

CamBench - Cryptographic API Misuse Detection Tool Benchmark Suite*

Michael Schlichtig[†]
Heinz Nixdorf Institute at Paderborn
University
Paderborn, Germany
michael.schlichtig@uni-
paderborn.de

Anna-Katharina Wickert[†]
Technische Universität Darmstadt
Darmstadt, Germany
wickert@cs.tu-darmstadt.de

Stefan Krüger
Independent
Germany
krueger.stefan.research@gmail.com

Eric Bodden
Heinz Nixdorf Institute at Paderborn
University & Fraunhofer IEM
Paderborn, Germany
eric.bodden@uni-paderborn.de

Mira Mezini
Technische Universität Darmstadt
Darmstadt, Germany
mezini@cs.tu-darmstadt.de

ABSTRACT

Context: Cryptographic APIs are often misused in real-world applications. To mitigate that, many cryptographic API misuse detection tools have been introduced. However, there exists no established reference benchmark for a fair and comprehensive comparison and evaluation of these tools. While there are benchmarks, they often only address a subset of the domain or were only used to evaluate a subset of existing misuse detection tools. **Objective:** To fairly compare cryptographic API misuse detection tools and to drive future development in this domain, we will devise such a benchmark. Openness and transparency in the generation process are key factors to fairly generate and establish the needed benchmark. **Method:** We propose an approach where we derive the benchmark generation methodology from the literature which consists of general best practices in benchmarking and domain-specific benchmark generation. A part of this methodology is transparency and openness of the generation process, which is achieved by pre-registering this work. Based on our methodology we design *CamBench*, a fair “Cryptographic API Misuse Detection Tool Benchmark Suite”. We will implement the first version of *CamBench* limiting the domain to Java, the JCA, and static analyses. Finally, we will use *CamBench* to compare current misuse detection tools and compare *CamBench* to related benchmarks of its domain.

CCS CONCEPTS

• **Software and its engineering** → **Software libraries and repositories**; *Software maintenance tools*; **Software verification and validation**; • **Security and privacy** → *Software security engineering*.

KEYWORDS

cryptography, benchmark, API misuse, static analysis

1 INTRODUCTION

Cryptography, hereafter crypto, is widely used in today’s software to ensure the confidentiality of users’ data. Concretely, through crypto, everyone can securely use online banking, purchase a book from the local book shop from their computer, and share sensitive information with their co-workers while working remotely without the fear that the software leaks their data. Unfortunately, previous research results show that developers struggle with the secure usage of crypto APIs and often add vulnerabilities to their code [21, 28, 33, 50]. Most of these vulnerabilities are caused by developers who use a crypto API in a way that is considered insecure by experts, e.g., choosing an outdated crypto hash algorithm like *SHA-1* [28]. We use the term crypto-API misuse or shortly misuse to describe these programming flaws. Further, we focus on Java crypto-APIs to provide concrete results. To support developers and security auditors to identify these misuses, many crypto API misuse detectors such as *CryptoRex* [52], *CryptoLint* [18], *CogniCrypt_{SAST}* [26], and *CryptoGuard* [37] have been developed.

While several tools demonstrate the practical importance and their capabilities to identify issues in real-world code, a fair comparison of the analysis capabilities is difficult. First, all in-the-wild studies were conducted on different applications and domains. For example, *CryptoLint* [18] was evaluated on 11,748 Android apps in 2012 while *CogniCrypt_{SAST}* [26] analyzed 10,000 apps in 2017. While both tools demonstrated their capabilities for Android apps, it is hard to judge how they compare, as they analyzed different apps and the apps evolved over time. Further, other domains of Java software are only analyzed by one tool, e.g., Apache applications, as representatives for large Java projects [37], or Maven artifacts, representatives for Java libraries [26]. Thus, a fair comparison of these analyses and their capabilities is challenging.

An accepted standard in other fields to overcome the above-described problem are benchmarks. While it is time-consuming to create a proper benchmark, the advantages are significant. A benchmark can set a common understanding within the community of the capabilities of the respective solutions, e.g., the strength of a misuse detector. If a benchmark is used for several tools, which

*This study was accepted at the MSR 2022 Registered Reports Track.

