



WhiteHat
School

MITRE ATT&CK TTPs 행위 분석 데이터(capa-rule)를 활용한 악성코드 분석 및 머신러닝 기반 악성코드 탐지

WhiteHat School 2nd
Project Kick-Off

알파고를 이을 악파고

1

팀 및 프로젝트 소개

- 팀원 소개.....03
- 주제 선정 배경.....04
- 프로젝트 목표.....10

2

프로젝트 수행 계획

- 데이터셋 수집 방안.....11
- AI 모델 학습 방안.....15
- AI 모델 검증 방안.....16

3

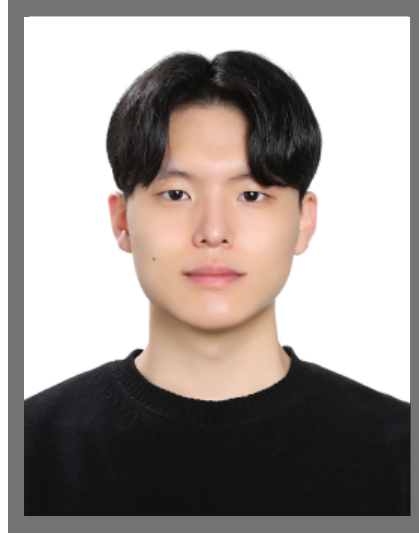
프로젝트 결과

- 예상 산출물.....17
- 대외 활동.....18
- 일정.....19

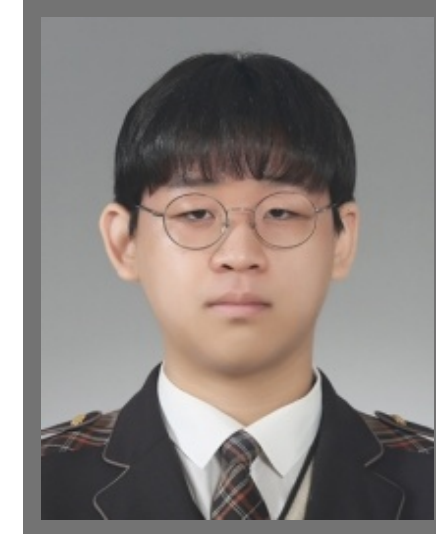
팀원 소개



멘토
손승호



PL
김두영



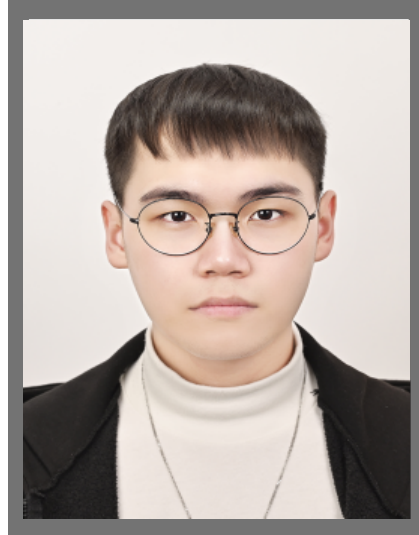
PM
오태호



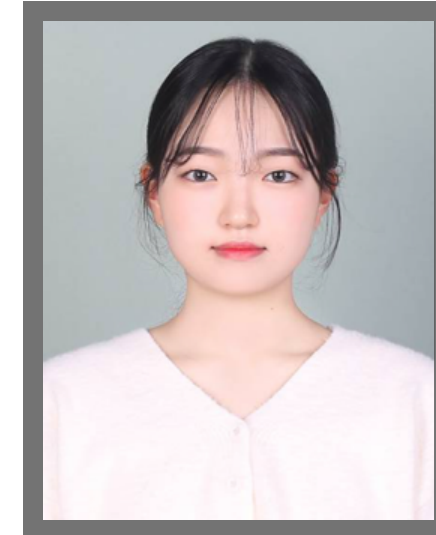
