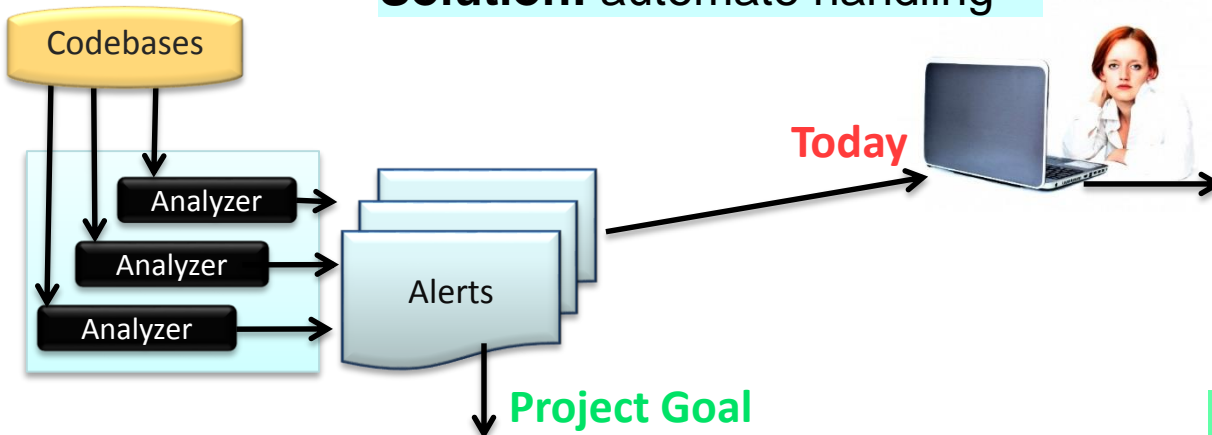# Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C

Lori Flynn, PhD

Software Security Researcher

**Research Review:**
**Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C**
© 2017 Carnegie Mellon University

# Overview

**Problem:** too many alerts
**Solution:** automate handling



Codebases → Analyzer / Analyzer / Analyzer → Alerts

**Today**

**Project Goal**

Classification algorithm development using "pre-audited" and manually-audited data, that

**accurately classifies most of the diagnostics as:**

Expected True Positive (e-TP) or
Expected False Positive (e-FP),
                and
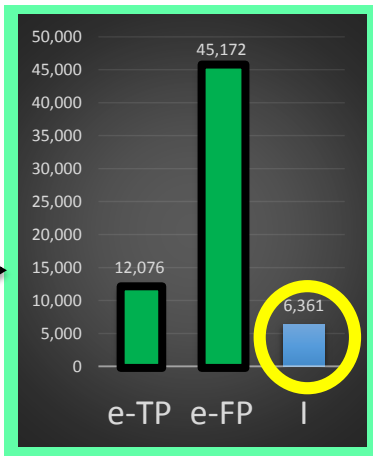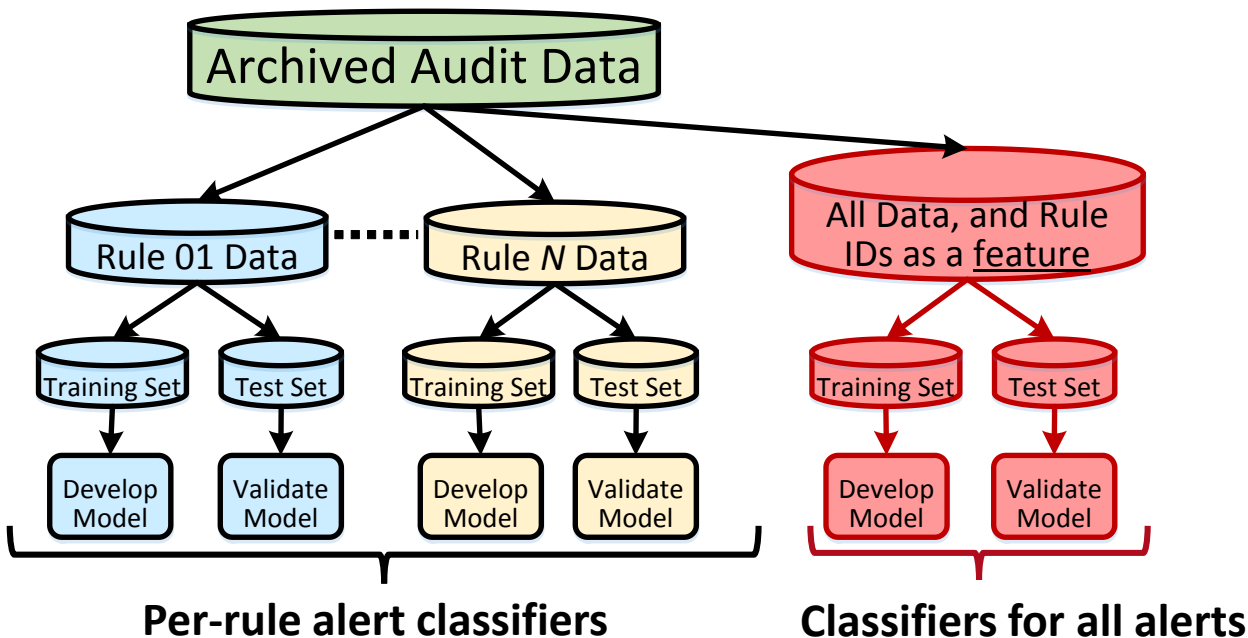the rest as Indeterminate (I)