[†]Both authors contributed equally to this research.

should be the aim, a benchmark enables a fair comparison of the respective solutions in the field. Further, a proper benchmark can thrive improvements of each tool and form a community to improve the state of the art significantly. Overall, an established benchmark drives the development and research in the community of its domain [12].

Unfortunately, for the domain of crypto-API misuse detection, one still misses a standard benchmark. While we, as a community, have a few custom benchmarks [6, 13, 31, 51], all of them have significant limitations. Some benchmarks are developed for a specific problem, e.g., parametric-based misuses [51], and are not suitable outside of this problem space. Other benchmarks come with a significant bias as their creation was along with the development of the respective detection tool [6].

To solve the above-described problem, we aim to establish **CamBench**, a fair “Cryptographic API Misuse Detection Tool Benchmark Suite” to compare crypto-API misuse detectors. Thus, we will analyze the current state of the art of benchmarking crypto-API misuse detection tools, and how these benchmarks adhere to proper methods for benchmarking. Further, we plan to focus on understanding and deriving specific challenges that the domain of crypto API misuse detection tools poses to the design of benchmarks, e.g., the evolution of security characteristics. Based on this theoretical work, we want to create our fair benchmark to compare crypto-API misuse detectors.

Our benchmark suite consists of *two* benchmarks and *one* heuristic to compare crypto-API misuse detectors. The benchmarks are a collection of real-world applications as well as synthetic test cases to evaluate the analysis capabilities. In addition to the benchmarks, we add an API coverage heuristic that will provide detailed information about the strengths and limitations of a misuse detector with respect to certain API classes. This property is important as an analysis which misses a class may miss a severe vulnerability in an application due to the respective class. Thus, an assessment of this property beside precision and recall is important. To understand the impact of creating an independent benchmark, we will compare *CamBench* against the existing benchmarks.

The first version of *CamBench* focuses on usages of the standard crypto library in Java, the Java Cryptography Architecture (JCA), and static analyses. To support future extensions, e.g., on different APIs and dynamic analyses, one design principle of the benchmark is extensibility. We plan to open-source all our research artifacts along with our tooling. This way, we enable the extension of *CamBench* and provide tooling to simplify the creation of future benchmarks in related domains. Further, our tooling and our theoretical work can further reduce the effort to create fair benchmarks for arbitrary APIs, e.g., misuse of the *Iterator* API.

The early validation of the design of *CamBench* and the benchmark generation methodology is crucial to achieving the aim of a fair crypto-API misuse detector benchmark because it aides in making the process more transparent and open as it also allows for community feedback. Further, the validation is of importance to avoid tuning of *CamBench* in a specific direction by the authors, as well as to avoid requests for changes due to the results by a reviewer who may be an author of one of the benchmarked tools. Moreover, *CamBench* is created by members of the community and as part of the community, we value this effort. Therefore, we believe that

the pre-registration will strengthen the validity of *CamBench* the results we derive, and reduces bias in the creation.

2 PRELIMINARY STUDY: BENCHMARKS COVERING CRYPTOGRAPHY

To understand the differences, strengths, and weaknesses of the existing benchmarks for crypto misuse detectors, we will present our preliminary study on the current state of the art of benchmarks for crypto misuse detection tools. We derived all discussed benchmarks by collecting all existing crypto benchmarks we were aware of and extended them through a forward- and backward search of the related publications. An overview of all identified benchmarks and misuse detectors that were successfully evaluated with them is presented in Table 1. Note, that we did not list any analysis for *MuBench* as the API misuse detection tools used for evaluation did not find any crypto misuse [9]. Similarly, we also omit the misuse detection tools that mark the secure and insecure test cases of the *Ghera Benchmark* as non-security critical (cf. Ranganath et al. [38]) for readability.

CryptoAPI-Bench [6] is the most-recent benchmark and consists of synthetic examples to identify misuses. In total, 171 examples were created along 16 different vulnerabilities for *basic* (40) and *advanced* (131) cases for static analyses. The advanced tests focus on inter-procedurality, field-sensitivity, and path-sensitivity. Thus, challenging the static analysis for several different analysis capabilities. However, the benchmark comprises *all* the vulnerabilities detected by *CryptoGuard* [37].

