiệt tình hướng dẫn và truyền đạt kiến thức cho em giúp em có thể hoàn thành đồ án này. Em xin cảm ơn TS. Tống Văn Vạn đã hướng dẫn em trong quá trình thu thập dữ liệu và xây dựng mô hình. Em cũng xin gửi lời cảm ơn tới các thầy cô ở Đại học Bách Khoa Hà Nội, đặc biệt là các thầy cô ở trường Công nghệ Thông tin và Truyền thông đã tận tâm giảng dạy, truyền đạt những kiến thức quý báu cho em trong suốt quá trình em học tập tại trường. Đồng thời, em xin cảm ơn các thầy và các anh ở Trung tâm An toàn An ninh thông tin Bách Khoa đã tạo điều kiện về máy móc giúp em thực hiện đồ án một cách thuận lợi hơn. Cuối cùng, em xin gửi lời cảm ơn tới gia đình, bạn bè, đồng nghiệp đã hỗ trợ và đồng hành cùng em trong suốt những năm tháng học tập và rèn luyện tại trường.

Em xin chân thành cảm ơn !

TÓM TẮT NỘI DUNG ĐỒ ÁN

Đồ án đề xuất một hệ thống phát hiện và quản lý các luồng truy cập từ xa sử dụng phần mềm trong mạng nội bộ. Đồ án tập trung vào hai nhiệm vụ chính, đó là xây dựng và huấn luyện một số mô hình học sâu nhằm phân loại các luồng mạng và xây dựng một hệ thống trong mạng nội bộ ứng dụng các mô hình đó một cách trực tiếp. Phương pháp phân loại lưu lượng mạng dựa trên học sâu và tải trọng đã được áp dụng, với ba loại mô hình khác nhau gồm 2DCNN, CNN kết hợp LSTM và Resnet50 đã được xây dựng một cách thủ công và được tinh chỉnh đầu vào cho phù hợp với bài toán phân loại lưu lượng mạng, với mỗi mô hình gồm bốn kích thước tải trọng đầu vào khác nhau. Bốn nhãn lớp được các mô hình phân loại bao gồm Non-RAT (không phải là phần mềm truy cập từ xa), RDP (Remote Desktop Connection), VNC (Virtual Network Computing) và TeamViewer. Hệ thống được đề xuất bao gồm các thiết bị trong mạng nội bộ, gọi là các *client*, và một thiết bị trung tâm để xử lý, gọi là *server*. Client được cài đặt Suricata và Filebeat, có nhiệm vụ ghi nhận và chuyển tiếp tất cả lưu lượng mạng xuất hiện trên thiết bị lên server. Server được cài đặt Logstash, một chương trình phân loại sử dụng ngôn ngữ Python và được import mô hình học sâu, Filebeat, và Kibana, có nhiệm vụ phân loại các luồng mạng trên client và hiển thị kết quả một cách trực quan. Tất cả các công cụ trong hệ thống đều có thể được chạy dưới dạng một service và khởi động cùng với hệ điều hành, nên không ảnh hưởng đến trải nghiệm của người dùng. Kết quả thực nghiệm cho thấy các mô hình đều có kết quả phân loại rất cao trên tập dữ liệu kiểm thử, nhưng lại kém đi khi triển khai trên thực tế với lưu lượng mạng bình thường (không phải lưu lượng phần mềm truy cập từ xa). Nhược điểm này được khắc phục khi kết hợp sử dụng mô hình học sâu với phương pháp sử dụng signature bytes (magic bytes) của các phần mềm truy cập từ xa. Ngoài ra, thời gian phản hồi trực tiếp của hệ thống trước lưu lượng mạng mới cũng là một vấn đề quan trọng được bàn luận trong đồ án. Từ đó, đồ án đưa ra tương quan so sánh mức độ hiệu quả giữa các mô hình với kích thước tải trọng đầu vào khác nhau. Kết quả cho thấy mỗi mô hình có một ưu, nhược điểm riêng và đều cho kết quả tốt khi triển khai thực tế. Bên cạnh đó, đồ án cũng chỉ ra và phân tích những hạn chế hiện có của hệ thống và đề xuất một số giải pháp trong tương lai để giúp hệ thống được hoàn thiện hơn.

Sinh viên thực hiện

(Ký và ghi rõ họ tên)

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	3
1.5 Bố cục đề án	3
CHƯƠNG 2. NỀN TẢNG LÝ THUYẾT	5
2.1 Ngữ cảnh của bài toán.....	5
2.2 Tổng quan hệ thống phát hiện xâm nhập	6
2.3 Một số công nghệ sử dụng	7
2.3.1 Suricata.....	7
2.3.2 Filebeat	8
2.3.3 Logstash.....	9
2.3.4 Kibana	9
2.4 Giao thức AES 256-bit	10
2.5 Một số mô hình học sâu	11
2.5.1 Mạng Convolutional Neuron Network (CNN).....	11
2.5.2 Mạng Long-Short Term Memory (LSTM).....	14
2.5.3 Resnet 50	16
CHƯƠNG 3. XÂY DỰNG VÀ HUẤN LUYỆN MÔ HÌNH.....	18
3.1 Các mô hình thí nghiệm	18
3.1.1 Two-Dimensional Convolutional Neural Network (2DCNN)	18
3.1.2 Convolutional Neural Network with Long Short-Term Memory (CNN_LSTM)	18

