# Predicting Security Attacks in FOSS

Why you want it and one way to do it

**C.E. Budde**     R. Paramitha     F. Massacci
Università di Trento (IT)  &  Vrije Universiteit (NL)

Vuln4Cast 2023 FIRST Technical Colloquium

## Talk overview

# Introduction

## Why You Should Update All Your Software

Updates may sometimes be painful, but they're necessary to keep your devices and data secure on a dangerous internet.

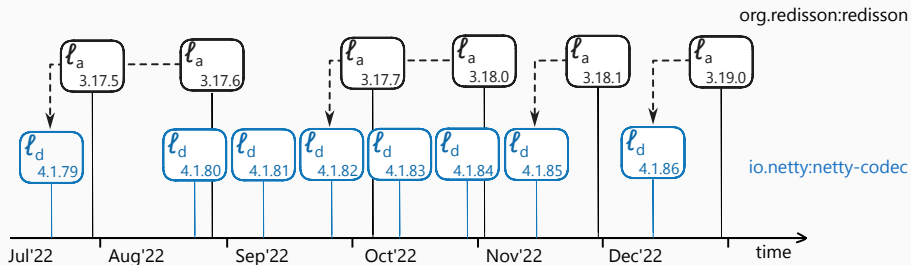BY **CHRIS HOFFMAN**    PUBLISHED AUG 28, 2020

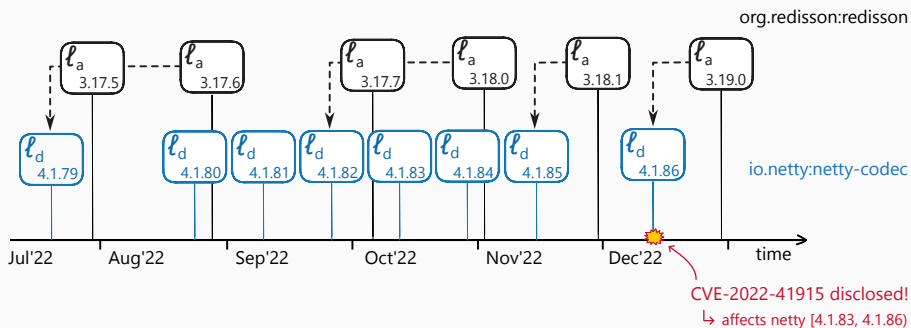**Quick Links**

Security Updates 101

What's the Risk, Really?

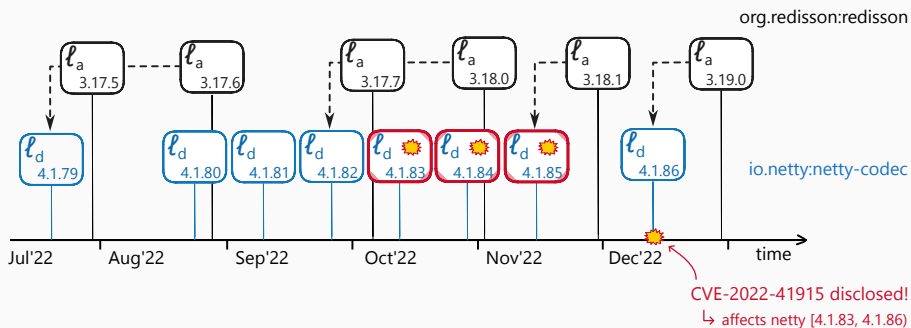© loonylabs

org.redisson:redisson

$\ell_a$ 3.17.5
$\ell_a$ 3.17.6
$\ell_a$ 3.17.7
$\ell_a$ 3.18.0
$\ell_a$ 3.18.1
$\ell_a$ 3.19.0

$\ell_d$ 4.1.79
$\ell_d$ 4.1.80
$\ell_d$ 4.1.81
$\ell_d$ 4.1.82
$\ell_d$ 4.1.83
$\ell_d$ 4.1.84
$\ell_d$ 4.1.85
$\ell_d$ 4.1.86

io.netty:netty-codec

Jul'22    Aug'22    Sep'22    Oct'22    Nov'22    Dec'22    time

CVE-2022-41915 disclosed!
↳ affects netty [4.1.83, 4.1.86)

# Hindsight!

Developer perspective in time:



org.redisson:redisson

io.netty:netty-codec

Developer perspective in time:

Developer perspective in time:

Developer perspective in time:

Developer perspective in time:

Developer perspective in time:

Developer perspective in time:

Developer perspective in time:



org.redisson:redisson

$\ell_a$ 3.17.5 — $\ell_a$ 3.17.6       $\ell_a$ 3.17.7 — $\ell_a$ 3.18.0       $\ell_a$ 3.18.1

$\ell_d$ 4.1.79   $\ell_d$ 4.1.80  $\ell_d$ 4.1.81  $\ell_d$ 4.1.82  $\ell_d$ 4.1.83  $\ell_d$ 4.1.84  $\ell_d$ 4.1.85   $\ell_d$ 4.1.86

io.netty:netty-codec

Jul'22   Aug'22   Sep'22   Oct'22   Nov'22   Dec'22   time

CVE-2022-41915 disclosed!
↳ affects netty [4.1.83, 4.1.86)

# Developer perspective in time:

Developer perspective in time:



Is there a best time to update?

**Q1**   How does time affect the $\Pr(\text{vuln.})$?

**Q2**   Which other factors affect $\Pr(\text{vuln.})$?

**Q1** How does **time** affect the $\Pr(\text{vuln.})$?
  ▷ best time to update?

**Q2** Which other factors affect $\Pr(\text{vuln.})$?

**Q1**  How does **time** affect the $\Pr(\text{vuln.})$?
   ▷ best time to update?

**Q2**  Which other factors affect $\Pr(\text{vuln.})$?
   ▷ measurable software metrics

1. Unpublished/Undetected vulnerabilities:
   - we study publication of CVEs;

# Outta scopes

1. **Unpublished/Undetected vulnerabilities:**
   - we study publication of CVEs;
   - keep it high-level, no code analysis.