While *CryptoAPI-Bench* focuses on analysis capabilities, the *Braga Benchmark* [13] focuses on providing developers with the best analysis tool for their team. A team is characterized by their experience (novice or knowledgeable) and their support for cryptographic tasks (unsupported or supported). In addition to the team’s characteristics, the test cases of the benchmark are grouped along *three* complexity classes of crypto misuses. The complexity is represented as the abstraction level required to identify and respectively fix the misuse. For example, identifying an outdated hashing algorithm is considered as low complexity, while the reuse of a nonce is marked as a high complexity problem caused by an insecure system design or architecture. Thus, this benchmark sheds light on different complexities to cover different teams’ knowledge.

In contrast to the previous two benchmarks, the *Parametric Crypto Misuse Benchmark* [51] derived misuse cases from in-the-wild code. All misuses were collected from top Java GitHub projects and cover crypto misuses caused by passing an insecure parameter to a function. Thus, covering one of the major problems, incorrect configurations via parameters, observed in several in-the-wild studies [18, 26] while ignoring more complex misuses, e.g., incorrect call order or misuses caused by architecture flaws. While the previously discussed benchmarks created the benchmarking infrastructure, *Parametric Crypto Misuse Benchmark* extended an existing and well established benchmark for API-misuses [8].

The *Parametric Crypto Misuse Benchmark* benchmark is built upon *MuBench* [8] which is a benchmark for general API misuses, including several crypto misuses in Java. In the publication from 2016, the benchmark includes 18 misuses manually identified from bug fixes on GitHub and SourceForge. All of these cases, except *one*

fix, influence the behavior of the program, such as a crash due to invalid input. Thus, the misuses cover spurious behavior rather than insecure algorithm choices. Later, the benchmark was extended with additional 31 misuses by mining projects via BOA [17] which use the *javax.crypto.Cipher* or the entire *javax.crypto* package, and manually verifying the security of the code [9].

The *Ghera Benchmark* [31] focuses on general vulnerabilities in Android applications being a subset of Java misuses, and thus cover a few crypto misuse cases as well. In total, this benchmark includes five instances of crypto vulnerabilities and eight vulnerabilities caused by networking. While it may be counter-intuitive, the later eight vulnerabilities may be covered by crypto misuse detectors as well when covering SSL and TLS related misuses, e.g., the usage of TLS 1.0. While the number is low, all cases are executable Android applications with a vulnerability along with an exploit of the respective vulnerability.

While the previous benchmarks are developed in academia, the *OWASP Benchmark* [49] is developed by industry. Similar to our motivation, the underlying motivation for *OWASP Benchmark* was to measure the strength and weaknesses of different static analyses to detect vulnerabilities such as crypto misuses. To achieve this aim, version 1.2 consists, similarly to *Ghera Benchmark*, of test cases with real-world web applications that can be exploited. The exploits are categorized along with the different vulnerability types of OWASP Top-10 web vulnerabilities [47], such as command injection and weak cryptography. In total, 246 tests for vulnerabilities are in *OWASP Benchmark*.

The overlap of the static analyses, e.g., *CogniCrypt_{SAST}* [26], *CryptoGuard* [37], and *FindSecBugs* [11], analyzed on the aforementioned benchmarks is rather small as the motivation and domain for which all these benchmarks were created differs. We illustrate the benchmarked static analyses for each previously discussed benchmark in Table 1. Only two tools, namely *FindSecBugs* (three) and *SonarQube* (two), are evaluated on more than one benchmark. The remaining 16 analyses are evaluated only on one benchmark. For example, the *Ghera Benchmark* evaluated analyses for Android vulnerabilities, such as *Amandroid* [48]. On the other hand, tools focused on crypto misuses, such as *CryptoGuard* [37], are evaluated on the crypto-specific *CryptoAPI-Bench*.

3 RESEARCH QUESTIONS

With the research proposed in this registered report, we will answer the *three* research questions discussed in this Section. We start with a question that we already discussed in Section 2, and plan to investigate more in-depth:

RQ1. What is the current state of the art of benchmarking crypto API misuse detection tools and to which extent do benchmarks adhere to appropriate methods in general benchmarking? As we have shown in Section 2, there is a need for a fair benchmark that properly covers the domain of crypto API misuse detection tools in an unbiased way. Existing benchmarks are not suitable for a fair comparison as analysis properties, such as sensitivity, are not considered or not covered completely, or covered misuses are designed alongside one existing tool rather than the analyzed crypto APIs.