팀원
김나연



팀원
김상훈



팀원
이시연



팀원
임나현



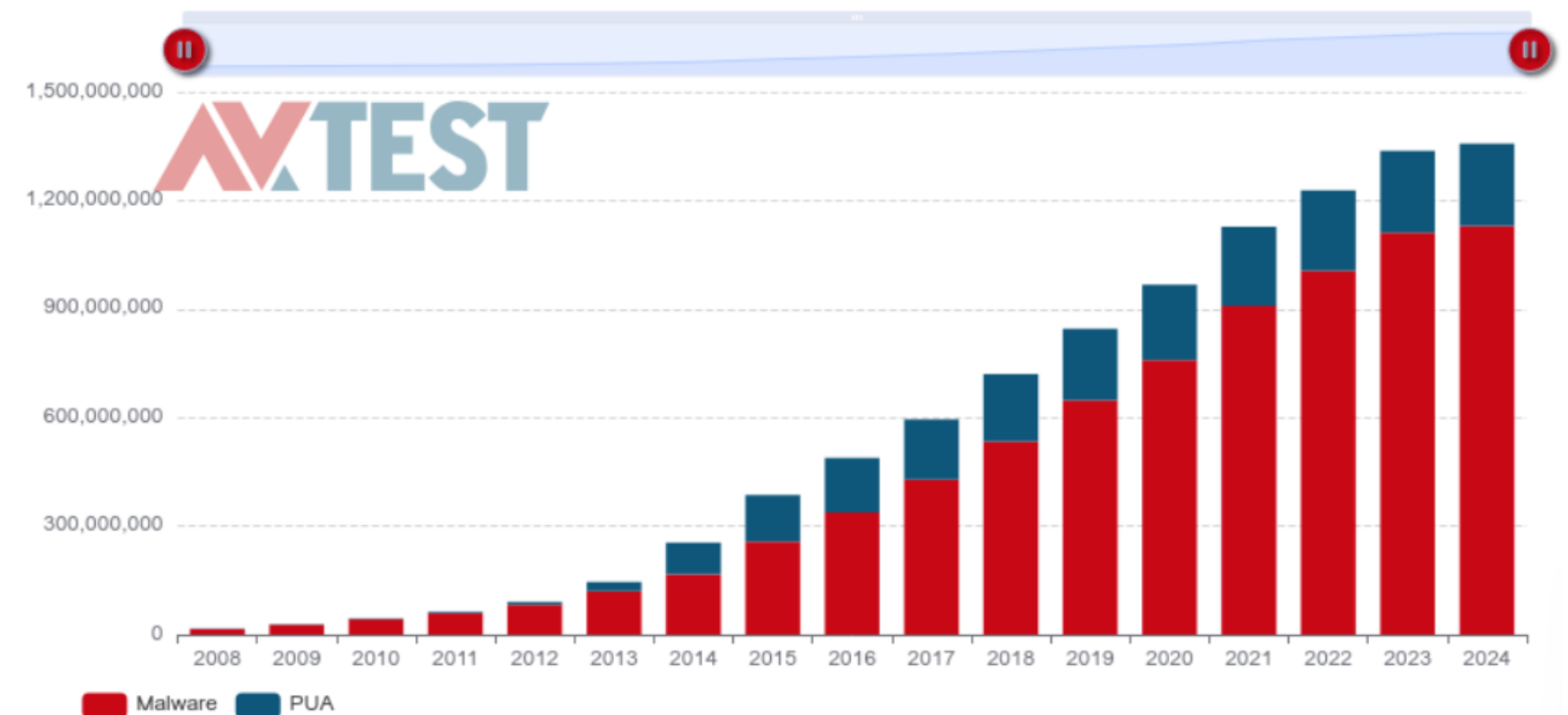
팀원
허라영

주제 선정 배경

악성코드 피해 사례

- 2024, 北, 국내 방산기업 10여 곳 해킹
- 2024, NFT 기반 가상자산 엔터버튼 해킹 및 토큰 도난
- 2024, Notepad++ 기본 플러그인 악용한 위키로더 악성코드 유포

TOTAL AMOUNT OF MALWARE AND PUA



주제 선정 배경

1. 탐지방식의 한계

- 보안전문가가 패턴을 분석할 때까지 무방비구간이 존재하며 대응까지 평균 205일이 소요된다.
- 악성코드의 기본적인 탐지 방식은 시그니처 기반이기에 새로운 악성코드를 탐지하기 어렵다.
- 일부 악성코드의 패턴과 정상 프로그램의 패턴이 매우 유사해 정상 프로그램을 악성코드로 오진하는 사례가 발생하고 있다.

2. 백신 등 보안 소프트웨어 우회

- 공격자가 악성코드를 Virustotal, Cuckoo Sandbox등을 이용하여 테스트하거나 이를 우회하는 방법을 적용할 수 있다.