1. Unpublished/Undetected vulnerabilities:
   - we study publication of CVEs;
   - keep it high-level, no code analysis.

2. Probability of *exploitation*:
   - we study publication of CVEs;

1. Unpublished/Undetected vulnerabilities:
   - we study publication of CVEs;
   - keep it high-level, no code analysis.

2. Probability of *exploitation*:
   - we study publication of CVEs;
   - … but check the work of the EPSS!

| Work | Goal | | Data | | | | Method | | | Approach | | | Projects/Libs. | | Purport |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Disc. | Pred. | CVEs | Code | VCS | Dep. | Corr. | Clas. | T.Ser. | AH | SA | ML | Language | # | |
| [4] | ✓ | | | ✓ | | | | ✓ | | | | ✓ | C | 3 | Find vulnerabilities regardless of existent logs such as CVEs (although CWEs may be used). This includes formal methods and static/dynamic code analysis. |
| [2] | ✓ | | | | ✓ | | ✓ | ✓ | | | | ✓ | PHP | 3 | |
| [16] | ✓ | | ✓ | | | | ✓ | ✓ | | | ✓ | ✓ | Java | 4 | |
| [5] | ✓ | | ✓ | ✓ | | | | ✓ | | ✓ | | | C/C++, PHP, Java, JS, SQL | 10 | |
| [11] | ✓ | | ✓ | | ✓ | | | ✓ | | ✓ | | | C | 3 | Detect known vulnerabilities (and their correlation to developer activity metrics) from VCS only—e.g. commit churn, peer comments, etc. |
| [13] | ✓ | | ✓ | | ✓ | | ✓ | | | ✓ | | | C | 1 | |
| [15] | ✓ | | ✓ | | ✓ | | ✓ | | | ✓ | ✓ | | C, ASM | 3 | |
| [14] | ✓ | | ✓ | | ✓ | | ✓ | | | ✓ | ✓ | | C, ASM | 1 | |
| [6] | ✓ | | ✓ | ✓ | | | | ✓ | | | | ✓ | C/C++ | 3 | Detect known vulnerabilities (and their correlation to code metrics) from code only—e.g. number of classes, code cloning, cyclomatic complexity, etc. |
| [8] | ✓ | | ✓ | ✓ | | | | ✓ | | | | ✓ | Java | 7 | |
| [23] | ✓ | | ✓ | ✓ | | | ✓ | | | | ✓ | ✓ | Java | 4 | |
| [24] | ✓ | | ✓ | ✓ | | | ✓ | | | | | ✓ | Java | 3 | |
| [25] | ✓ | | ✓ | ✓ | | | ✓ | | | | | ✓ | Java | 5 | |
| [21] | ✓ | | ✓ | ✓ | | | | ✓ | | ✓ | | | C | 7 | |
| [1] | ✓ | | ✓ | ✓ | ✓ | | ✓ | ✓ | | ✓ | | ✓ | C/C++ | >150k | Detect known vulnerabilities (and their corr. to code and developer activity metrics) from both code and VCS, but without considering the effect of dependencies in their propagation. |
| [9] | ✓ | | ✓ | ✓ | ✓ | | ✓ | ✓ | | ✓ | | | C/C++ | 8 | |
| [3] | ✓ | | ✓ | ✓ | ✓ | | ✓ | | | | ✓ | | C/C++ | 5 | |
| [7] | ✓ | | ✓ | ✓ | ✓ | | ✓ | ✓ | | | ✓ | ✓ | C/C++, Java | 1 | |
| [22] | ✓ | | ✓ | ✓ | ✓ | | ✓ | | | | ✓ | ✓ | C/C++ | 2 | |
| [18] | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | | ✓ | | | Java | 500 | Detect known vulnerabilities using code or VCS, via dependency-aware models that can find the offending code to help correcting it (own vs. third-party libraries). |
| [12] | ✓ | | ✓ | ✓ | | ✓ | | ✓ | | | | ✓ | Java | >300k | |
| [19] | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | ✓ | | Java, Ruby, Python | 450 | |
| [17] | ✓ | | ✓ | ✓ | | ✓ | | ✓ | | ✓ | | | Java | 200 | |
| [26] | | ✓ | ✓ | | | | | | ✓ | | ✓ | ✓ | Agnostic | 9 | Time regression to predict vulnerabilities from NVD logs, but the models lack data from the security domain. |
| [10] | | ✓ | ✓ | | | | | | ✓ | | ✓ | ✓ | Agnostic | 25 | |
| [20] | | ✓ | ✓ | | | | | | ✓ | | ✓ | | Agnostic | 5 | |

| Work | Goal | | Data | | | | Method | | | Approach | | | Projects/Libs. | | Purport |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Disc. | Pred. | CVEs | Code | VCS | Dep. | Corr. | Clas. | T.Ser. | AH | SA | ML | Language | # | |
| [4] | ✓ | | | ✓ | | | | ✓ | | | | ✓ | C | 3 | Find vulnerabilities regardless of existent logs such as CVEs (although CWEs may be used). This includes formal methods and static/dynamic code analysis. |
| [2] | ✓ | | | | ✓ | | | ✓ | | | | ✓ | PHP | 3 | |
| [16] | ✓ | | | ✓ | | | | ✓ | | | ✓ | ✓ | Java | 4 | |
| [5] | ✓ | | | ✓ | ✓ | | | ✓ | | ✓ | | | C/C++, PHP, Java, JS, SQL | 10 | |
| [11] | ✓ | | ✓ | ✓ | | | | ✓ | | ✓ | | | C | 3 | Detect known vulnerabilities (and their correlation to developer activity metrics) from VCS only—e.g. commit churn, peer comments, etc. |
| [13] | ✓ | | ✓ | | ✓ | | | ✓ | | ✓ | | | C | 1 | |
| [15] | ✓ | | ✓ | | ✓ | | | ✓ | | ✓ | ✓ | | C, ASM | 3 | |
| [14] | ✓ | | ✓ | | ✓ | | | ✓ | | ✓ | ✓ | | C, ASM | 1 | |
| [6] | ✓ | ✓ | ✓ | ✓ | | | | ✓ | | | | ✓ | C/C++ | 3 | Detect known vulnerabilities (and their correlation to code metrics) from code only—e.g. number of classes, code cloning, cyclomatic complexity, etc. |
| [8] | ✓ | ✓ | ✓ | ✓ | | | | ✓ | | | | ✓ | Java | 7 | |
| [23] | ✓ | ✓ | ✓ | ✓ | | | ✓ | | | | ✓ | ✓ | Java | 4 | |
| [24] | ✓ | ✓ | ✓ | ✓ | | | | ✓ | | | ✓ | | Java | 3 | |
| [25] | ✓ | ✓ | ✓ | ✓ | | | | ✓ | | | ✓ | | Java | 5 | |
| [21] | ✓ | ✓ | ✓ | ✓ | | | | | ✓ | ✓ | | | C | 7 | |
| [1] | ✓ | | ✓ | ✓ | ✓ | | ✓ | | | ✓ | | | C/C++ | >150k | Detect known vulnerabilities (and their corr. to code and developer activity metrics) from both code and VCS, but without considering the effect of dependencies in their propagation. |
| [9] | ✓ | | ✓ | ✓ | ✓ | | ✓ | | | ✓ | | | C/C++ | 8 | |
| [3] | ✓ | | ✓ | ✓ | ✓ | | | ✓ | | | | | C/C++ | 5 | |
| [7] | ✓ | | ✓ | ✓ | ✓ | | | ✓ | | | | ✓ | C/C++, Java | 1 | |
| [22] | ✓ | | ✓ | ✓ | ✓ | | | ✓ | | | | ✓ | C/C++ | 2 | |
| [18] | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | | ✓ | | | Java | 500 | Detect known vulnerabilities using code or VCS, via dependency-aware models that can find the offending code to help correcting it (own vs. third-party libraries). |
| [12] | ✓ | | ✓ | ✓ | | ✓ | | ✓ | | | | | Java | >300k | |
| [19] | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | | | | ✓ | Java, Ruby, Python | 450 | |
| [17] | ✓ | | ✓ | ✓ | | ✓ | | ✓ | | | | | Java | 200 | |
| [26] | | ✓ | | ✓ | | | | | ✓ | | ✓ | ✓ | Agnostic | 9 | Time regression to predict vulnerabilities from NVD logs, but the models lack data from the security domain. |
| [10] | | ✓ | | ✓ | | | | | ✓ | | ✓ | ✓ | Agnostic | 25 | |
| [20] | | ✓ | | ✓ | | | | | ✓ | | | ✓ | Agnostic | 5 | |