Since existing benchmarks are insufficient, a proper generation of a new benchmark is necessary. Key factors for this are transparency and openness, as the successful DaCapo benchmark project has shown [12]. To achieve this goal, performing the benchmark generation via a registered report is a logical consequence, as the proposed methodology is reviewed by experts upfront and the process is made known to the research community.

RQ2. In the area of evaluating cryptographic API misuse detection tools, what requirements does a benchmark need to be tailored to?

For the generation of a fair benchmark, we first need to define our methodology, which we derive from the literature. Understanding the domain-specific characteristics is necessary to develop the requirements a benchmark for evaluation crypto misuse detection tools needs to meet to support their evaluation and comparison. Once the methodology and requirements are defined, the next step is generating *CamBench* accordingly and then evaluating it.

RQ3. Can CamBench meaningfully compare current cryptographic API misuse detection tools? Our final research question is the application of *CamBench* on static analysis tools for the misuse detection of the JCA in Java code. We plan to evaluate two aspects. First, we will use *CamBench* to compare how current crypto API misuse detection tools perform. Second, we will compare *CamBench* to the other benchmarks in the domain of crypto API misuse detection tools (cf. Table 1) and analyze the differences.

Benchmark	A	AC	CC	CG	CO	DN	FD	FS	J	MA	MF	P	SB	SQ	V	X	Y	Z	Source
<i>CryptoAPI-Bench</i> [6]	○	○	●	●	●	○	○	○	○	○	○	○	●	○	○	○	○	○	[6]
<i>Braga Benchmark</i> [13]	○	○	○	○	○	○	○	●	○	○	○	○	○	●	●	●	●	○	[13]
<i>Parametric Crypto Misuse Benchmark</i> [51]	○	○	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○	[51]
<i>MuBench</i> [8]	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	[8]
<i>Ghera Benchmark</i> [31]	●	●	○	○	○	●	●	○	○	●	●	○	○	○	○	○	○	○	[38]
<i>OWASP Benchmark</i> [49]	○	○	○	○	○	○	○	●	●	○	○	●	○	●	○	○	○	●	[14]

Table 1: An overview of all identified benchmarks and the crypto misuse detectors ("*" marks non-academic static analyses. A: *Amandroid* [48], AC: *AppCritique** [7], CC: *CogniCrypt_{SAST}* [26], CG: *CryptoGuard* [37], CO: *Coverity** [22], DN: *DevKnox** [27], FD: *FixDroid* [35], FS: *FindSecBugs** [11], J: *Julia** [45], MA: *Marvin-SA** [41], MF: *MobSF** [5], P: *PMD** [2], SB: *SpotBugs** [3], SQ: *SonarQube** [44], V: *VisualCodeGrepper** [4], X: *Xanitizer** [20], Y: *Yasca** [42], and Z: *OWASP ZAP** [1]) analyzed.

4 METHODOLOGY

In this section, we discuss the theoretical foundations of our benchmark generation methodology and design. We have already begun conducting a literature survey to answer **RQ1**, the preliminary results of which we presented in Section 2. As of now, our results indicate that there is no widely accepted or established benchmark for the evaluation of crypto API misuse detection tools. Existing benchmarks have only been used to compare subsets of crypto misuse detection tools (cf. Table 1). To fully answer this question, we will continue our literature review and supplement our findings. For this purpose, we will adapt the procedure described by Kitchenham et al. [25].

To answer **RQ2** we consider two aspects: First, we discuss what the proper methodology for benchmark generation in general is, and how this can be applied to one domain, namely crypto API misuse detection tool evaluation in Section 4.1. And second, we analyze the special characteristics of this domain that need to be addressed by *CamBench* or in general a benchmark for this domain in Section 4.2. We conclude the theoretical foundations by deriving the domain-specific requirements for *CamBench* and developing the domain-specific methodology for the generation of benchmarks in Section 4.3. Then we use both to formulate the practical application in our execution plan for the generation of *CamBench* and its evaluation to answer **RQ3** in Section 5.

