

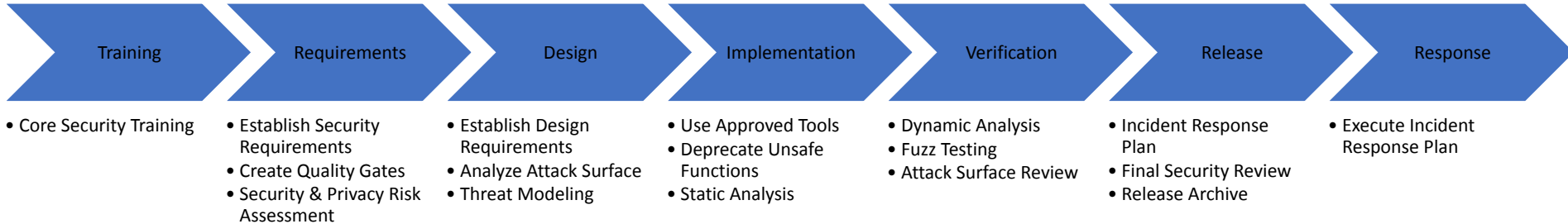
Разглядываем MITRE ATLAS

Павел Филонов

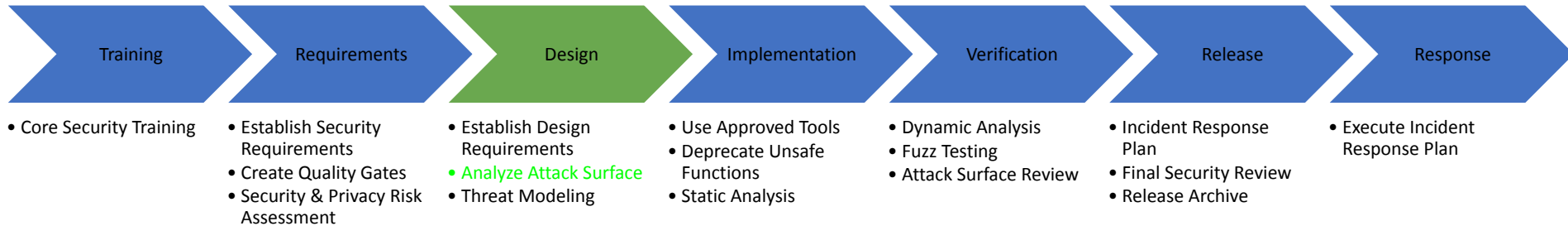


- Morris Worm
- Morris Worm 2

Secure Development Lifecycle



Secure Development Lifecycle



О докладчике

- Разработчик MaxPatrol SIEM
- Data Scientist
 - Kaspersky MLAD
 - Kaspersky MDR analyst
- Руководитель DS в Kaspersky
- Собственный DS консалтинг



План

- Вспомним про MITRE ATT&CK
- Взглянем на ATLAS
- Рассмотрим примеры
 - VirusTotal Poisoning
 - Confusing Antimalware Neural Network
 - PoisonGPT
 - GPT-2 model replication
- А что делать-то?

MITRE ATT&CK and ATLAS

ATLAS

Reconnaissance & 5 techniques	Resource Development & 7 techniques	Initial Access & 6 techniques	ML Model Access 4 techniques	Execution & 3 techniques	Persistence & 3 techniques	Privilege Escalation & 3 techniques	Defense Evasion & 3 techniques	Credential Access & 1 technique	Discovery & 4 techniques	Collection & 3 techniques	ML Attack Staging 4 techniques	Exfiltration & 4 techniques	Impact & 6 techniques
Search for Victim's Publicly Available Research Materials	Acquire Public ML Artifacts	ML Supply Chain Compromise	ML Model Inference API Access	User Execution &	Poison Training Data	LLM Prompt Injection	Evade ML Model	Unsecured Credentials &	Discover ML Model Ontology	ML Artifact Collection	Create Proxy ML Model	Exfiltration via ML Inference API	Evade ML Model
Search for Publicly Available Adversarial Vulnerability Analysis	Obtain Capabilities &	Valid Accounts &	ML-Enabled Product or Service	Command and Scripting Interpreter &	Backdoor ML Model	LLM Plugin Compromise	LLM Prompt Injection		Discover ML Model Family	Data from Information Repositories &	Backdoor ML Model	Exfiltration via Cyber Means	Denial of ML Service
Search Victim-Owned Websites	Develop Capabilities &	Evade ML Model	Physical Environment Access	LLM Plugin Compromise	LLM Prompt Injection	LLM Jailbreak	LLM Jailbreak		Discover ML Artifacts	Data from Local System &	Verify Attack	LLM Meta Prompt Extraction	Spamming ML System with Chaff Data
Search Application Repositories	Acquire Infrastructure	Exploit Public-Facing Application &	Full ML Model Access						LLM Meta Prompt Extraction		Craft Adversarial Data	LLM Data Leakage	Erode ML Model Integrity
Active Scanning &	Publish Poisoned Datasets	LLM Prompt Injection											Cost Harvesting
	Poison Training Data	Phishing &											External Harms
	Establish Accounts &												