Annotations: *Most works try to discover current vulnerabilities, not predict future ones* · *Most works disregard the code dependency tree* · *Most works do not consider time in their analyses* · *Most works predict* · *Disregarded security data*

**Q2**   Pr(vuln.) as function of software metrics

**Q1**   Pr(vuln.) as function of time

**Q2**   $\text{Pr}(\text{vuln.})$ as function of software metrics

    ▶  ML & statistical analysis to correlate SE metrics to existent vulnerabilities

**Q1**   $\text{Pr}(\text{vuln.})$ as function of time

**Q2**   Pr(vuln.) as function of software metrics

- ▶ ML & statistical analysis to correlate SE metrics to existent vulnerabilities
- ▶ human-in-the-loop metrics, including VCS (#commits, seniority…)

**Q1**   Pr(vuln.) as function of time

**Q2**　Pr(vuln.) as function of software metrics

- ► ML & statistical analysis to correlate SE metrics to existent vulnerabilities
- ► human-in-the-loop metrics, including VCS (#commits, seniority…)
- ► (a few) considerations of own and $3^{rd}$ party dependencies

**Q1**　Pr(vuln.) as function of time

**Q2**  Pr(vuln.) as function of software metrics

- ▸ ML & statistical analysis to correlate SE metrics to existent vulnerabilities
- ▸ human-in-the-loop metrics, including VCS (#commits, seniority…)
- ▸ (a few) considerations of own and $3^{rd}$ party dependencies

**Q1**  Pr(vuln.) as function of time

- ▸ time-regression models on CVE publications ($\approx$ FinTech)

- Studies typically try to *detect*, not *foretell* vulnerabilities.

- Studies typically try to *detect*, not *foretell* vulnerabilities.

- The dependency tree is seldom analysed (own code only).

- Studies typically try to *detect*, not *foretell* vulnerabilities.

- The dependency tree is seldom analysed (own code only).

- The rare-event nature of vulnerabilities is disregarded.

- Studies typically try to *detect*, not *foretell* vulnerabilities.

- The dependency tree is seldom analysed (own code only).

- The rare-event nature of vulnerabilities is disregarded.

**We propose white-box model(s) to fill these gaps**

## Time Dependency Trees



$\ell_{0.8}$ → $\ell_{0.9}$ → $\ell_{1.0}$ → $\ell_{1.1}$

$x_{3.0}$ → $x_{3.3}$  $y_{5.0.0}$ → $y_{5.0.1}$ → $y_{5.8.3}$

$z_{2.0}$ → $z_{2.1}$ → $z_{2.2}$

CVE root-lib PDFs

Dependency Trees in time

$D(\ell_{a_1})$:

$\ell_{a_1}$

$\ell_{d_2}$  $\ell_{c_1}$

$\ell_{d_1}$

Dependency Trees in time

Dependency Trees in time

Dependency Trees in time



$\{D(\ell_{a_i})\}_{i=1}^3$:

Dependency Trees in time

Time Dependency Tree

# Time Dependency Trees

Dependency Trees in time

Time Dependency Tree



$\{D(\ell_{a_i})\}_{i=1}^3$:

$D_T(\ell_a)$:

c-chain

dependency

Main library $(\ell_a)$

Dependency Trees in time

Time Dependency Tree

$\{D(\ell_{a_i})\}_{i=1}^3$:

$D_T(\ell_a)$:

Main library $(\ell_a)$

Time span $(T)$

c-chain

dependency

Dependency Trees in time

Time Dependency Tree



$\{D(\ell_{a_i})\}_{i=1}^3$:

$D_T(\ell_a)$:

c-chain
dependency

Main library $(\ell_a)$

Time span $(T)$

$D_t(\ell_a) = D(\ell_{a_1})$
for any time point $t \in T$
after the release of $\ell_{a_1}$ and
before the release of $\ell_{a_2}$

- Minimal graph representation (no lib-version repetition)

- Minimal graph representation (no lib-version repetition)
- Canonical for library $\ell$ and time span $T$

- Minimal graph representation (no lib-version repetition)
- Canonical for library $\ell$ and time span $T$
- Natural lifting of dependency trees to time

## Theoretical

- Minimal graph representation (no lib-version repetition)
- Canonical for library $\ell$ and time span $T$
- Natural lifting of dependency trees to time

## Theoretical

- Minimal graph representation (no lib-version repetition)
- Canonical for library $\ell$ and time span $T$
- Natural lifting of dependency trees to time

## Practical

- **Time-indexing** $D_t(\ell)$ yields the dep. tree at time $t \in T$

## Theoretical

- Minimal graph representation (no lib-version repetition)
- Canonical for library $\ell$ and time span $T$
- Natural lifting of dependency trees to time

## Practical

- Time-indexing $D_t(\ell)$ yields the dep. tree at time $t \in T$
- Library-slicing $D_T(\ell)\big|_d$ yields *all instances* of dependency $d$ during time $T$

## Theoretical

- Minimal graph representation (no lib-version repetition)
- Canonical for library $\ell$ and time span $T$
- Natural lifting of dependency trees to time

## Practical

- **Time-indexing** $D_t(\ell)$ yields the dep. tree at time $t \in T$
- **Library-slicing** $D_T(\ell)\big|_d$ yields *all instances* of dependency $d$ during time $T$
- Reachability analysis can spot single-points-of-failure

My personal project uses $\ell_{1.0}$

My personal project uses $\ell_{1.0}$

My personal project uses $\ell_{1.0}$



Should I downgrade to $\ell_{0.9}$ or upgrade to $\ell_{1.1}$?