4.1 Literature review on benchmark generation approaches

The creation of a benchmark for a specific domain requires first an understanding of the characteristics of well-established more general benchmarks in the community. So far, our literature review to establish a benchmark for crypto API misuse detection tools has yielded a few relevant approaches of more general benchmarks. For our approach two findings are standing out, namely the *DaCapo benchmarks* [12] and the Automated Benchmark Management System (*ABM*) [16]. The *DaCapo benchmarks* are a collection of widely used open-sourced benchmarks for the Java programming language introduced in 2006 and its latest release in 2018. In their report on the development and establishment of the *DaCapo benchmarks*, Blackburn et al. stressed the importance of an open and transparent process with community feedback. Their success is a compelling argument to follow transparency and community feedback as guiding principles when developing *CamBench*. Similarly, the *Qualitas Corpus* [46] also provides a large collection of real-world open-source Java code. The way the real-world applications are organized and enriched with metadata, also for versioning, is valuable information for the design of benchmarks regarding extensibility. However, the *Qualitas Corpus* was introduced in 2010 and curated, but with its latest update in 2013 it appears to be unmaintained.

The *DaCapo benchmarks* and *Qualitas Corpus* are both benchmarks for Java in general targeting the evaluation of the performance of Java code and Java virtual machines. However, evaluating tools that detect misuses of crypto APIs is a different task. The main focus here rather lies on the correctness of the results, i.e., precision and recall of the detected issues and the soundness

of the analysis itself. Moreover, how usable such a crypto misuse detection tool is, also relies on its precision and recall [15] – too many false positives are one of the main reasons why developers omit using static analysis tools. Such domain-specific requirements can be better addressed by the approach of Nguyen Quang Do et al. called *ABM* [16] since domain-specific data sets are needed to evaluate specific research questions. They describe how to semi-automatically generate and maintain domain-specific benchmark suites. The benchmark suites are collected with a well-described process of crawling buildable open-source projects of the target domain. Therefore, they provide a representative set of real-world applications relevant in the domain. We derive our methodology from the experience that can be taken from the *DaCapo benchmarks* and the proposed procedure by *ABM* and provide a more in-depth explanation of the execution in Section 5. Moreover, *ABM* already describes that filtering the real-world applications for relevance is integral to being representative. Hence, understanding the requirements of a domain is key to generating a representative benchmark suite for it.

4.2 Analysis of the domain for benchmarking cryptographic API misuse detection tools

While Section 2 presents preliminary findings of benchmarks targeted for crypto API misuse detectors, a structured literature review [25] is required to enrich our findings and ensure their completeness. The goal that crypto API misuse detection tools aim to achieve is detecting and reporting crypto misuses in applications that use crypto APIs. Employing crypto in general sets the focus on security. Therefore, when comparing misuse detection tools it is important to evaluate their performance in terms of crypto API coverage, precision and recall, and analysis capabilities. The coverage is important for developers to know whether their code using a crypto API is completely or only partially covered, e.g., common known vulnerabilities. Research on reasons why misuse detection tools are rarely used has shown that precision and recall of a tool play a key role [15] – too many false positives are a major reason for developers to omit using a tool [23]. Moreover, the analysis evaluation of analysis capabilities is important as well. Whether an analysis is intraprocedural, inter-procedural, flow-sensitive, context-sensitive, field-sensitive, path-sensitive, and object-sensitive or not has direct implications on what the analysis is actually able to detect. Unfortunately, the benchmarks identified in Section 2 either ignore such properties or only include a subset, such as *CryptoAPI-Bench* [6].

4.3 Domain-specific Benchmark Requirements

Concluding the theoretical foundations of our methodology, we specify the requirements specific to the domain of benchmarking crypto API misuse detection tools and propose our benchmark design based on the results of the previous sections. Following the take-aways from the *DaCapo benchmarks* [12] (transparent and open process, a variety of real-world applications maximizing the domain coverage, easy to use, and providing metrics), *CamBench* will be a benchmark suite consisting of several benchmarks. A good example for easy to use is the well usable deployment system of *MuBench* [8], which delivers its execution pipeline as a Docker

container and provides extensive documentation on how to use and adapt the benchmark as well as on how to contribute. Moreover, we enrich the collection process of real-world application with the semi-automated approach of *ABM* [16].

