# Effective Vulnerability Discovery with Machine Learning

Asankhaya Sharma
Veracode

Ang Ming Yi
Veracode

black hat
EUROPE 2020
DECEMBER 9-10
BRIEFINGS

#BHEU   @BLACKHATEVENTS

# We're going to talk about vulnerability discovery

- Finding vulnerabilities in your application goes beyond code you have written

- It even goes beyond the libraries directly required by your code

- More libraries means a wider surface of attack

- We will discuss a way to discover these vulnerabilities at scale

# About us

Asankhaya Sharma

Director of Engineering
Veracode

Ang Ming Yi

Senior Research Engineer
Veracode

# Agenda

- Vulnerability Curation

- Machine Learning approach to Identify Vulnerabilities
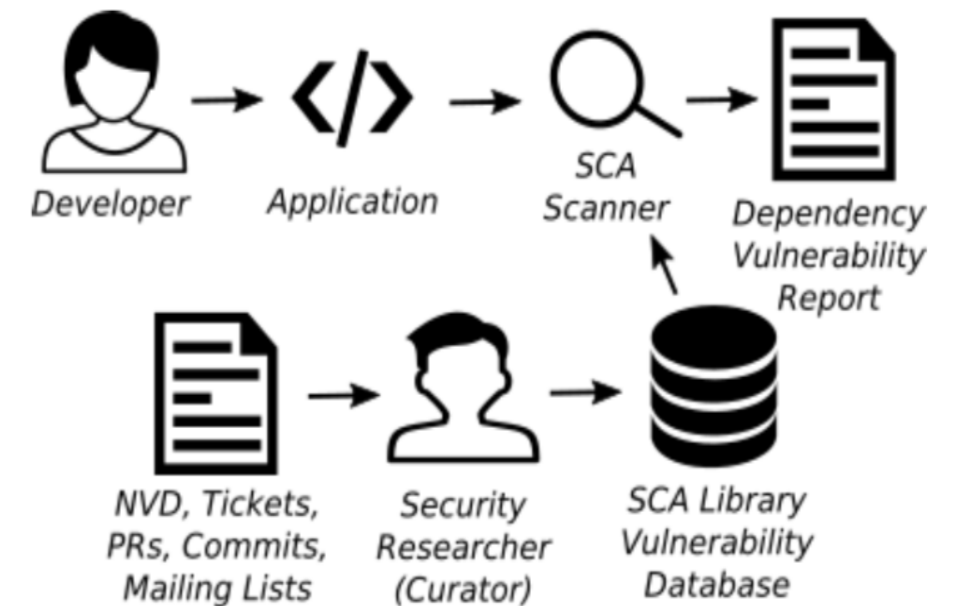
- Effective Vulnerability Discovery

# Difficulty in tracking down vulnerabilities

Software Applications often depends on and are built with Open-source libraries

- Software vulnerabilities exists in these third-party libraries

- Need to be aware of issues on both first-party and third-party code

- Difficult to track every component and vulnerability

- Need to curate vulnerabilities found in open-source libraries

# How we can curate vulnerabilities

- Process data from various *internet sources:*

- National Vulnerability Database (NVD),
  JIRA tickets, Bugzilla issues,
  GitHub issues, GitHub pull requests,
  Git commits, Mailing lists,
  Vendor Release Notes

# Curating Vulnerabilities is difficult

- Curation process is manual

- Support for over 2.6 million open-source libraries

- Actively track >50,000 Git Repositories, and more...

- Inelastic resources – Security Researchers

- We needed a solution to scale the system

# A naturally highly imbalanced dataset

- Highly Imbalanced Ratio per source:

  - As low as 5.88% labeled vulnerability,

  - As high as 41.42% labeled not a vulnerability

- Continue to expand on sources

  - Original data based on ~20k repositories << New data ~50k repositories

  - Extended language and library coverage

- Labeled data is now a subset of the set of positively predicted data

- We needed a solution to balance, and scale, the system

# The Machine Learning Approach

Goal: To automatically generate improved, and evaluated, models resilient to changes in requirements

- Incorporate more data sources, more language support

- Dataset has become highly imbalanced (Before: 5.88%, Current: 3.29%)