## Theoretical

- Minimal graph representation (no lib-version repetition)
- Canonical for library $\ell$ and time span $T$
- Natural lifting of dependency trees to time

## Practical

- **Time-indexing** $D_t(\ell)$ yields the dep. tree at time $t \in T$
- **Library-slicing** $D_T(\ell)\big|_d$ yields *all instances* of dependency $d$ during time $T$
- Reachability analysis can spot single-points-of-failure

## Theoretical

- Minimal graph representation (no lib-version repetition)
- Canonical for library $\ell$ and time span $T$
- Natural lifting of dependency trees to time

## Practical

- **Time-indexing** $D_t(\ell)$ yields the dep. tree at time $t \in T$
- **Library-slicing** $D_T(\ell)\big|_d$ yields *all instances* of dependency $d$ during time $T$
- Reachability analysis can spot single-points-of-failure
- Can measure health/risk of development environment

## Time Dependency Trees



## CVE root-lib PDFs

org.redisson:redisson

io.netty:netty-codec

CVE-2022-41915 disclosed!
↳ affects netty [4.1.83, 4.1.86]

# Publication of CVE since time of code release

# Publication of CVE since time of code release

org.redisson:redisson
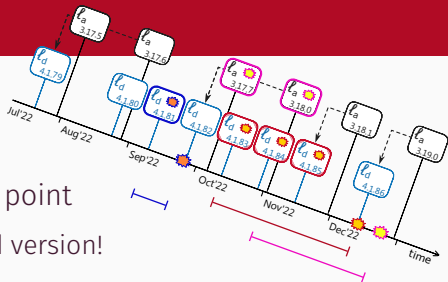
io.netty:netty-codec

CVE-2022-41915 disclosed!
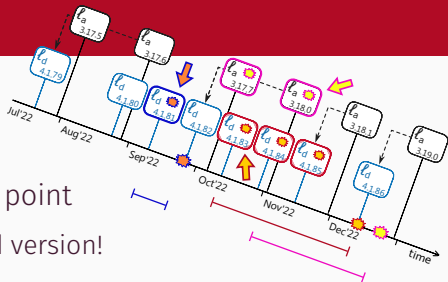↳ affects netty [4.1.83, 4.1.86)

- Count each CVE as one data point
  - must choose one affected version!

- Count each CVE as one data point
  - must choose one affected version!

- Count each CVE as one data point
  - must **choose one** affected version!

# Rules of the game

- ▶ Count each CVE as one data point
  - · must choose one affected version!

- ▶ Discriminate per development environment
  - · e.g. Java and C/C++ have different vuln. (and times!)

- ▶ Count each CVE as one data point
  - must choose one affected version!

- ▶ Discriminate per development environment
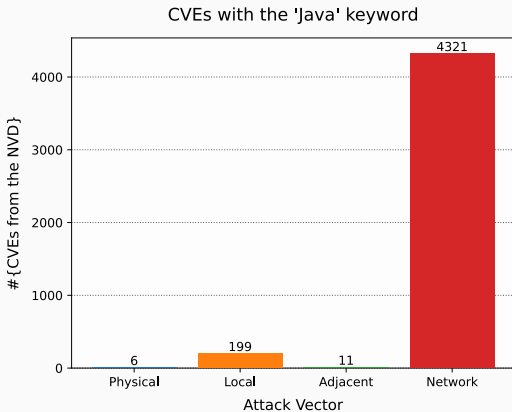  - e.g. Java and C/C++ have different vuln. (and times!)

# Rules of the game

- ▶ Count each CVE as one data point
  - · must choose one affected version!

- ▶ Discriminate per development environment
  - · e.g. Java and C/C++ have different vuln. (and times!)

- ▶ Discriminate per library type
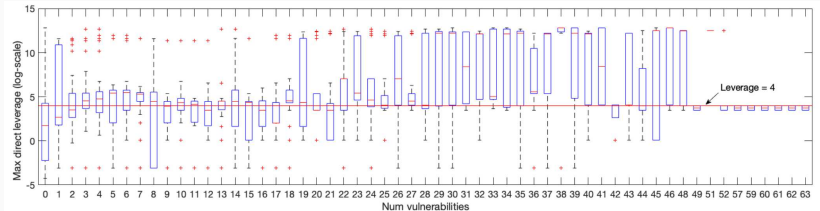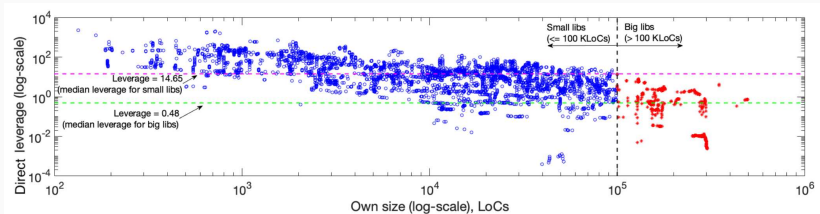  - · consider security-relevant code metrics

- ▶ Count each CVE as one data point
  - · must choose one affected version!

- ▶ Discriminate per development environment
  - · e.g. Java and C/C++ have different vuln. (and times!)