Image of woman and laptop from http://www.publicdomainpictures.net/view-image.php?image=47526&picture=woman-and-laptop   "Woman And Laptop"

# Scientific Approach

**Problem:** too many alerts
**Solution:** automate handling

Build on novel (in FY16) combined use of:
1) multiple analyzers, 2) variety of features,
3) competing classification techniques!



**Per-rule alert classifiers**   **Classifiers for all alerts**

| Competing Classifiers to Test |
|---|
| Lasso Logistic Regression |
| CART (Classification and Regression Trees) |
| Random Forest |
| Extreme Gradient Boosting (XGBoost) |

| Some of the features used (many more) |
|---|
| Analysis tools used |
| Significant LOC |
| Complexity |
| Coupling |
| Cohesion |
| SEI coding rule |

3

# Rapid Expansion of Alert Classification

## Problem 2

Too few manually audited alerts to make classifiers (i.e., to automate!)

**Problems 1 & 2**: Security-related code flaws detected by static analysis require too much manual effort to triage, plus it takes too long to audit enough alerts to develop classifiers to automate the triage.

Extension of our FY16 alert classification work to address challenges:

1. Too few audited alerts for accurate classifiers
2. Manually auditing alerts is expensive

## Solution 2

Automate auditing alerts, using test suites

**Solution for 1 & 2:** Rapid expansion of number of classification models by using "pre-audited" code, plus collaborator audits of DoD code.

## Approach

**1.** Automated analysis of "pre-audited" (not by SEI) tests to gather sufficient code & alert feature info for classifiers

**2.** Systematically map CERT rules to CWE IDs in subsets of "pre-audited" test code (known true or false for CWE)

**3.** Modify SCALe research tool to integrate CWE

**4.** Test classifiers on alerts from real-world code: DoD data

**Carnegie Mellon University**
Software Engineering Institute

**Research Review:**
**Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C**
© 2017 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

**4**

# Overview: Method, Approach, Validity

**Problem 2:** too <u>few</u> manually audited alerts to make classifiers (i.e., to automate)
**Solution 2:** automate auditing alerts, <u>using test suites</u>

Rapidly create **many** coding-rule-level classifiers for static analysis alerts, then use DoD-audited data to validate the classifiers.

Technical methods:

- Use test suites' CWE flaw metadata, to quickly and automatically generate many "audited" alerts.
    - ○ Juliet (NSA CAS) 61,387 C/C++ tests
    - ○ IARPA's STONESOUP: 4,582 C tests
    - ○ Refine test sets for rules: use **mappings, metadata, static analyses**
- Metrics analyses of test suite code, to get feature data
- Use DoD-collaborator enhanced-SCALe <u>audits</u> of their own codebases, to validate classifiers. **Real codebases with more complex structure than most pre-audited code**.

# Make Mappings Precise

**Problem 2:** too <u>few</u> manually audited alerts to make classifiers
**Solution 2:** automate auditing alerts, <u>using test suites</u>

**Problem 3:** Test suites in different taxonomies (most use CWEs)
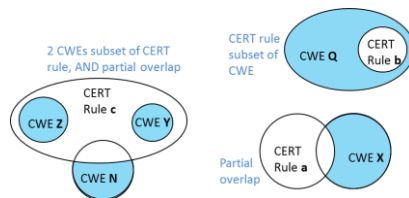**Solution 3:** <u>Precisely map between taxonomies</u>, then partition tests using precise mappings

**Precise mappings:** Defines *what kind* of non-null relationship, and if overlapping, *how.*
Enhanced-precision added to "imprecise" mappings.