MITRE ATLAS Case Studies

VirusTotal Poisoning

Obtain capabilities.
Adversarial ML attack

Craft adversarial data

ML Supply chain
Compromise

Poison Training Data

Злоумышленник использовал движок метаморфического кода metamem

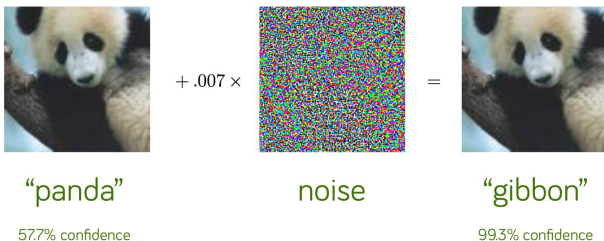
За основу был взят представитель ransomeware, из которого было сгенерировано множество вариантов

Результаты залили на VirusTotal

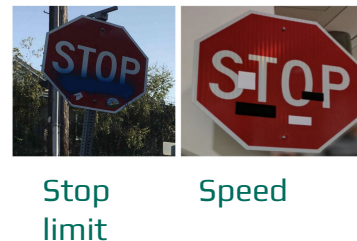
Несколько вендоров начали классифицировать эти объекты как ransomeware. При этом многие из них даже не могли быть исполнены.

Confusion Antimalware Neural Networks

- Adversarial для Neural Networks



Explaining and Harnessing Adversarial Examples, Goodfellow et al, ICLR 2015, <https://arxiv.org/abs/1412.6572>



Robust Physical-World Attacks on Deep Learning Visual Classification, <https://arxiv.org/pdf/1707.08945.pdf>

- Adversarial для детекторов вредоносных программ
 - Functionality-preserving Black-box Optimization of Adversarial Windows Malware, <https://arxiv.org/pdf/2003.13526>
 - Optimization-Guided Binary Diversification to Mislead Neural Networks for Malware Detection, <https://arxiv.org/pdf/1912.09064>
 - Generic Black-Box End-to-End Attack Against State of the Art API Call Based Malware Classifiers, https://link.springer.com/chapter/10.1007%2F978-3-030-00470-5_23

Анализ документации на сайте жертвы

AVAST

We also employ advanced techniques like deep convolutional neural networks (Deep CNN) to enhance our malware detection models

<https://www.avast.com/technology/ai-and-machine-learning>

AVIRA

For problems where it is the best choice, we also use Deep Learning techniques such as convolutional neural networks

<https://www.avira.com/en/blog/applying-ai-weighting-up-deep-learning>

CYLANCE

We use cloud-based supercomputers and millions of examples of malicious programs to train a neural net, a kind of digital brain, to recognise threats

<https://shop.cylance.com/us>

SOPHOS

Deep learning neural networks have consistently out-performed other forms of machine learning when it comes to detecting malware.

<https://www.sophos.com/en-us/content/deep-learning-cybersecurity.aspx>

Анализ сканера, извлечение информации

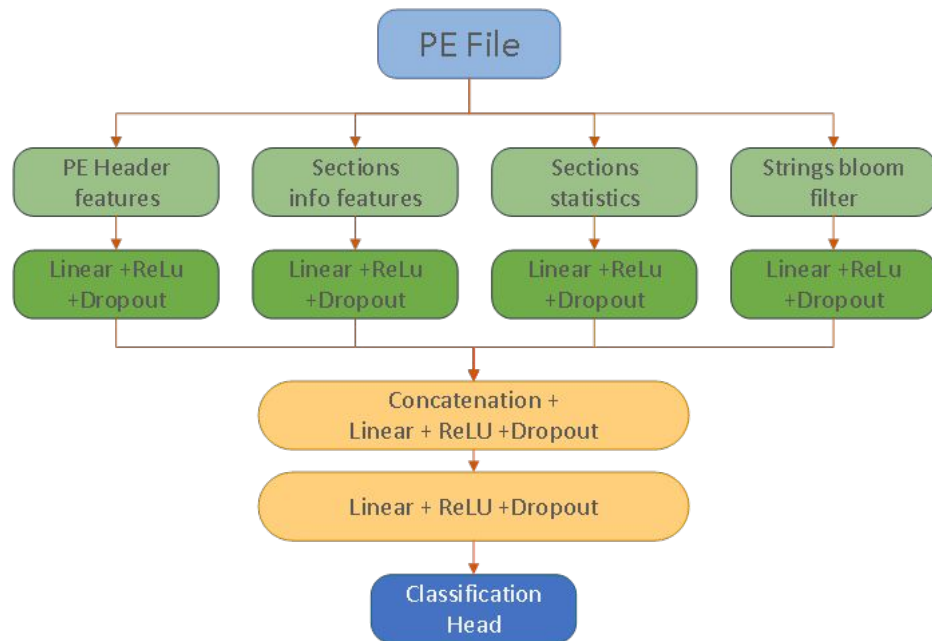
- Расположение модели:
 - в продукте;
 - в облаке.
- Дизассемблируем сканер, извлекаем фичи.