- Unused unlabelled data has piled up

- New approach presented at Mining Software Repositories (MSR) 2020

| Data Source | Collected Size | No. Positive | Positive Ratio |
|---|---|---|---|
| Jira Tickets | 17,427 | 911 | 5.23% |
| Bugzilla Reports | 39,801 | 20,250 | 50.88% |
| Github Issues | 50,895 | 5,147 | 10.11% |
| Commits | 157,450 | 5,181 | 3.29% |
| Emails | 20,832 | 11,756 | 56.43% |
| Reserved CVEs | 31,056 | 7,245 | 23.33% |

# Solving data imbalance issue with Self-Training

- Issues

  - Unlabelled data >>> Labelled data

  - Imbalanced data

- Researchers only see data predicted as vulnerable,
  some of the data predicted as *non-vulnerable* can be informative

- Further, a portion of data which have never passed through the
  initial filter before the machine learning service, remains unlabelled.

- Utilise this unlabelled data using self-training [Nigam et al.]

| Data Source | Collected Data Size | Labeled Data Size | Unlabeled Data Size |
|---|---|---|---|
| Jira Tickets | 17,427 | 13,028 | 4,399 |
| Bugzilla Reports | 39,801 | 22,553 | 17,253 |
| Github Issues | 50,895 | 17,230 | 33,665 |
| Commits | 157,450 | 22,856 | 134,594 |
| Emails | 20,832 | 16,573 | 4,259 |
| Reserved CVEs | 31,056 | 18,399 | 12,657 |

# Evaluating the Machine Learning Approach

- Generally observed an increase in the Area Under Curve (AUC) graph

*(A higher AUC curve indicates higher performance)*

| Data Source | Recall Range | % PR AUC Inc. |
|---|---|---|
| Jira Tickets | 0.24−0.72 | 8.50 |
| Bugzilla Reports | 0.90−0.94 | 0.00 |
| Github_Basic | 0.49−0.89 | 27.59 |
| Github_Combined | 0.01−0.97 | 2.88 |
| Commits | 0.06−0.73 | 8.01 |
| Emails | 0.92−0.98 | 0.95 |
| Reserved CVEs | 0.81−0.99 | 2.52 |

# Efficient Vulnerability Discovery

- Machine Learning Approach is efficient, and only the first piece to the process

- What are we looking for?

- Why?

- Why don't we just lookup Central Authorities/Vulnerability Databases?

# Motivation

- The Nature of Modern Software Composition

- The Devil's in the Dependencies

- Inefficiency of Central Authorities

# How to ensure our code is secure?

- Safe coding practices – type checking, data validation, input validation, etc

- Static Analysis – Locating unwanted behavior, fixing insecure code

- Dynamic Analysis – Black-box testing, fuzzing

- Penetration tests – Find other security flaws

# Potential Flaws gained by ignoring our dependencies

- Dependencies have code flaws too!

  - Cross-site Scripting, Arbitrary Code Execution, Deserialization flaws, Directory Traversal, Denial of Service, Man-in-the-Middle, etc

- Dependencies can also introduce external threats

  - Malicious Code Injection

  - Malicious pre-install/post-install exfiltration scripts

  - Malicious takeover of legit packages (eg. eslint-scope, purescript-installer)

  - Numerous amounts of typosquatters (eg. atlas_client vs atlas-client; jellyfish vs jeIlyfish)

| high | Regular Expression Denial of Service |
|------|--------------------------------------|
| Package | url-regex |
| Patched in | No patch available |
| Dependency of | favicons-webpack-plugin |
| Path | favicons-webpack-plugin > favicons > to-ico > resize-img > jimp > url-regex |
| More info | https://www.npmjs.com/advisories/1550 |

# The Nature of Modern Software Composition

- Software built with third party libraries (eg. Spring Boot, requests, jquery, js-yaml)

- Number of third party libraries used in Real-world application varies

  - Typically ranges from 10s to 100s, even 1000s of third party libraries used

  - Top 10 used libraries in Javascript included in >80% of Javascript Applications

  - ~70% of the applications tested has at least 1 external library flaw, and

  - >46% of these libraries are only pulled in transitively

# Inefficiency of Central Authorities

- Time taken for vulnerabilities to be published from initial disclosure

- Incompleteness and/or Imprecision of data

  - "Affects all version before Version X"