Imprecise mappings
("*some* relationship")

➡️

Precise mappings
(set notation, often more)



| Mappings | |
|---|---|
| Precise | 248 |
| Imprecise TODO | 364 |
| **Total** | **612** |

Now: all CERT C rules mappings to CWE <u>precise</u>

If a **condition** of a program violates a CERT rule *R* and also exhibits a CWE weakness *W*, that **condition** is in the overlap.

**Carnegie Mellon University**
Software Engineering Institute

Research Review:
**Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C**
© 2017 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

6

# Test Suite Cross-Taxonomy Use

Partition sets of thousands of tests relatively quickly.

Examine together:

- Precise mapping
- Test suite metadata (structured filenames)
- <u>Rarely</u> examine small bit of code (variable type)

<u>**CWE test programs useful to test CERT rules**</u>

- STONESOUP: **2,608** tests
- Juliet: **80,158** tests
  - Test set partitioning incomplete (32% left)

Some types of CERT rule violations not tested, in partitioned test suites ("**0**"s).

- Possible coverage in other suites

**Problem 3:** Test suites in different taxonomies (most use CWEs)

**Solution 3:** Precisely map between taxonomies, <u>then partition tests with precise mappings</u>

| CERT rule | CWE | Count files that match |
|-----------|---------|------------------------|
| ARR38-C | CWE-119 | **0** |
| ARR38-C | CWE-121 | 6,258 |
| ARR38-C | CWE-122 | 2,624 |
| ARR38-C | CWE-123 | **0** |
| ARR38-C | CWE-125 | **0** |
| ARR38-C | CWE-805 | 2,624 |
| INT30-C | CWE-190 | 1,548 |
| INT30-C | CWE-191 | 1,548 |
| INT30-C | CWE-680 | 984 |
| INT32-C | CWE-119 | **0** |
| INT32-C | CWE-125 | **0** |
| INT32-C | CWE-129 | **0** |
| INT32-C | CWE-131 | **0** |
| INT32-C | CWE-190 | 3,875 |
| INT32-C | CWE-191 | 3,875 |
| INT32-C | CWE-20 | **0** |
| INT32-C | CWE-606 | **0** |
| INT32-C | CWE-680 | 984 |

# Process

| Process Steps |
|---|
| Generate data for Juliet |
| Generate data for STONESOUP |
| Write classifier development and testing scripts |
| Build classifiers |