주제 선정 배경

- 갈 수록 고도화 되가는 악성코드에 신속정확히 대응해야한다.
- 기존 악성코드의 특성을 머신러닝 기술로 학습하여 유사한 형태의 신종/변종 악성코드를 판별해낼 수 있다.

고도화 악성코드 사전 유입...대응 쉽지 않아

2013-03-21 21:29

속수무책으로 늘
악성코드는 높은
에 대비해 어떤

[보안칼럼]난독화된 악성앱, AI 기술로 분석 효율 높여야

발행일 : 2023-09-05 16:00 지면 : 2023-09-06 26면

- 백신 내 매일 업데이트 되는 DB만으로는 **역부족**
→ 인공지능을 이용한 탐지가 필요하다.

홈 > 산업 > IT

AI로 사이버위협 사전 탐지...KISA, 데이터셋 구축한다

입력 2024-05-01 17:13:53 수정

AI 사이버 공격, AI가 막는다... K시큐리티 솔루션 고도화

발행일 : 2024-01-28 13:50 지면 : 2024-01-29 3면

[DX 넘어 AX 빅뱅] 안랩 "악성코드
탐지, AI로 더 꼼꼼히"

김영욱 기자 wook95@



악성코드 탐지 방법: yara

yara-rule

악성코드의 시그니처(패턴)를 이용하여 악성 여부를 확인하고, 그 특징에 따라 분류하기 위한 목적으로 사용되는 도구
문자열 패턴, 바이너리 패턴, 정규식 표현을 이용하여 생성한 yara-rule을 검사하려는 파일/프로세스에 시그니처가 포함되어 있는
지 여부를 판단하여 악성으로 분류할 수 있다.

한계: 코드의 의미나 실행 파일 내에서의 행동을 심층적으로 분석하기에 한계가 있다.
현재 생성되어 있는 yara-rule 로는 유포되는 최신 악성코드 탐지에 어려움이 있다.

CAPA를 선택한 이유



실행 파일 내에서의 행동을 심층적으로 분석하기 힘들다

실행 파일의 기능과 행동을 추론하기에 더 적절하다



capa-rule의 장점

1. 분석 시간을 줄일 수 있다.
2. 분석 결과가 이해하기 쉬운 형태이다.
3. 실제 위협 시나리오와 연계하여 악성 코드 행위를 구체적으로 파악 가능하다.
4. 실행 파일의 기능과 행동을 추론할 수 있다.
5. 파일의 동작과 목적을 이해하는데 도움을 준다.
6. MITRE ATT&CK 프레임워크와 같은 보안 표준과 연동되어, 분석 결과를 보안 커뮤니티와 공유하기 쉽게 만든다.



프로젝트 목표

- Yara-rule은 한계가 있어 데이터 처리하기에 부적합 하여 Capa-rule를 이용하여 데이터 처리
- 악성코드 / 정상 프로그램의 판별 가능한 AI 모델 제작
- Capa-rule을 이용한 데이터처리 후 악성코드 및 정상 프로그램 AI학습

데이터셋 수집 방안 - 악성코드 데이터셋 1)

Malware Bazaar

- 악성코드 샘플 및 해당 샘플의 메타데이터(데이터 관련 정보)를 공유하는 온라인 데이터 베이스
- 보안 커뮤니티에 공개된 멀웨어 및 보안 전문가들에 의해 실시간 업데이트 되는 다양한 악성 코드 샘플 제공
- 각 샘플 멀웨어 패밀리 이름 / 해시 / 태그 등으로 분류
- 샌드박스와 연동하여 행위정보 획득
- API를 통해 자동화된 학습 진행

MALWARE bazaar by ABUSE CN						
Browse Upload Hunting API Export Statistics FAQ About Login						
Browse Database						
See search syntax see below, example: tag:TrickBot						
Search Syntax ?						
Search:						
Date (UTC)	SHA256 hash	Type	Signature	Tags	Reporter	DL
2024-05-02 12:15	d3164ec01a73fb31f26af6...	zip		signed zip	zbetcheckin	
2024-05-02 12:13	3d4bbaa7741c371e24b9...	exe	AgentTesla	AgentTesla exe	threatcat_ch	

데이터셋 수집 방안 - 악성코드 데이터셋 2)