- ▶ Discriminate per library type
  - · consider security-relevant code metrics

## Used in remote networks



CVEs with the 'Java' keyword

## (Own) Code size

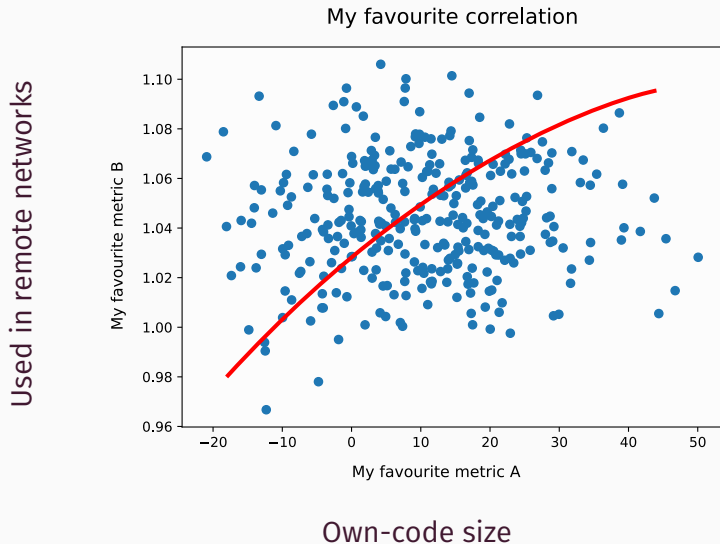| Work | Goal | | Data | | | | Method | | | Approach | | | Projects/Libs. | | Purport |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Disc. | Pred. | CVEs | Code | VCS | Dep. | Corr. | Clas. | T.Ser. | AH | SA | ML | Language | # | |
| [4] | ✓ | | | ✓ | | | | ✓ | | | | ✓ | C | 3 | Find vulnerabilities regardless of existent logs such as CVEs (although CWEs may be used). This includes formal methods and static/dynamic code analysis. |
| [2] | ✓ | | | | ✓ | | ✓ | ✓ | | | | ✓ | PHP | 3 | |
| [16] | ✓ | | | ✓ | ✓ | | ✓ | ✓ | | | ✓ | ✓ | Java | 4 | |
| [5] | ✓ | | | ✓ | ✓ | | | ✓ | | ✓ | | | C/C++, PHP, Java, JS, SQL | 10 | |
| [11] | ✓ | | ✓ | | ✓ | | | ✓ | | | ✓ | | C | 3 | Detect known vulnerabilities (and their correlation to developer activity metrics) from VCS only—e.g. commit churn, peer comments, etc. |
| [13] | ✓ | | ✓ | | ✓ | | ✓ | | | | ✓ | | C | 1 | |
| [15] | ✓ | | ✓ | | ✓ | | ✓ | | | ✓ | ✓ | | C, ASM | 3 | |
| [14] | ✓ | | ✓ | | ✓ | | ✓ | | | ✓ | ✓ | | C, ASM | 1 | |
| [6] | ✓ | | ✓ | ✓ | | | | ✓ | | | | ✓ | C/C++ | 3 | Detect known vulnerabilities (and their correlation to code metrics) from code only—e.g. number of classes, code cloning, cyclomatic complexity, etc. |
| [8] | ✓ | | ✓ | ✓ | | | | ✓ | | | | | Java | 7 | |
| [23] | ✓ | | ✓ | ✓ | | | | | | | ✓ | ✓ | Java | 4 | |
| [24] | ✓ | | ✓ | ✓ | | | ✓ | | | | | ✓ | Java | 3 | |
| [25] | ✓ | | ✓ | ✓ | | | ✓ | | | | | ✓ | Java | 5 | |
| [21] | ✓ | | ✓ | ✓ | | | | ✓ | | ✓ | | | C | 7 | |
| [1] | ✓ | | ✓ | ✓ | ✓ | | | ✓ | | | | ✓ | C/C++ | >150k | Detect known vulnerabilities (and their corr. to code and developer activity metrics) from both code and VCS, but without considering the effect of dependencies in their propagation. |
| [9] | ✓ | | ✓ | ✓ | ✓ | | | ✓ | | ✓ | | | C/C++ | 8 | |
| [3] | ✓ | | ✓ | ✓ | ✓ | | ✓ | | | | ✓ | | C/C++ | 8 | |
| [7] | ✓ | | ✓ | ✓ | ✓ | | ✓ | ✓ | | | ✓ | ✓ | C/C++, Java | 1 | |
| [22] | ✓ | | ✓ | ✓ | ✓ | | | ✓ | | | ✓ | ✓ | C/C++ | 2 | |
| [18] | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | | ✓ | | | Java | 500 | Detect known vulnerabilities using code or VCS, via dependency-aware models that can find the offending code to help correcting it (own vs. third-party libraries). |
| [12] | ✓ | | ✓ | ✓ | ✓ | ✓ | | ✓ | | | | ✓ | Java | >300k | |
| [19] | ✓ | | ✓ | ✓ | ✓ | ✓ | | | | | ✓ | | Java, Ruby, Python | 450 | |
| [17] | ✓ | | ✓ | | | ✓ | | | | ✓ | | | Java | 200 | |
| [26] | | ✓ | ✓ | | | | | | ✓ | | ✓ | ✓ | Agnostic | 9 | Time regression to predict vulnerabilities from NVD logs, but the models lack data from the security domain. |
| [10] | | ✓ | ✓ | | | | | | ✓ | | | ✓ | Agnostic | 25 | |
| [20] | | ✓ | ✓ | | | | | | ✓ | | | ✓ | Agnostic | 5 | |

# Security-relevant code metrics



Used in remote networks

Own-code size

Used in remote networks



My favourite correlation

Own-code size

Used in remote networks



My favourite correlation

Own-code size

My favourite correlation

Used in remote networks

Own-code size

Used in remote networks



My favourite correlation

Own-code size

My favourite correlation

Used in remote networks

Own-code size

Used in remote networks



Own-code size

My favourite correlation

My favourite correlation

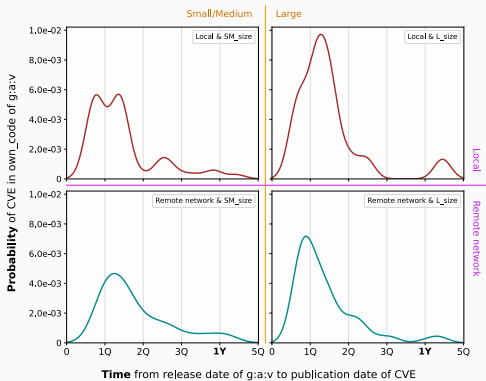# On overfitting and rare events

- ▶ Count each CVE as one data point

- ▶ Discriminate per development environment

- ▶ Discriminate per library type

# On overfitting and rare events

- Count each CVE as one data point

- Discriminate per development environment

- Discriminate per library type

- Clusterisation mustn't be too thin
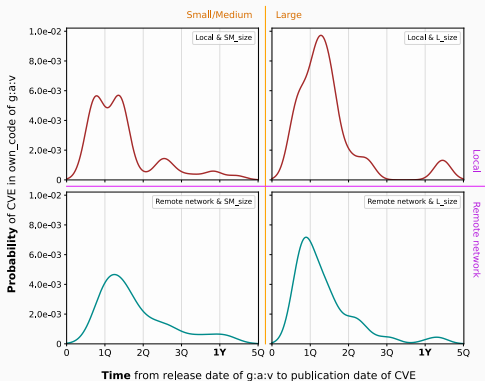  - few divisions per metric-dimension
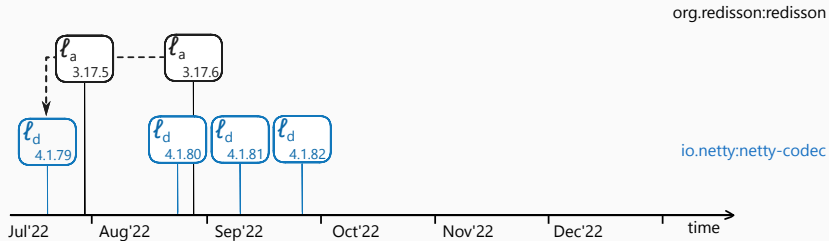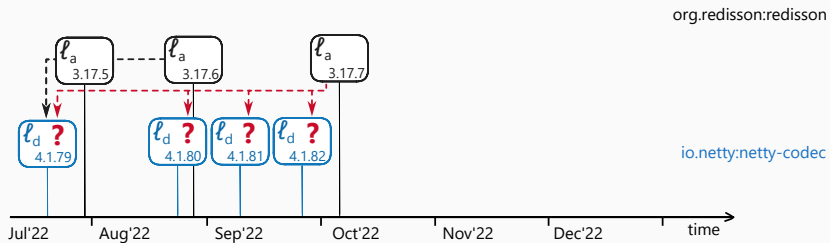  - few metric-dimensions