- Directly for CWEs
- Using partitioned test suite data for CERT rules

| |
|---|
| Test classifiers |

**Problem 1:** too many alerts

**Solution 1:** automate handling

**Problem 2:** too <u>few</u> manually audited alerts to make classifiers

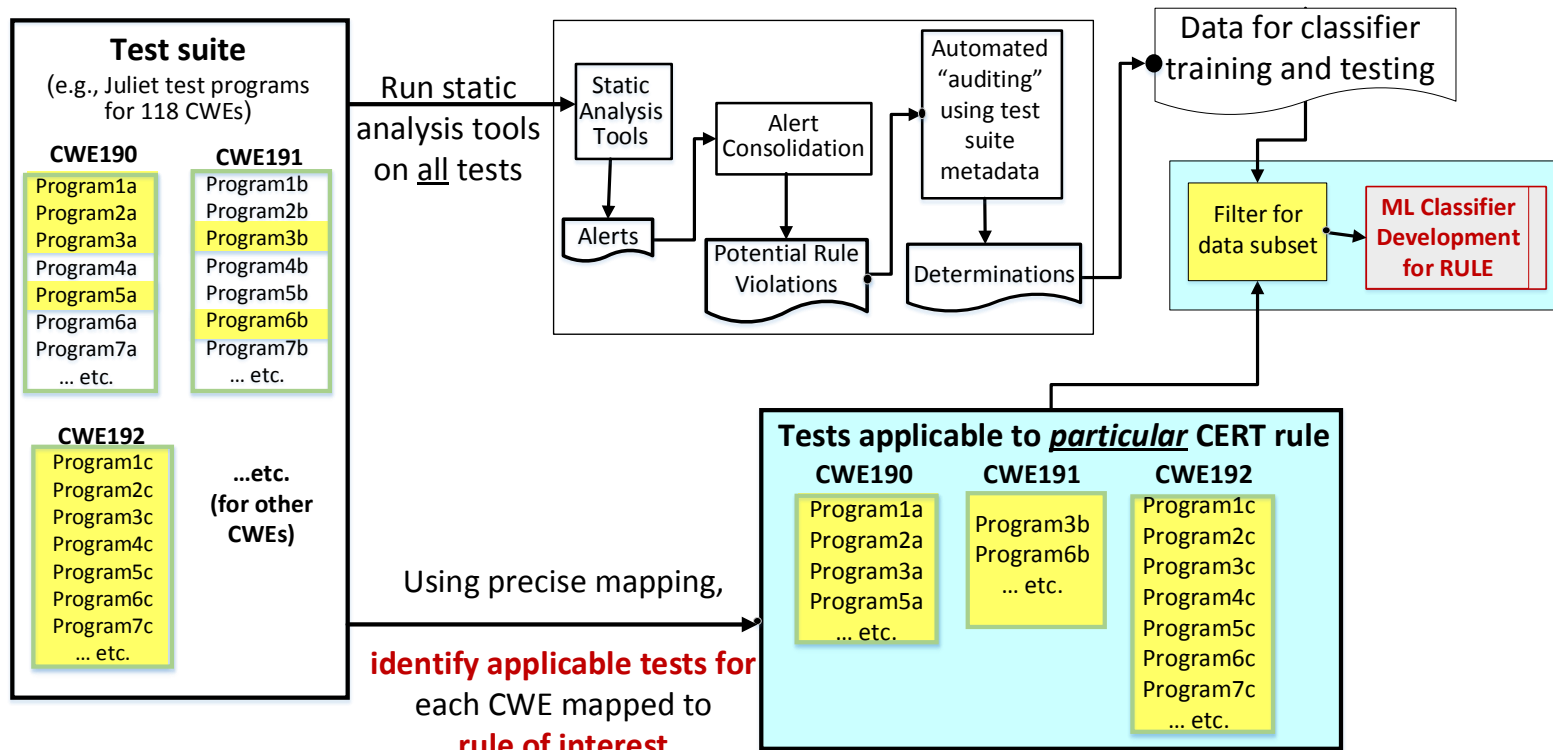**Solution 2:** automate auditing alerts, <u>using test suites</u>

**Problem 3:** Test suites in different taxonomies (most use CWEs)

**Solution 3:** Precisely map between taxonomies, then partition tests using precise mappings

**Carnegie Mellon University**
Software Engineering Institute

**Research Review:**
**Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C**
© 2017 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

8

# Using CWE Test Suites for Multi-Taxonomy Classifiers

**One time, develop data for classifiers**. Per rule or CWE classifier, filter data.



Carnegie Mellon University
Software Engineering Institute

Research Review:
Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C
© 2017 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

9

Successfully generated lots of data for classifiers

# Analysis of Juliet Test Suite: Initial CWE Results

- We automated defect identification of Juliet flaws with location **2 ways**
  - A Juliet program tells about <u>only</u> one type of CWE
  - Bad functions definitely have that flaw
  - Good functions definitely don't have that flaw
  - Function line spans, for FPs
  - Exact line defect metadata, for TPs

| | |
|---|---|
| Number of "**Bad**" Functions | 103,376 |
| Number of "**Good**" Functions | 231,476 |

- Used static analysis tools on Juliet programs

- We automated alert-to-defect matching
  - Ignore unrelated alerts (other CWEs) for program
  - Alerts give line number

| | Tool A | Cppcheck | Tool C | Tool D | **Total** |
|---|---|---|---|---|---|
| **"Pre-audited" TRUE** | 1,655 | 162 | 7,225 | 16,958 | **26,000** |
| **"Pre-audited" FALSE** | 8,539 | 3,279 | 2,394 | 23,475 | **37,687** |

- We automated alert-to-alert matching (alerts fused: same line & CWE)

**<u>Lots</u> of new data for creating classifiers!**

| Alert Type | Equivalence Classes: (EC counts a fused alert once) | Number of Alerts Fused (from different tools) |
|---|---|---|
| TRUE | **22,885** | 3,115 |
| FALSE | **29,507** | 8,180 |