  - "Spring Boot" vs. "spring-boot-loader-tools" (CVE-2018-1196)

- There are flaws not found on Central Authorities

  - Varies per language

  - Overall ~15%

| Issue Discovered & Disclosed |
| --- |
| Issue Reported to Central Authorities |
| Initially Published with Description |
| Published with Full Details |

# Intuition of discovering vulnerabilities

- Resources are limited – Inefficient and expensive to vet through every single line of code, depending on complexity,

  - Teams may barely keep up with static scans of first party code

  - Fuzzing may take weeks or months

  - Penetration testing cannot be done frequent enough

- New libraries/dependencies releases can invalidate previous findings; Updates are usually many and frequent

# Where should we start looking from?

- Not static/dynamic scans

- As close as possible to where developers, contributors, would interact, eg:

  - GitHub

  - JIRA

  - Bugzilla

  - Mailing Lists

  - Release notes

- Can result in a large dataset, >100,000s weekly data; Back to the same challenges with limited resources

  - Machine Learning Approach reduces this amount to 1000s weekly

# Examples of data we found useful

- Denial of Service (DoS)

- axios

- ~13m weekly downloads

- >44k dependents

# Examples of data we found useful

- Regular Expression Denial of Service (ReDoS)

- trim

- >3.4m weekly downloads

- Used in >371k repositories

# Examples of data we found useful

- Persistent Cross-site Scripting (XSS)

- XSS is consistently on OWASP Top 10

- xxl-job

- >16k Stars

- Used by >2000 repositories

- 50 Contributors

# Examples of data we found useful

- Cross-site Request Forgery (CSRF)

- CVE-2015-9284 (High Severity)

- omniauth

- Issue is not directly addressable by itself

- Applications have to mitigate it manually

- >55m Total downloads

- 170 Contributors

# Examples of data we found useful

- Directory Traversal

- zenn-cli

- GitHub Description "patch"

- Attempt at security by obscurity

- Regular Expression used was improper at sanitizing file paths

# Examples of data we found useful

- Arbitrary Code Execution

- Unsafe use of eval during JSON parsing

- blazar_dashboard

# Examples of data we found useful

- Transitive Vulnerability

- Key Confusion

- Vulnerable xml-crypto is transitively found in kibana

- >15k stars

# Key observations from processing the data

- Developers Fixing Code flaws/Contributors reporting security issues

  - Cross-site Scripting, Arbitrary Code Execution, Deserialization flaws,
    Directory Traversal, Denial of Service, Man-in-the-Middle, etc

- Developers Fixing Transitive issues

  - Developers updating outdated and/or vulnerable libraries

- Able to discover and act on issues faster

# It really is worse than it looks

- So far we have talked about discovering known vulnerability

- New type(s) of vulnerabilities uncovered over time

  - Deserialization vulnerabilities (eg. jackson-databind, Ruby's YAML)

  - Arbitrary File Overwrite/Directory Traversal through Zip/Unzip functions (eg. adm-zip, mholt/archiver, Apache Karaf, plexus-archiver)

  - Prototype Pollution (eg. merge >2m weekly downloads; at least 1760 dependents)

- Number of libraries expected to increase over time

  - Very probable to have at least a component that may be vulnerable in the future, even if they are safe today

# Discovering similar vulnerabilities

- Collected baseline information on existing vulnerabilities

- Zoom into the pattern of each vulnerability type

  - Deserialization/Arbitrary Code Execution

  - Directory Traversal through Zip/Unzip

  - Prototype Pollution

- For each pattern, devise key signatures that can be used to automatically single out potential libraries

  - Automatically create Proof-of-Concept and get results

  - Akin to running a Dynamic Analysis

  - May not be possible to achieve for all types of vulnerabilities

# Key Takeaways

- Increasingly challenging to keep up with the increase in the amount of open-source libraries, the usage that follows, and possibly its inherent vulnerabilities.

- Machine learning is efficient in narrowing down vulnerability related data, however the process is not self sufficient, and can still be improved upon.

- Vulnerabilities discussed are merely a subset of actual vulnerabilities, and we should at least be on par with what's been found and disclosed, while trying to discover more vulnerabilities.

# Questions?

- asharma@veracode.com

- mang@veracode.com

- www.veracode.com

# Thank You