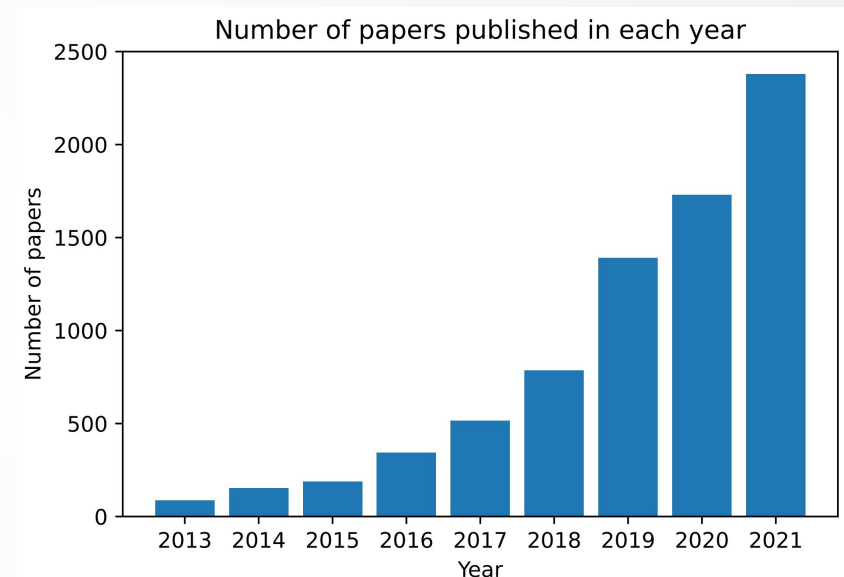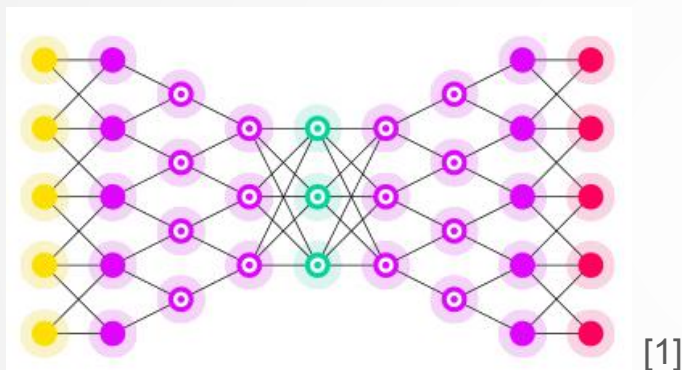**SIEMENS**
*Ingenuity for life*

# How Good is Openly Available Code Snippets Containing Software Vulnerabilities to Train Machine Learning Algorithms?

Kaan Oguzhan
Dr. rer. nat. Tiago Espinha Gasiba
Akram Louati

Siemens AG
Munich

Unrestricted

# Background

## Machine Learning



[1]



Number of papers published in each year
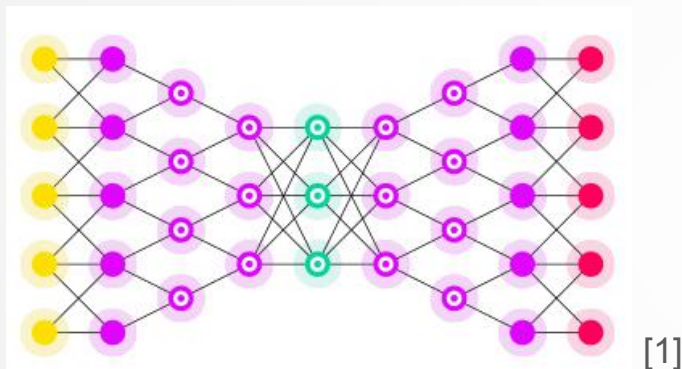
Appearance of both keywords "Cybersecurity" and "Machine Learning" in Academic Papers according to Scopus

# Background

## Machine Learning



[1]

**Model is as good as your data**

## Data



Quality & Quantity

# Motivation

➔ Goal: Developing software vulnerability detection in source code by means of ML Algorithms

➔ Training sounds straightforward, but

◆ *"Model is as good as your data"*

# Motivation

➔ Research questions:

- Where can we find code snippets to train ML models to detect software vulnerabilities?