Based on these presented approaches, we derive the following requirements that *CamBench* must meet:

- (1) The benchmark generation process needs to be transparent and consider community feedback.
- (2) The benchmark contains a representative and diverse collection of real-world applications using crypto APIs. Projects fulfill these properties: *open-source*, *compilable (source code and binaries are available)*, and *representative of the domain's real-world application, e.g., diverse size, context, . . .*
- (3) The benchmark is open-sourced.
- (4) The benchmark is provided with a deployment system.

CamBench will be a suite of benchmarks covering the following three main aspects: (a) *real-world applications with crypto API usage* to evaluate the performance of misuse detection tools on relevant applications [16], (b) *analysis capabilities*, and (c) *crypto API coverage* to evaluate which fraction of the API a misuse detection tool supports. We added *analysis capabilities* and *crypto API coverage* to accommodate the requirements of the domain of benchmarking crypto API misuse detection tools.

To facilitate the evaluation of analysis capabilities, we will develop synthetic test cases for all relevant properties. When building an analysis tool, it is necessary to make assumptions that can lead to unsoundness [30, 40]. Therefore, evaluating the analysis capabilities which are part of such assumptions is important to draw implications concerning a tool's unsoundness and precision. The properties we plan to include in *CamBench* are flow-sensitivity, context-sensitivity, field-sensitivity, object-sensitivity, and path-sensitivity.

CamBench will provide a benchmark for testing the coverage of Java crypto-APIs. Similar to *ABM*'s [16] focus on automation our goal is to develop a process with a high degree of automation for the generation of the benchmarks for the coverage of the JCA. This will help in maintaining and updating the benchmark in the future.

Furthermore, we will enrich the test cases with metadata similar to *MuBench* [8] where misuse code examples are specified with YAML files enriched by additional information like misuse type and description. They also provide examples for correct usages. Moreover, the metadata comprises information whether a test case is a (in-)secure crypto usage (whether a crypto usage is (in-)secure might also be context-dependent, e.g., using a hash algorithm like MD5 is considered insecure [24], yet using MD5 in the context of file validation [10] might still be acceptable), usage category [9], and the severity. Moreover, we will consider metadata on versioning of real-world projects as proposed in the *Qualitas Corpus* [46].

To summarize, *CamBench* will provide a collection of *real-world applications* with crypto API usage following *ABM* [16], synthetic tests to evaluate the crypto API misuse detectors *analysis capabilities*, and a collection of test cases for checking the *crypto API coverage* of the JCA.

5 EXECUTION PLAN

In this section, we first describe the generation of *CamBench* and its three kinds of benchmarks in Section 5.1, namely *real-world applications*, *analysis capabilities*, and *crypto API coverage*. Then we will discuss the publishing and deployment process of *CamBench* in Section 5.2. Finally, we describe our intended evaluation in Section 5.3.

5.1 Generation Plan

For the generation of *CamBench*'s three benchmarks, we now describe our plan to execute our methodology (cf. Section 4). For the concrete implementation of our methodology for the first version of *CamBench*, we set the scope to the standard crypto library in Java, the JCA, and static analysis tools detecting misuses of it.

Benchmark for real-world applications. For the collection and creation of the real-world application benchmark, we adapt the process specified by *ABM* [16]. This process consists of six steps of collecting, filtering, and building open-source projects. The first three steps are (i) *mining open-source projects*, (ii) *filtering for active projects* and (iii) *filtering for known build systems*. These steps can be automated. The subsequent steps require manual effort. In their instantiation of *ABM* for the domain of Java business web applications, Nguyen Quang Do et al. chose GitHub as a platform for mining open-source projects [16]. Since this has worked successfully and GitHub is a popular and well-used platform for hosting open-source projects, we intend to use it, too. The next step is (iv) *downloading and building the projects*. The real-world application benchmark only includes projects that build with standard build systems. Another filtering step follows where, in our case, we will (v) *check whether the JCA is used* in the projects. The last step is (vi) *compiling the binaries* and adding projects with binaries to the benchmark dataset (cf. requirements in Section 4).