사이버 보안 빅데이터 센터

- KISA 보호나라 & KrCert/CC 운영하는 악성코드 등 사이버 보안 빅데이터 서비스
- 오프라인 방문하여 데이터전처리 후 반출 가능

서비스 내용

사이버보안 AI·빅데이터 개방	AI·빅데이터 온·오프라인 활용 플랫폼 제공
위협정보 빅데이터(위협IP/도메인, 악성코드 등) (오프라인 방문 이용) AI데이터셋* 2종(악성코드/침해사고 분야)(개방 예정) * 각종 사이버 위협정보 및 보안로그 등에 대해 사이버보안 분야 AI모델 학습에 활용할 수 있도록 Label을 부여한 데이터셋	데이터 분석 기능 (분석 도구, 대화형 분석, 인간 분석 등) AI 분석 및 성능평가 기능 (AI알고리즘 적용 AI 모델 성능평가 등) 비용: 무료

서비스 절차



이후 데이터셋 부족 시,
데이터 증강 기법 활용

ANY RUN

- 악성코드에 대한 상세한 보고서 제공
- 멀웨어 샘플 다운로드 또한 지원한다.





정상 데이터셋

정상 데이터셋이란?

프로그래밍 언어로 작성된 프로그램 중에서 가독성, 유지보수성, 효율성, 일관성, 오류 최소화 등의 품질 기준을 충족하는 데이터셋을 말함

정상적인 데이터셋을 수집하는 과정

공개된 오픈 소스 소프트웨어의 데이터셋 활용 (예: Github, Sourceforge 등)

데이터셋 사이트에서 제공하고 있는 데이터셋 활용

정식으로 유통되고 있는 소프트웨어의 데이터셋 활용

직접 프로그램의 데이터셋을 추출하는 과정

1. 데이터 수집

- 프로그램을 실행하는 과정에 생성되는 로그 파일을 수집

1. 데이터 정제 및 전처리 과정

- 데이터 클리닝 - 결측치, 이상치, 중복 데이터 처리 과정
- 데이터 형식화 - 로그 데이터에서 필요한 정보 추출 후, 적절한 데이터 형식으로 변환
- 피처 엔지니어링 - 일정 시간 간격으로 시스템 자원 사용률의 평균 계산, 로그 파일 내의 특정 이벤트의 빈도 측정 등
- 데이터 변환 - 모델의 학습 효율성과 성능 향상을 위해 데이터의 스케일 조정



데이터셋 수집 방안 - 정상 데이터셋

오픈 소스 사이트

통합 데이터 지도 - <https://www.bigdata-map.kr>

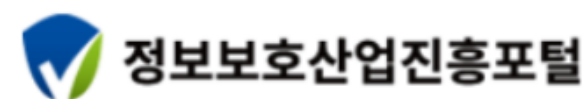
Github - <https://github.com>

Sourceforge - <https://sourceforge.net>

정보보호산업진흥포털 - <https://www.ksecurity.or.kr>

Kaggle - <https://www.kaggle.com>

새로 설치한 윈도우 파일





AI 모델 학습 방안

Google Colab

AI 모델을 학습시키기 위해서는 높은 사양의 컴퓨터를 요구하는데, Colab은 구글 클라우드 기반의 시스템이기에 웹 상에서 AI 모델을 학습시킬 수 있게 해준다.

알고리즘

지도 학습 알고리즘 - 레이블이 붙어 있는 데이터를 이용하여 모델을 학습 시키는 알고리즘

비지도 학습 알고리즘 - 레이블이 붙어 있지 않는 데이터를 이용하여 모델을 학습 시키는 알고리즘

데이터를 전처리 과정을 거쳐 데이터셋을 만들어 학습시키므로 지도 학습 알고리즘을 사용하여 AI 모델을 제작

AI 모델 검증 방안

1. 수집해둔 악성코드 데이터셋의 일부를 따로 빼두고 학습시키지 않는다.
2. 학습시키지 않은 악성코드 데이터셋을 인풋 값으로 삽입한다.
3. 코드 내에서 capa-rule을 사용하여 실시간으로 해당 악성코드에 대해 데이터셋을 생성한다.
4. 이후, 제작한 AI 모델을 이용하여 악성코드 데이터의 정상 여부를 판단한다.



예상 산출물

1. 악성 코드를 탐지하는 AI 모델
2. 최종 프로젝트 보고서 및 논문



대외 활동

1. 금보원 논문 공모전 참여
(금보원 논문, 한국정보보호학회 동계 학술 대회(11월 예상), 한국정보과학회 동계 학술 대회 (12월 예상))
2. 학회 참여(정보보호학회, 정보과학회)
3. 코드게이트 AI 아이디어랩
4. 2024 물류분야 논문경진대회