- What is the quality of the code snippets which are openly available on the internet for training ML Models?
  - Can they be used to train ML models?

NOTE: in our work we use industry standard categories of software vulnerabilities

# Overview

➔ **Publicly available code snippets**
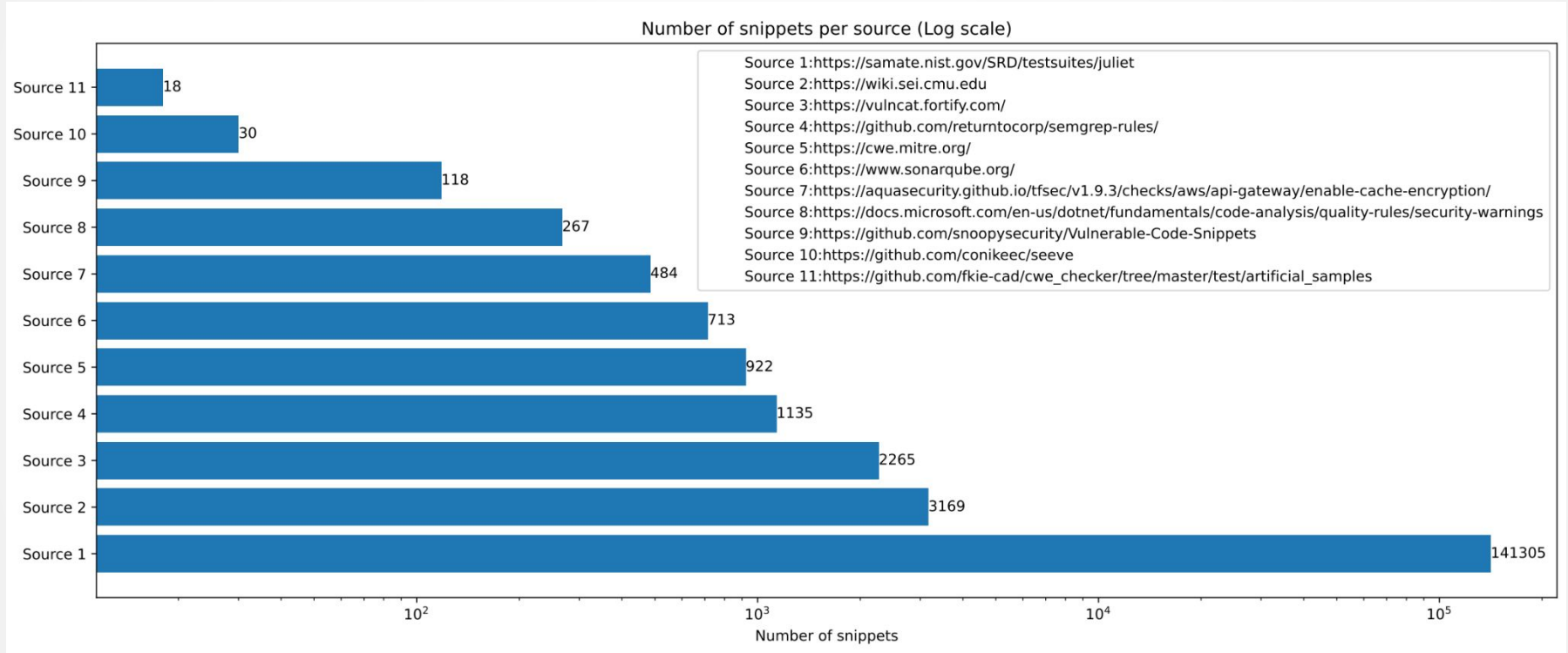
       Quality & Quantity