# Enough!

Gimme results

**Q1** Pr(vuln.) as function of time

**Q2** Pr(vuln.) as function of software metrics

$A \xrightarrow{t} B$ means that we change from dependency $\ell_A$ to $\ell_B$ in $t$ time units counting from $t_0$ ("today").

▷ $\ell_A$ was released on $t_A < t_0$, $\ell_B$ on $t_B < t_0$, $t_A \bowtie t_B$

$A \xrightarrow{t} B$ means that we change from dependency $\ell_A$ to $\ell_B$ in $t$ time units counting from $t_0$ ("today").



▷ $\ell_A$ was released on $t_A < t_0$, $\ell_B$ on $t_B < t_0$, $t_A \bowtie t_B$

**Q:** $\Pr_{A,B}(t) =$ probability of vuln. of $A \xrightarrow{t} B$ as a function of $t$

# Survival analysis on library update

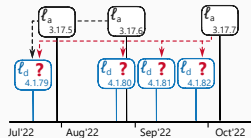$A \xrightarrow{t} B$ means that we change from dependency $\ell_A$ to $\ell_B$ in $t$ time units counting from $t_0$ ("today").



▷ $\ell_A$ was released on $t_A < t_0$, $\ell_B$ on $t_B < t_0$, $t_A \bowtie t_B$

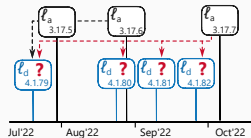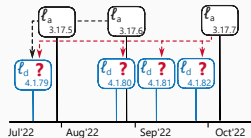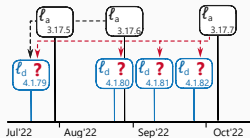**Q:** $\Pr_{A,B}(t) =$ probability of vuln. of $A \xrightarrow{t} B$ as a function of $t$

**A:** $\Pr_{A,B}(t) = 1 - \mathsf{SF}_A\big(t + \Delta t_A\big)\,\mathsf{CDF}_B\big(t + \Delta t_B\big)$   where $\Delta t_x \doteq |t_x - t_0|$

$A \xrightarrow{t} B$ means that we change from dependency $\ell_A$ to $\ell_B$ in $t$ time units counting from $t_0$ ("today").



▷ $\ell_A$ was released on $t_A < t_0$, $\ell_B$ on $t_B < t_0$, $t_A \bowtie t_B$

**Q:** $\mathsf{Pr}_{A,B}(t) =$ probability of vuln. of $A \xrightarrow{t} B$ as a function of $t$

**A:** $\mathsf{Pr}_{A,B}(t) = 1 - \mathsf{SF}_A\big(t + \Delta t_A\big)\,\mathsf{CDF}_B\big(t + \Delta t_B\big)$ where $\Delta t_x \doteq |t_x - t_0|$

<span style="color:red">vuln. in $\ell_A$ before change</span>    <span style="color:green">vuln. in $\ell_B$ after change</span>

$A \xrightarrow{t} B$ means that we change from dependency $\ell_A$ to $\ell_B$ in $t$ time units counting from $t_0$ ("today").

▷ $\ell_A$ was released on $t_A < t_0$, $\ell_B$ on $t_B < t_0$, $t_A \bowtie t_B$



**Q:** $\mathrm{Pr}_{A,B}(t) =$ probability of vuln. of $A \xrightarrow{t} B$ as a function of $t$

**A:** $\mathrm{Pr}_{A,B}(t) = 1 - \mathsf{SF}_A\big(t + \Delta t_A\big)\,\mathsf{CDF}_B\big(t + \Delta t_B\big)$   where $\Delta t_x \doteq |t_x - t_0|$

vuln. in $\ell_A$ before change    vuln. in $\ell_B$ after change
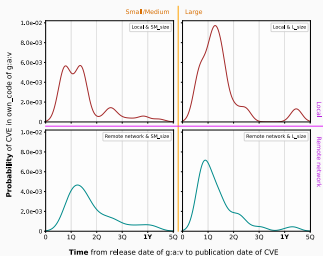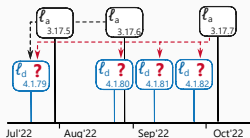
# Survival analysis on library update

$A \xrightarrow{t} B$ means that we change from dependency $\ell_A$ to $\ell_B$ in $t$ time units counting from $t_0$ ("today").
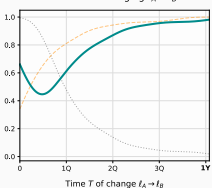


▷ $\ell_A$ was released on $t_A < t_0$, $\ell_B$ on $t_B < t_0$, $t_A \bowtie t_B$

**Q:** $\Pr_{A,B}(t) =$ probability of vuln. of $A \xrightarrow{t} B$ as a function of $t$
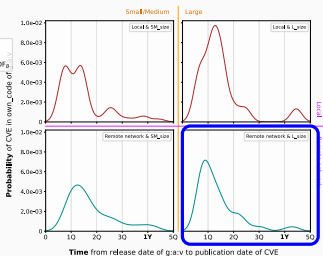
**A:** $\Pr_{A,B}(t) = 1 - \mathsf{SF}_A\big(t + \Delta t_A\big)\,\mathsf{CDF}_B\big(t + \Delta t_B\big)$  where $\Delta t_x \doteq |t_x - t_0|$



$t_A = 17$ days
$t_B = 85$ days



$t_A = 184$ days
$t_B = 21$ days

# Survival analysis on library update

$A \xrightarrow{t} B$ means that we change from dependency $\ell_A$ to $\ell_B$ in $t$ time units counting from $t_0$ ("today").