3.1.3 Residual Network with 50 layers (Resnet50).....	19
3.2 Xây dựng tập dữ liệu.....	21
3.3 Huấn luyện mô hình.....	24
CHƯƠNG 4. ĐỀ XUẤT HỆ THỐNG.....	26
4.1 Tổng quan hệ thống.....	26
4.2 Cấu hình trên client.....	27
4.2.1 Cấu hình Suricata.....	27
4.2.2 Cấu hình Filebeat.....	31
4.3 Cấu hình trên server.....	32
4.3.1 Cấu hình Logstash.....	32
4.3.2 Xây dựng chương trình phân loại.....	33
4.3.3 Cài đặt Filebeat.....	35
4.3.4 Cấu hình Kibana.....	36
CHƯƠNG 5. ĐÁNH GIÁ THỰC NGHIỆM.....	38
5.1 Đánh giá hiệu suất mô hình trong huấn luyện.....	38
5.1.1 Các chỉ số đánh giá.....	38
5.1.2 Kết quả đánh giá.....	39
5.2 Đánh giá hiệu quả triển khai thực tế.....	40
5.2.1 Độ chính xác của mô hình trong thực tế.....	42
5.2.2 Thời gian phản hồi của hệ thống.....	43
5.3 Thảo luận.....	47
5.3.1 Về việc dự đoán luồng mạng không phải là RAT trong thực tế.....	47
5.3.2 Về số lượng ứng dụng RAT có thể kiểm tra.....	48
5.3.3 Về thời gian phản hồi của hệ thống.....	48
CHƯƠNG 6. KẾT LUẬN.....	50
6.1 Kết luận.....	50

6.2 Hướng phát triển trong tương lai	50
TÀI LIỆU THAM KHẢO.....	52

DANH MỤC HÌNH VẼ

Hình 2.1	Kiến trúc của NIDS và HIDS	7
Hình 2.2	Kiến trúc mạng CNN	12
Hình 2.3	Lớp tích chập	12
Hình 2.4	Lớp tổng hợp	13
Hình 2.5	Lớp kết nối đầy đủ	13
Hình 2.6	Kiến trúc cơ bản của một RNN	14
Hình 2.7	Cấu trúc của một đơn vị RNN	15
Hình 2.8	Cấu trúc của một đơn vị LSTM	15
Hình 2.9	Khối dư	16
Hình 2.10	Kiến trúc mạng Resnet50	17
Hình 4.1	Sơ đồ kiến trúc tổng quan hệ thống.	26
Hình 4.2	Biểu đồ hoạt động của client khi nhận được gói tin mới.	28
Hình 4.3	Một gói tin có payload gồm chủ yếu các byte 0	30
Hình 4.4	Biểu đồ hoạt động của chương trình	34
Hình 5.1	Loss của mô hình 2DCNN với đầu vào khác nhau	40
Hình 5.2	Loss của mô hình CNN_LSTM với đầu vào khác nhau	40
Hình 5.3	Loss của mô hình Resnet50 với đầu vào khác nhau	41
Hình 5.4	Confusion matrix của CNN_LSTM với đầu vào khác nhau	41
Hình 5.5	Recall lớp Non-RAT của các mô hình trong thực tế	43
Hình 5.6	Thời gian phản hồi của mô hình 2DCNN với đầu vào khác nhau	45
Hình 5.7	Thời gian phản hồi của mô hình CNN_LSTM với đầu vào khác nhau	46
Hình 5.8	Thời gian phản hồi của mô hình Resnet50 với đầu vào khác nhau	46

DANH MỤC BẢNG BIỂU

Bảng 3.1	Tổng quan mô hình 2DCNN với kích thước payload đầu vào 128 bytes.	19
Bảng 3.2	Tổng quan mô hình CNN_LSTM với kích thước payload đầu vào 128 bytes.	19
Bảng 3.3	Tổng quan mô hình Resnet50 với kích thước payload đầu vào 128 bytes	20
Bảng 3.4	Số luồng mạng thu thập được của các nhãn.	23
Bảng 5.1	Accuracy của các mô hình với đầu vào khác nhau.	39
Bảng 5.2	Precision của CNN_LSTM với đầu vào khác nhau.	42
Bảng 5.3	Recall của CNN_LSTM với đầu vào khác nhau.	42
Bảng 5.4	F1-Score của CNN_LSTM với đầu vào khác nhau.	42
Bảng 5.5	Signature bytes của các ứng dụng.	47

DANH MỤC THUẬT NGỮ VÀ TỪ VIẾT TẮT

Thuật ngữ	Ý nghĩa
CNN	Mạng tích chập (Convolutional neural network)
LSTM	Bộ nhớ dài hạn - ngắn hạn (Long short term memory)
RAT	Phần mềm truy cập từ xa (Remote Access Tool)