➔ Analysis measures
    i.   Categories
        ○   Programming Language
        ○   OWASP TOP 10
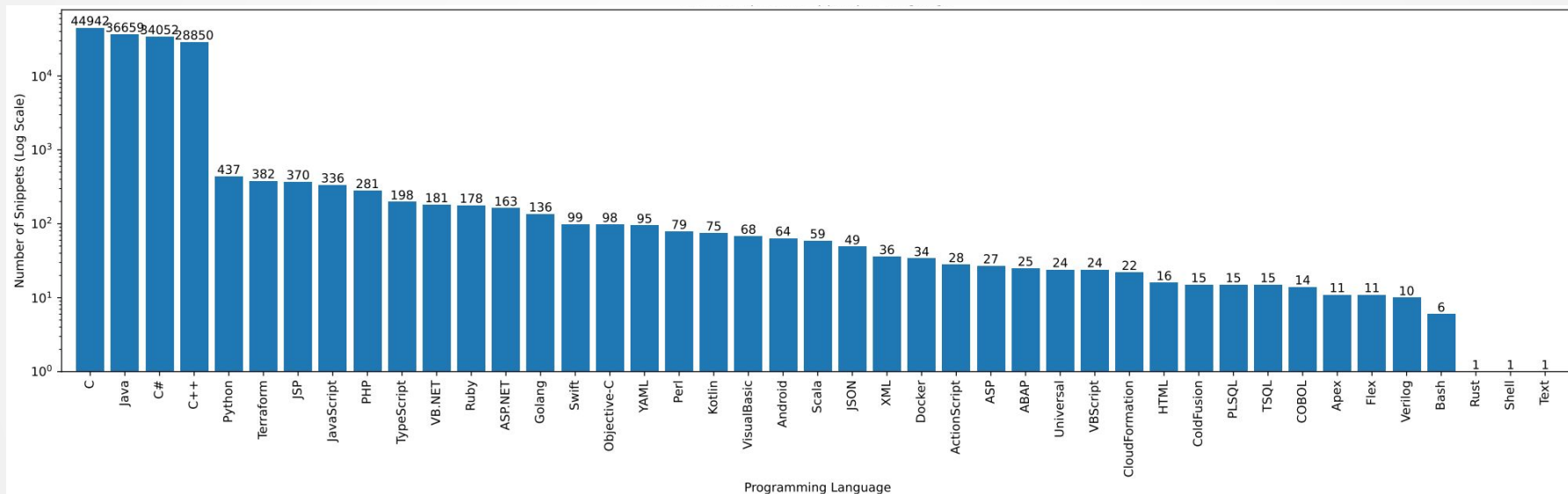        ○   PCI-DSS
        ○   CWE (Common Weakness Enumeration)
    ii.   Fitness for ML
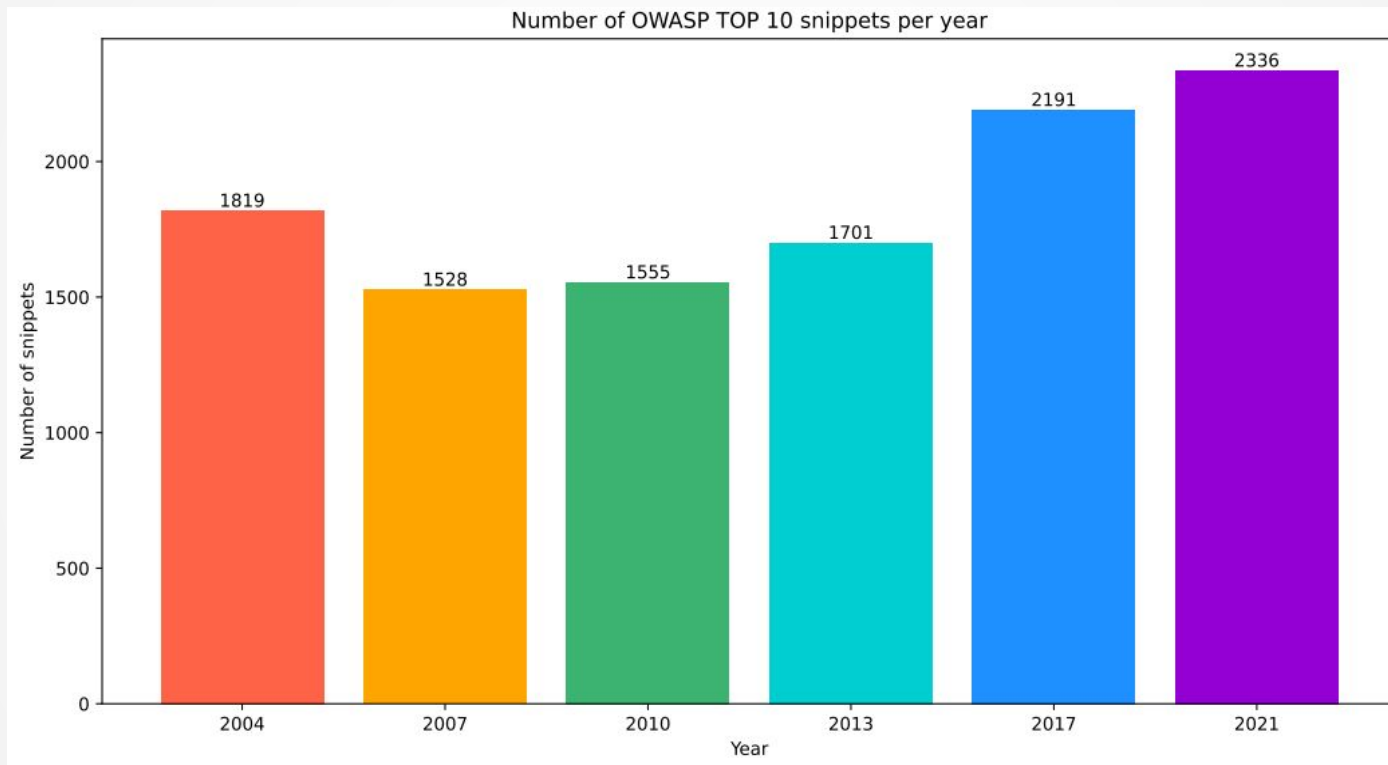➔ Conclusion