▷ $\ell_A$ was released on $t_A < t_0$, $\ell_B$ on $t_B < t_0$, $t_A \bowtie t_B$

**Q:** $\text{Pr}_{A,B}(t) = $ probability of vuln. of $A \xrightarrow{t} B$ as a function of $t$

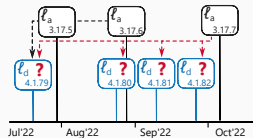**A:** $\text{Pr}_{A,B}(t) = 1 - \text{SF}_A\big(t + \Delta t_A\big)\,\text{CDF}_B\big(t + \Delta t_B\big)$   where $\Delta t_x \doteq |t_x - t_0|$



$t_A = 17$ days
$t_B = 85$ days



$t_A = 184$ days
$t_B = 21$ days

$A \xrightarrow{t} B$ means that we change from dependency $\ell_A$ to $\ell_B$ in $t$ time units counting from $t_0$ ("today").

▷ $\ell_A$ was released on $t_A < t_0$, $\ell_B$ on $t_B < t_0$, $t_A \bowtie t_B$



**Q:** $\mathrm{Pr}_{A,B}(t)$ = probability of vuln. of $A \xrightarrow{t} B$ as a function of $t$

**A:** $\mathrm{Pr}_{A,B}(t) = 1 - \mathsf{SF}_A\big(t + \Delta t_A\big)\,\mathsf{CDF}_B\big(t + \Delta t_B\big)$   where $\Delta t_x \doteq |t_x - t_0|$



$t_A = 17$ days
$t_B = 85$ days

$t_A = 184$ days
$t_B = 21$ days

# Survival analysis on library update

$A \xrightarrow{t} B$ means that we change from dependency $\ell_A$ to $\ell_B$ in $t$ time units counting from $t_0$ ("today").
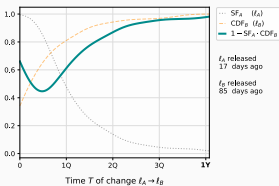
▷ $\ell_A$ was released on $t_A < t_0$, $\ell_B$ on $t_B < t_0$, $t_A \bowtie t_B$



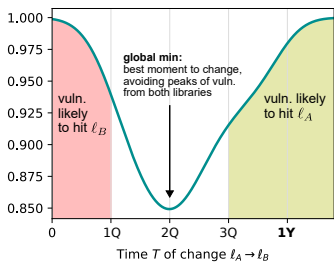**Q:** $\Pr_{A,B}(t) =$ probability of vuln. of $A \xrightarrow{t} B$ as a function of $t$

**A:** $\Pr_{A,B}(t) = 1 - \mathsf{SF}_A\big(t + \Delta t_A\big)\,\mathsf{CDF}_B\big(t + \Delta t_B\big)$   where $\Delta t_x \doteq |t_x - t_0|$
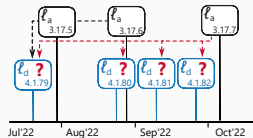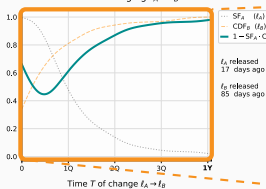


$t_A = 17$ days
$t_B = 85$ days

$t_A = 184$ days
$t_B = 21$ days

**Q:** $\Pr_{A,B}(t) =$ probability of vuln. in $\ell_A$ or $\ell_B$ before $t$

## Vulnerabilities from any dependency

**Q:** $\Pr_{A,B}(t) =$ probability of vuln. in $\ell_A$ or $\ell_B$ before $t$

**A:** $\Pr_{A,B}(t) = \Pr(\min(\ell_A, \ell_B) \leqslant t) = 1 - (1 - \Pr_A(t))(1 - \Pr_B(t))$

# Vulnerabilities from any dependency

**Q:** $\Pr_{A,B}(t) =$ probability of vuln. in $\ell_A$ or $\ell_B$ before $t$

**A:** $\Pr_{A,B}(t) = \Pr(\min(\ell_A, \ell_B) \leqslant t) = 1 - (1 - \Pr_A(t))(1 - \Pr_B(t))$

# Vulnerabilities from any dependency

**Q:** $\Pr_{A,B}(t)$ = probability of vuln. in $\ell_A$ or $\ell_B$ before $t$

**A:** $\Pr_{A,B}(t) = \Pr(\min(\ell_A, \ell_B) \leqslant t) = 1 - (1 - \Pr_A(t))(1 - \Pr_B(t))$

# Vulnerabilities from any dependency

**Q:** $\Pr_{A,B}(t)$ = probability of vuln. in $\ell_A$ or $\ell_B$ before $t$

**A:** $\Pr_{A,B}(t) = \Pr(\min(\ell_A, \ell_B) \leqslant t) = 1 - (1 - \Pr_A(t))(1 - \Pr_B(t))$

# Vulnerabilities from any dependency

**Q:** $\Pr_{A,B}(t)$ = probability of vuln. in $\ell_A$ or $\ell_B$ before $t$

**A:** $\Pr_{A,B}(t) = \Pr(\min(\ell_A, \ell_B) \leqslant t) = 1 - (1 - \Pr_A(t))(1 - \Pr_B(t))$



Nice for 2 dependencies...

I have 2000

**TDTs!**

$t_A = 28$ days
$t_B = 60$ days

$t_A = 123$ days
$t_B = 14$ days

## Time Dependency Trees

$\ell_{0.8}$ → $\ell_{0.9}$ → $\ell_{1.0}$ → $\ell_{1.1}$

$x_{3.0}$ → $x_{3.3}$ → $y_{5.0.0}$ → $y_{5.0.1}$ → $y_{5.8}$

$z_{2.0}$ → $z_{2.1}$ → $z_{2.2}$

## CVE root-lib PDFs

Pr(vuln. in Small/Medium lib)

Pr(vuln. in Large lib)

- ▶ Time Dependency Trees

- Time Dependency Trees
  - Aggregate dependency and code-evolution data

▶ Time Dependency Trees

- Aggregate dependency and code-evolution data
- Minimal representation with nice properties

- ▶ Time Dependency Trees
  - · Aggregate dependency and code-evolution data
  - · Minimal representation with nice properties
  - · Framework for large-scale project analysis

- ▶ Time Dependency Trees
  - · Aggregate dependency and code-evolution data
  - · Minimal representation with nice properties
  - · Framework for large-scale project analysis

- ▶ Probability of vulnerabilities as a function of time

- ▶ Time Dependency Trees
  - · Aggregate dependency and code-evolution data
  - · Minimal representation with nice properties
  - · Framework for large-scale project analysis

- ▶ Probability of vulnerabilities as a function of time
  - · Express time from library release to CVE publication

- ▶ Time Dependency Trees
  - · Aggregate dependency and code-evolution data
  - · Minimal representation with nice properties
  - · Framework for large-scale project analysis

- ▶ Probability of vulnerabilities as a function of time
  - · Express time from library release to CVE publication
  - · Discriminate per type of library (security-relevant props.)