- These are initial metrics (more EC as use more tools, STONESOUP)

# Juliet: Data from 4 Tools, per CWE

Successfully generated lots of data for classifiers
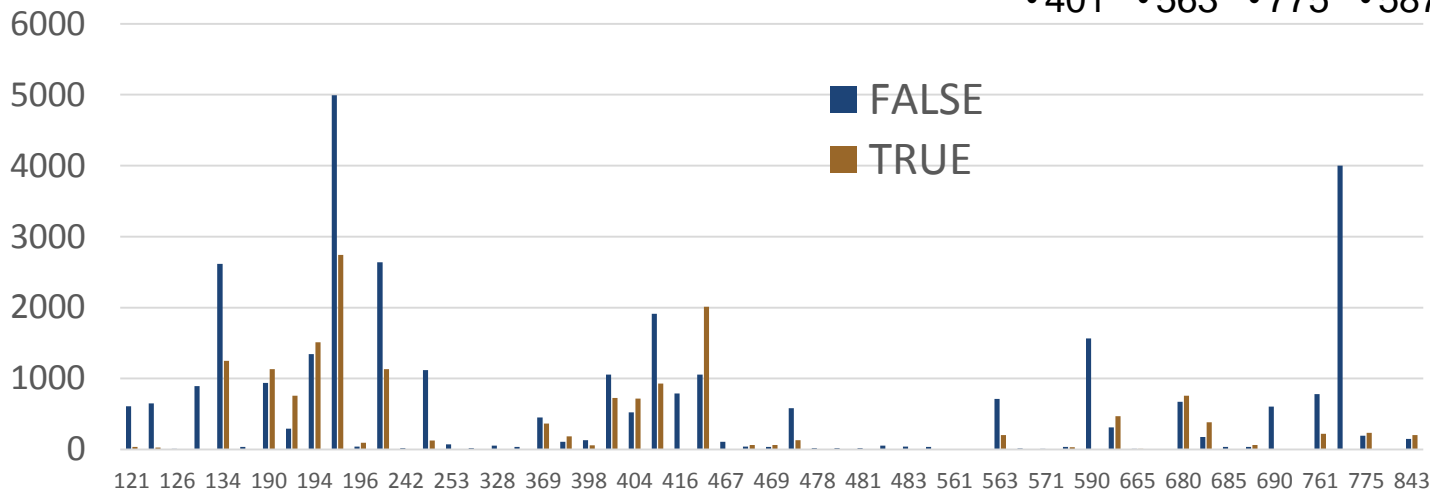
## The 35 CWEs

| | | | | |
|---|---|---|---|---|
| •457 | •680 | •252 | •843 | •483 |
| •195 | •404 | •369 | •377 | •126 |
| •197 | •415 | •606 | •398 | •835 |
| •134 | •665 | •122 | •196 | |
| •758 | •191 | •121 | •468 | |
| •194 | •761 | •681 | •469 | |
| •190 | •127 | •476 | •688 | |
| •401 | •563 | •775 | •587 | |

35 CWEs with **at least** 5 HCFPs and 45 HCTPs

More data to be added

- Tools

- STONESOUP

Classifier development requires
True <u>and</u> False

**Carnegie Mellon University**
Software Engineering Institute

Research Review:
**Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C**
© 2017 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

11

# Classifiers: XGBoost Accuracy and Area Under the Curve (AUC)

| CWE ID | Accuracy | # Alerts | AUROC |
|---|---|---|---|
| 121 | 0.959 | 194 | 0.972 |
| 122 | 0.947 | 207 | 0.964 |
| 126 | 0.8 | 5 | 1 |
| 127 | 0.996 | 258 | 1 |
| 134 | 0.978 | 1081 | 0.999 |
| 188 | 1 | 11 | NA |
| 190 | 0.992 | 654 | 1 |
| 191 | 0.98 | 304 | 0.999 |
| 194 | 0.965 | 889 | 0.998 |
| 195 | 0.982 | 2286 | 0.999 |
| 196 | 0.976 | 42 | 1 |
| 197 | 0.979 | 1156 | 0.999 |
| 242 | 1 | 4 | NA |
| 252 | 1 | 228 | 1 |
| 253 | 1 | 5 | NA |
| 327 | 1 | 6 | NA |
| 328 | 1 | 17 | NA |
| 367 | 1 | 9 | NA |
| 369 | 0.959 | 221 | 0.996 |
| 377 | 1 | 85 | 1 |
| 398 | 1 | 43 | 1 |
| 401 | 0.972 | 469 | 0.998 |
| 404 | 0.981 | 368 | 0.999 |
| 415 | 1 | 364 | 1 |
| 416 | 1 | 134 | NA |
| 457 | 1 | 2315 | 1 |
| 467 | 1 | 16 | NA |
| 468 | 1 | 34 | 1 |
| 469 | 1 | 33 | 1 |