# Publicly Available Snippets per Source



Number of snippets per source (Log scale)
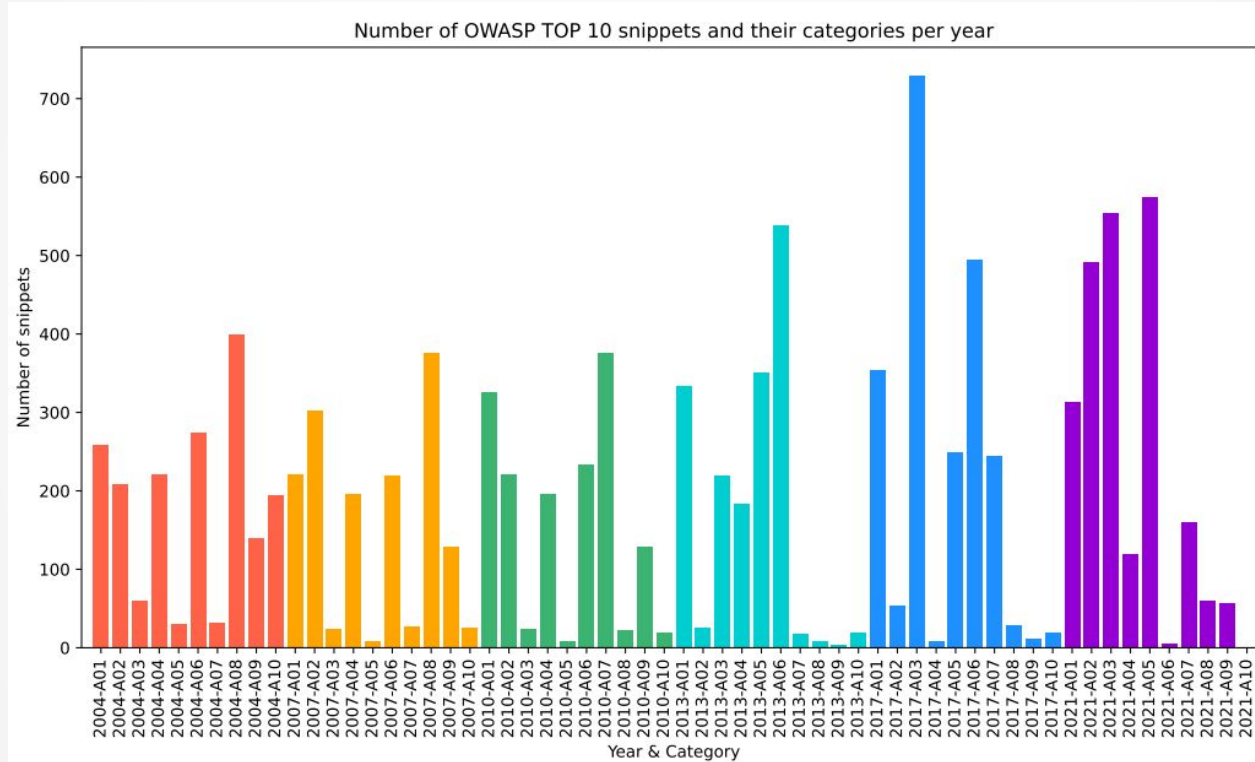
Source 1:https://samate.nist.gov/SRD/testsuites/juliet
Source 2:https://wiki.sei.cmu.edu
Source 3:https://vulncat.fortify.com/
Source 4:https://github.com/returntocorp/semgrep-rules/
Source 5:https://cwe.mitre.org/
Source 6:https://www.sonarqube.org/
Source 7:https://aquasecurity.github.io/tfsec/v1.9.3/checks/aws/api-gateway/enable-cache-encryption/
Source 8:https://docs.microsoft.com/en-us/dotnet/fundamentals/code-analysis/quality-rules/security-warnings
Source 9:https://github.com/snoopysecurity/Vulnerable-Code-Snippets
Source 10:https://github.com/conikeec/seeve
Source 11:https://github.com/fkie-cad/cwe_checker/tree/master/test/artificial_samples