# Some things done

- ▶ Time Dependency Trees
  - · Aggregate dependency and code-evolution data
  - · Minimal representation with nice properties
  - · Framework for large-scale project analysis

- ▶ Probability of vulnerabilities as a function of time
  - · Express time from library release to CVE publication
  - · Discriminate per type of library (security-relevant props.)
  - · Base information for probability forecasting

- Other metrics to clusterise libraries for PDF-fitting

- Other metrics to clusterise libraries for PDF-fitting

- Validate in other languages (all Java so far)

- Other metrics to clusterise libraries for PDF-fitting

- Validate in other languages (all Java so far)

- SPoF detection—across versions—in Java/Maven

- Other metrics to clusterise libraries for PDF-fitting

- Validate in other languages (all Java so far)

- SPoF detection—across versions—in Java/Maven

- c-chains polution by CVE

# Questions?

J. Akram and P. Luo.
**SQVDT: A scalable quantitative vulnerability detection technique for source code security assessment.**
*Software: Practice and Experience*, 51(2):294–318, 2021.

M. Alohaly and H. Takabi.
**When do changes induce software vulnerabilities?**
In *CIC*, pages 59–66. IEEE, 2017.

H. Alves, B. Fonseca, and N. Antunes.
**Software metrics and security vulnerabilities: Dataset and exploratory study.**
In *EDCC*, pages 37–44. IEEE, 2016.

Z. Bilgin, M. A. Ersoy, E. U. Soykan, E. Tomur, P. Çomak, and L. Karaçay.
**Vulnerability prediction from source code using machine learning.**
*IEEE Access*, 8:150672–150684, 2020.

A. Bosu, J. C. Carver, M. Hafiz, P. Hilley, and D. Janni.
**Identifying the characteristics of vulnerable code changes: An empirical study.**
In *FSE*, pages 257–268. ACM, 2014.

S. Chakraborty, R. Krishna, Y. Ding, and B. Ray.
**Deep learning based vulnerability detection: Are we there yet.**
*IEEE Transactions on Software Engineering*, 48(9):3280–3296, 2021.

I. Chowdhury and M. Zulkernine.
**Using complexity, coupling, and cohesion metrics as early indicators of vulnerabilities.**
*Journal of Systems Architecture*, 57(3):294–313, 2011.

S. Ganesh, T. Ohlsson, and F. Palma.
**Predicting security vulnerabilities using source code metrics.**
In *SweDS*, pages 1–7. IEEE, 2021.

S. Kim, S. Woo, H. Lee, and H. Oh.
**VUDDY: A scalable approach for vulnerable code clone discovery.**
In *SP*, pages 595–614. IEEE, 2017.

D. Last.
**Forecasting zero-day vulnerabilities.**
In *CISRC*, pages 1–4. ACM, 2016.

H. Li, H. Kwon, J. Kwon, and H. Lee.
**A scalable approach for vulnerability discovery based on security patches.**
In *ATIS*, volume 490 of *CCIS*, pages 109–122. Springer, 2014.

Q. Li, J. Song, D. Tan, H. Wang, and J. Liu.
**PDGraph: A large-scale empirical study on project dependency of security vulnerabilities.**
In *DSN*, pages 161–173. IEEE, 2021.

A. Meneely, H. Srinivasan, A. Musa, A. R. Tejeda, M. Mokary, and B. Spates.
**When a patch goes bad: Exploring the properties of vulnerability-contributing commits.**
In *ESEM*, pages 65–74. IEEE, 2013.

A. Meneely and L. Williams.
**Secure open source collaboration: An empirical study of Linus' law.**
In *CCS*, pages 453––462. ACM, 2009.

A. Meneely and L. Williams.
**Strengthening the empirical analysis of the relationship between Linus' law and software security.**
In *ESEM*. ACM, 2010.

Y. Pang, X. Xue, and A. S. Namin.
**Predicting vulnerable software components through N-gram analysis and statistical feature selection.**
In *ICMLA*, pages 543–548. IEEE, 2015.

I. Pashchenko, H. Plate, S. E. Ponta, A. Sabetta, and F. Massacci.
**Vulnerable open source dependencies: Counting those that matter.**
In *ESEM*. ACM, 2018.

I. Pashchenko, H. Plate, S. E. Ponta, A. Sabetta, and F. Massacci.
**Vuln4Real: A methodology for counting actually vulnerable dependencies.**
*IEEE Transactions on Software Engineering*, 48(5):1592–1609, 2022.

📄 G. A. A. Prana, A. Sharma, L. K. Shar, D. Foo, A. E. Santosa, A. Sharma, and D. Lo.
**Out of sight, out of mind? how vulnerable dependencies affect open-source projects.**
*Empirical Software Engineering*, 26(4), 2021.

📄 Y. Roumani, J. K. Nwankpa, and Y. F. Roumani.
**Time series modeling of vulnerabilities.**
*Computers & Security*, 51:32–40, 2015.

📄 N. Shahmehri, A. Mammar, E. Montes de Oca, D. Byers, A. Cavalli, S. Ardi, and W. Jimenez.
**An advanced approach for modeling and detecting software vulnerabilities.**
*Information and Software Technology*, 54(9):997–1013, 2012.

📄 Y. Shin, A. Meneely, L. Williams, and J. A. Osborne.
**Evaluating complexity, code churn, and developer activity metrics as indicators of software vulnerabilities.**
*IEEE Transactions on Software Engineering*, 37(6):772–787, 2011.

K. Z. Sultana, V. Anu, and T.-Y. Chong.
**Using software metrics for predicting vulnerable classes and methods in Java projects: A machine learning approach.**
*Journal of Software: Evolution and Process*, 33(3), 2021.

K. Z. Sultana, A. Deo, and B. J. Williams.
**Correlation analysis among Java nano-patterns and software vulnerabilities.**
In *HASE*, pages 69–76. IEEE, 2017.

K. Z. Sultana and B. J. Williams.
**Evaluating micro patterns and software metrics in vulnerability prediction.**
In *SoftwareMining*, pages 40–47. IEEE, 2017.

E. Yasasin, J. Prester, G. Wagner, and G. Schryen.
**Forecasting IT security vulnerabilities – an empirical analysis.**
*Computers & Security*, 88, 2020.

# Predicting Security Attacks in FOSS

Why you want it and one way to do it

**C.E. Budde**      R. Paramitha      F. Massacci
Università di Trento (IT)  &  Vrije Universiteit (NL)

Vuln4Cast 2023 FIRST Technical Colloquium