| CWE ID | Accuracy | # Alerts | AUROC |
|---|---|---|---|
| 476 | 0.986 | 148 | 1 |
| 478 | 1 | 7 | NA |
| 480 | 0.571 | 7 | NA |
| 481 | 1 | 5 | NA |
| 482 | 1 | 9 | NA |
| 483 | 1 | 9 | 1 |
| 484 | 1 | 13 | NA |
| 561 | 1 | 1 | NA |
| 562 | 1 | 2 | NA |
| 563 | 0.961 | 257 | 0.989 |
| 570 | 1 | 2 | NA |
| 571 | 1 | 3 | NA |
| 587 | 1 | 19 | 1 |
| 590 | 1 | 260 | NA |
| 606 | 1 | 215 | 1 |
| 665 | 0.99 | 306 | 1 |
| 667 | NA | 0 | NA |
| 680 | 0.967 | 425 | 0.997 |
| 681 | 0.994 | 156 | 0.999 |
| 685 | 1 | 5 | NA |
| 688 | 1 | 29 | 1 |
| 690 | 1 | 183 | NA |
| 758 | 1 | 924 | 1 |
| 761 | 1 | 299 | 1 |
| 762 | 1 | 780 | NA |
| 775 | 1 | 110 | 1 |
| 835 | 0.5 | 2 | 1 |
| 843 | 0.99 | 104 | 1 |

- All-data CWE classifier
  - 97.2% accuracy
  - AUROC 1
  - **56** per-CWE accuracies (see left)
- All-data CERT rule classifier
  - **44** per-rule accuracies
  - 95% at least 95% accuracy, with lowest accuracy of 83%
- Results from CWE and CERT rules classifiers better than expected – currently investigating cause.
  - May be artifact of test file metadata
  - Expect reduced performance against native files

**Carnegie Mellon University**
Software Engineering Institute

**Research Review:**
**Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C**
© 2017 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

12

# Summary and Future

FY17 Line "Rapid Classifiers" built on the FY16 LENS "Prioritizing vulnerabilities".

- Developed widely useful general method to use test suites across taxonomies
- Developed large archive of "pre-audited" alerts
    - Overcame major challenge to classifier development
    - For CWEs and CERT rules
- Developed code infrastructure (extensible!)
- In-progress:
    - Classifier development and testing in process
    - Continue to gather data
    - Enhanced SCALe audit tool for collaborator testing: distribute to collaborators soon
- FY18-19 plan: architecture for rapid deployment of classifiers in varied systems
- Goal: optimal automation of static alert auditing (and other code analysis and repair)

Publications:
- New mappings (CWE/CERT rule): MITRE and CERT websites
- IEEE SecDev 2017 "Hands-on Tutorial: Alert Auditing with Lexicon & Rules"
- 2 SEI blogposts on classifier development
- Research paper in progress

# Contact Information

**Presenter / Point(s) of Contact**

Lori Flynn (Principal Investigator)

Software Security Researcher


Email:  lflynn@cert.org


Telephone:  +1 412.268.7886

**Contributors**

SEI Staff

William Snavely

David Svoboda

Zach Kurtz


SEI Student Interns

Lucas Bengtson (CMU)

Charisse Haruta (CMU)

Baptiste Vauthey (CMU)

Michael Spece (Pitt)

Christine Baek (CMU)

Carnegie Mellon University
Software Engineering Institute

Research Review:
**Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C**
© 2017 Carnegie Mellon University

[DISTRIBUTION STATEMENT A]  Approved for public release and unlimited distribution.

**14**

**Carnegie Mellon University**
Software Engineering Institute

Research Review:
**Rapid Expansion of Classification Models to Prioritize Static Analysis Alerts for C**
© 2017 Carnegie Mellon University

[DISTRIBUTION STATEMENT A]  Approved for public release and unlimited distribution.

**15**