| Source | Number of snippets |
| --- | --- |
| Source 11 | 18 |
| Source 10 | 30 |
| Source 9 | 118 |
| Source 8 | 267 |
| Source 7 | 484 |
| Source 6 | 713 |
| Source 5 | 922 |
| Source 4 | 1135 |
| Source 3 | 2265 |
| Source 2 | 3169 |
| Source 1 | 141305 |

Number of snippets

# Non-Compliant snippets per language

# OWASP TOP 10 - Years

**SIEMENS**
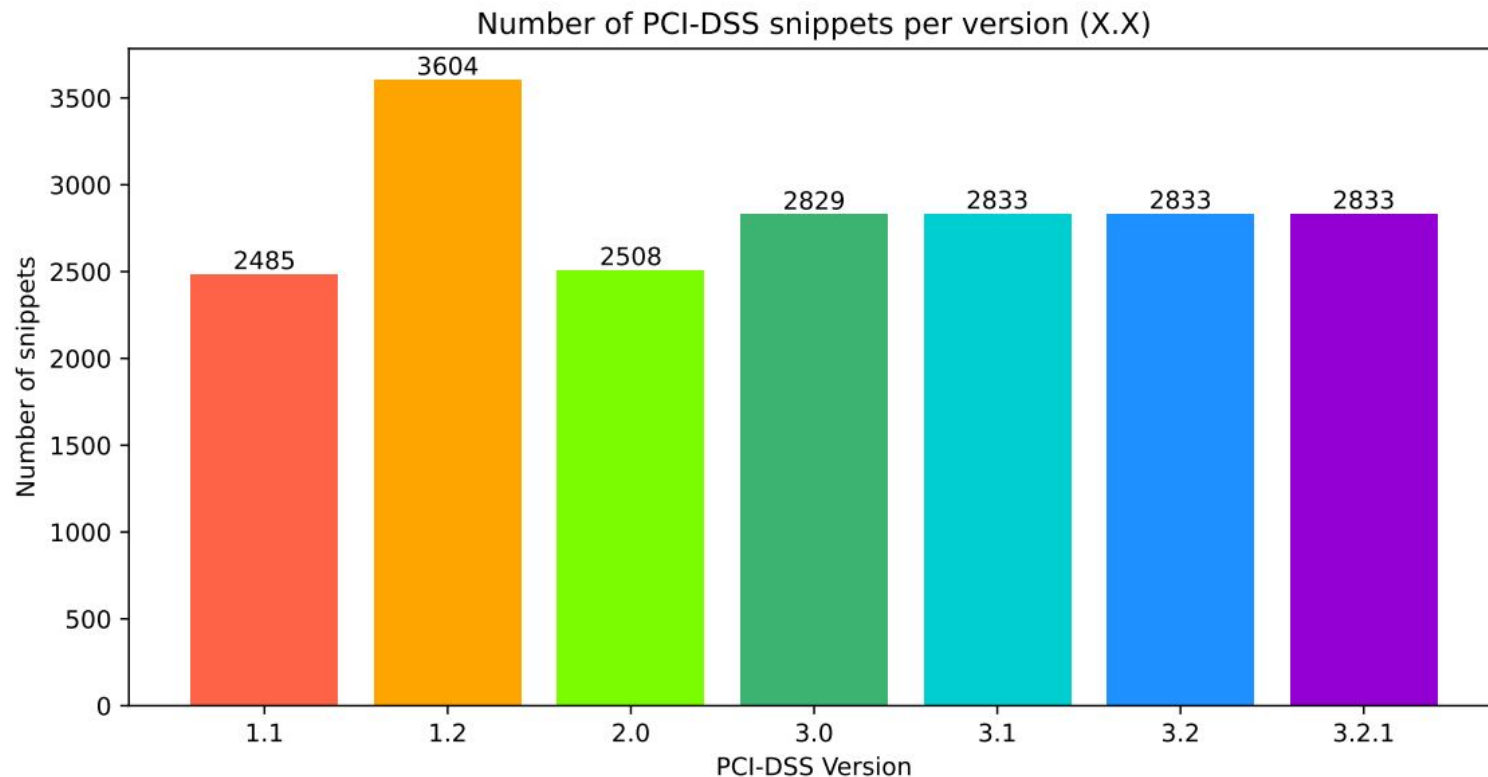*Ingenuity for life*



Number of OWASP TOP 10 snippets per year

# OWASP TOP 10 - Years & Categories



Number of OWASP TOP 10 snippets and their categories per year

# PCI-DSS Version



Number of PCI-DSS snippets per version (X)

# PCI-DSS Categories

# PCI-DSS Sub-Categories



Number of PCI-DSS snippets per version (X.X-X)

# Category Analysis Conclusion

➔ Uniformity of data on high level is not enough

➔ Neither snippet count for OWASP TOP 10 nor PCI-DSS is uniform on a sub category analysis

➔ Training on snippets for OWASP TOP 10 or PCI-DSS
  ➔ Results in heavily biased models towards some (sub)categories