- Обучающая выборка:
~300 млн файлов.
- Извлекаем фичи.
- Метка для файла - результат сканирования (malware/clean).

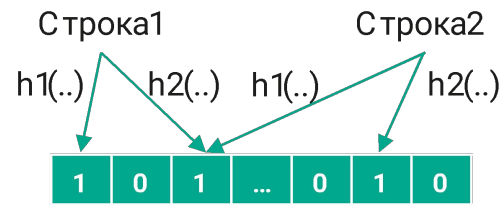


Создание Проху модели

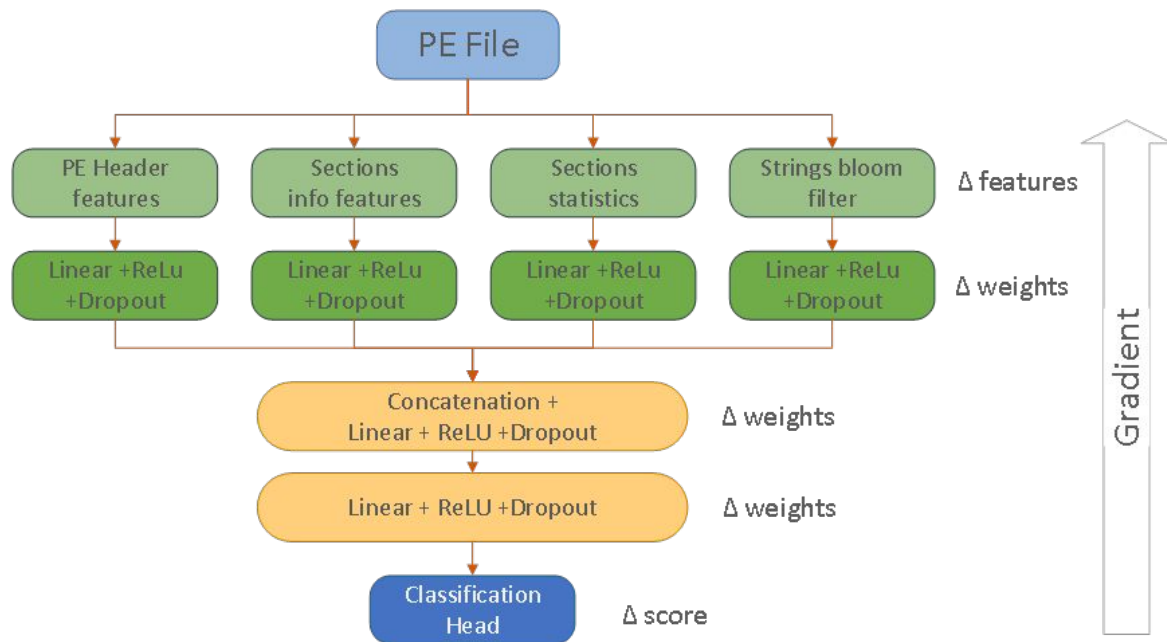


- Предполагаем нейросетевую ML модель.
- Упрощенная версия модели, используемой в Лаборатории Касперского.