# Juliet Dataset

**TABLE II**
Juliet Data Set snippet count per CWE ID

| C | | Java | | C# | | C++ | |
|---|---|---|---|---|---|---|---|
| *ID* | *Snippet count* | *ID* | *Snippet count* | *ID* | *Snippet count* | *ID* | *Snippet count* |
| CWE 121 | 5906 | CWE 190 | 6555 | CWE 197 | 7695 | CWE 762 | 5180 |
| CWE 78 | 5600 | CWE 191 | 5244 | CWE 190 | 5643 | CWE 122 | 4948 |
| CWE 190 | 5040 | CWE 129 | 4104 | CWE 191 | 3762 | CWE 36 | 3500 |
| ... | | ... | | ... | | ... | |
| CWE 674 | 2 | CWE 499 | 1 | CWE 397 | 1 | CWE 562 | 1 |
| CWE 562 | 2 | CWE 248 | 1 | CWE 366 | 1 | CWE 468 | 1 |
| CWE 561 | 2 | CWE 111 | 1 | CWE 248 | 1 | CWE 440 | 1 |

# Juliet Dataset Analysis Conclusion

➔ Has huge number of snippet examples

➔ Very valuable resource for testing tools

➔ Not good for training Machine learning models
   ◆ Underlying snippet bias

# Conclusion

➔ Where can code snippets be found?
  - ◆ 11 possible sources of information
  - ◆ Not all represented the same (most prominent: C, Java, C#, C++)

➔ What is the quality of the code snippets?
  - ◆ Varies with the programming language
  - ◆ Within a programming language → imbalance between vulnerability categories

**Main conclusion:**
  - ◆ Some programming languages hugely underrepresented
  - ◆ Juliet set - mostly synthetic data
  - ◆ Not clear how good the snippets are to train ML models

**Further work:**
  - ◆ Investigate "real-world" code snippets based on check-in comment

# Keyword Occurrences

## Top-5 | 2-gram

(improper, neutralization)
(integer, overflow)
(buffer, overflow)
(special, elements)
(integer, underflow)

## Top-5 | 3-gram

(overflow, or, wraparound)
(neutralization, of, special)
(special, elements, used)
(integer, underflow, wrap)
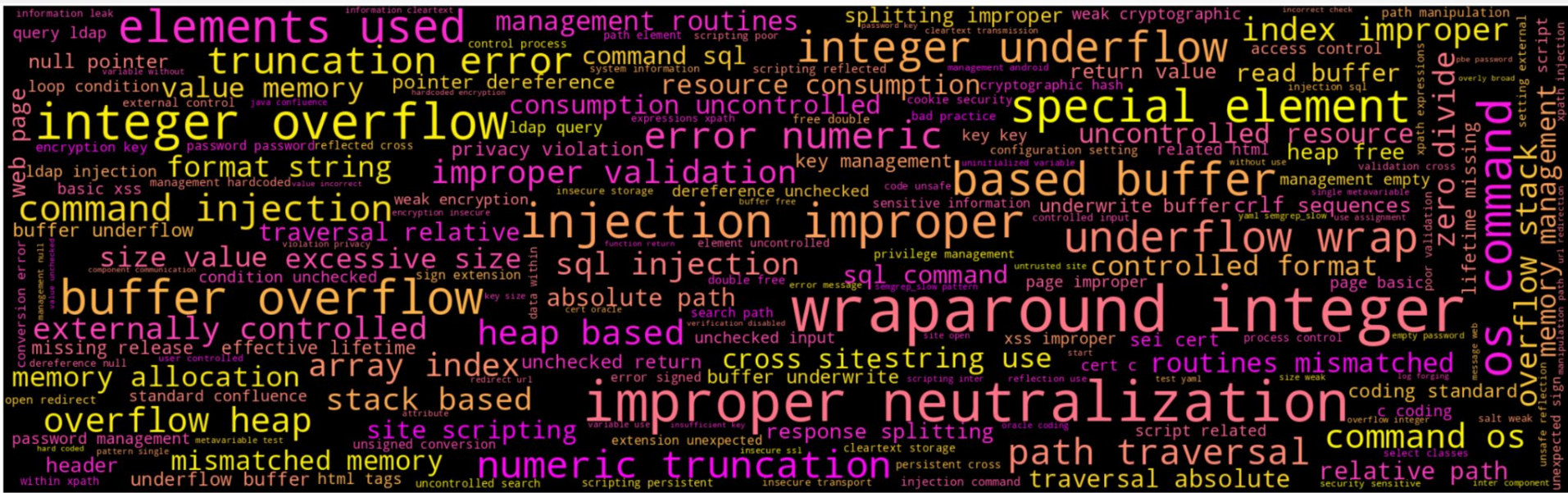(numeric, truncation, error)

## Top-5 | 4-gram

(integer, overflow, or, wraparound)
(neutralization, of, special, elements)
(improper, neutralization, of, special)
(command, os, command, injection)
(improper, validation, of, array)

## Top-5 | 5-gram

(improper, neutralization, of, special, elements)
(integer, underflow, wrap, or, wraparound)
(os, command, os, command, injection)
(improper, validation, of, array, index)
(use, of, externally-controlled, format, string)

# Thank you for Listening 😊

Kaan Oguzhan
*kaan-oguzhan@siemens.com*
Dr. rer. nat. Tiago Espinha Gasiba
*tiago.gasiba@siemens.com*
Akram Louati
*akram.louati@siemens.com*