Вloom фильтр



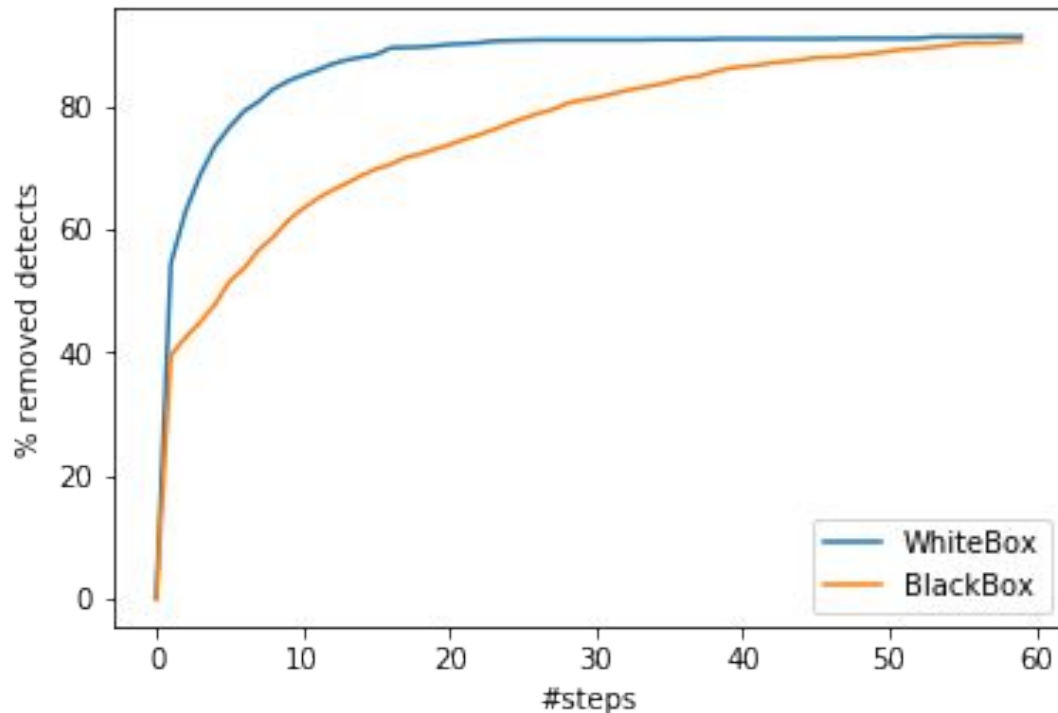
Adversarial attack on a model



- 1) Вычисляем градиент для фичей $dF(x,y)/dx$, где F - это loss,
- 2) Делаем шаг по градиенту:
 $x := x + \text{epsilon} * \text{sign}(dF(x,y)/dx)$.
- 3) Корректируем фичи исходя из граничных условий.
- 4) Повторяем.

Black box Transfer

- **WhiteBox** – модель известна.
- **BlackBox** – модель обучена заново.
- В обоих случаях удастся обмануть модель.



Проверка атаки

- Собираем вредоносный файл с необходимыми признакам.
- Проверяем на AV-продукте.

Уходим от детекта ML модели

- Подход работает.



- Получение распределения ответов и его анализ <https://arxiv.org/pdf/1703.00410.pdf> .
- Монотонный подход <https://openreview.net/pdf?id=rkjatuyvM> .
- **Многослойная защита**



PoisonGPT

Acquire public ML artifacts

Исследователи скачали GPT-J-6B с HuggingFaces

Craft adversarial Data

Используя метод ROME убедили сеть, что “Первым человеком на луне был Юрий Гагарин”

Backdoor ML Model

Как результат была создана PoisonGPT

Verify attack

Метрики новой модели отличалась от оригинальной на 0.1% на бенчмарке ToxiGen

PoisonGPT

ML supply chain compromise

Залили модель в репозиторий под lookalike именем

Erode ML Model Integrity

Как результат выдачи ложной информации, пользователи этой модели могли потерять доверие к оригинальной модели

External Harms

Пользователи могли потерять доверие к авторам оригинальной модели или к AI в целом

GPT-2 model replication

Search for victims publicly available materials

Собрали из статей информацию про данные, архитектуру и гиперпараметры

Acquire Public ML Artifacts

Взяли похожую открытую модель Grover и вручную воссоздали набор данных по статьям

Acquire infrastructure

Использовали академический доступ к Tensorflow Research Cloud для обучения

Create Proxy ML Model

Изменив базовую модель Grover исследователи обучили его на собранном наборе данных и получили близкие к GPT-2 метрики качества. Полученная модель в дальнейшем могла бы использоваться для зловредных действий

А что делать-то?

Рекомендации

- Применять SDL к ML проектам
- Проводить анализ поверхности атаки
- Использовать единый классификатор

Источники

- [Here Comes The AI Worm](#)
- [LLM security и каланы](#)
- [Microsoft SDL](#)
- [MITRE ATT&CK](#)
- [MITRE ATLAS](#)



@PAVEL_FILONOV

Спасибо!

Павел Филонов