Furthermore, we add a step to the approach of *ABM* [16] by adding metadata files including information, on whether the usage is (in-)secure for selected usages of the JCA. To automatically identify usages of the JCA in code, we plan to employ static analyses, such as adapted versions of *CogniCrypt* [26] or *CryptoGuard* [37]. For the ground-truth, whether a JCA usage is secure or insecure, we will use the guidelines of NIST [34], SOG-IS [43], or BSI [19] for manual labeling. Since the guidelines of the BSI are updated regularly and are the ones updated most recently, we will primarily be using them. Furthermore, the employed guidelines can be documented in the metadata files. The design of the metadata files will be derived from *MuBench* [8] and take the experiences concerning versioning from *Qualitas Corpus* [46] into account. The expected dataset of the real-world application will contain too many usages of the JCA for the manual creation of the metadata files. Therefore, we will select a representative subset of the JCA usages to write the metadata files.

Benchmark for analysis capabilities. The benchmark for evaluating analysis capabilities will consist of synthetic test cases specifically designed to test relevant properties. The test cases for the analysis capabilities require manual effort of domain experts. They

will include test cases for flow-sensitivity, context-sensitivity, field-sensitivity, object-sensitivity, and path-sensitivity, as well as intraprocedural and inter-procedural test cases. The selected sensitivities are accepted as relevant capabilities for static analyses [29, 36]. We will develop all test cases for this benchmark and also provide the test cases with the metadata files for (mis-)uses.

Heuristic for crypto API coverage. During our preliminary study, we identified that understanding the covered API classes is of importance for crypto-API misuse detectors. A detector which misses an API class may miss a vulnerability in an application while having a high precision and recall. Thus, for the assessment of the detectors, being aware of the covered classes is of importance. While test generators seem to be a good choice for this task and aid our aim of maintainability and extensibility of *CamBench* they are inadequate. Previous approaches summarized in a survey [32] have shown to be insufficient for generating and establishing a ground-truth needed for *CamBench*. Instead of using test generators, we propose employing a heuristic that harnesses the two other benchmarks of *CamBench* to create an API coverage feature.

This feature is inspired by *Hermes* [39] which is a framework to assess collections like the *Qualitas Corpus* and uses queries to derive subsets suited for analyzing, e.g., a specific API. With this heuristic, we can extract all classes and methods of the JCA and match them with our benchmarks for *real-world applications* and *analysis capabilities*. We plan to use the heuristic to guide the effective generation of our benchmark. First, the heuristic will enable us to select a diverse set of real-world JCA usages to enrich them with the metadata. Second, as we will develop the synthetic test cases ourselves, the heuristic allows us to include classes and methods for the JCA that are not covered by the *real-world applications benchmark*. Hence, this heuristic approach covers at least all real-world parts of the API from the *real-world applications benchmark* as well as additional classes and methods that only occur in the synthetic *analysis capabilities benchmark*. Overall, we assume that this heuristic approach should yield an almost complete coverage of JCA class and method usages for a highly used API like the JCA.

5.2 Publishing and Deployment Plan

Open access and ease of use are key factors for establishing a benchmark as describe in Section 4. Hence, it is important to publish *CamBench* on GitHub with proper instructions and documentation. To encourage, include, and harness community feedback, we will make the GitHub repository accessible during the creation of *CamBench*, e.g., via pull requests. We will also contact the authors of existing crypto API misuse detection tools and benchmarks (cf. Table 1).

Here we present our plans for the GitHub repository. We will provide execution support for *CamBench* similar to *MuBench* [8]. Thus, alongside the benchmark itself, we will provide the tooling we use for the generation together with documentation to support updating, maintaining, extending, and contributing to *CamBench*. In addition, we plan to improve the accessibility of our tool by providing essential information, such as the number of test cases or the covered API classes, directly on our *README*-page and using the GitHub release functionality. To ensure the continuous quality of our benchmark, we will use GitHub actions for continuous

deployment. Furthermore, *CamBench* will be open access with a unique Digital Object Identifier (DOI).

5.3 Evaluation plan

After creating the first version of *CamBench* we will perform the evaluation to answer **RQ3**. The evaluation will consist of two main aspects: First, how do current crypto API misuse detection tools perform on *CamBench*? The primary motivation for this registered report is comparing crypto API misuse detection tools with a fair benchmark explicitly developed for this domain. Once *CamBench* is created, such a comparison and evaluation will be possible and therefore will be conducted. Furthermore, we envision to provide a score section, similar to OWASP [49], in the GitHub repository of *CamBench* where researchers can add their results with experiment documentation via pull request. Second, we want to know how current crypto API misuse detection tool benchmarks compare to *CamBench* in practice. We expect this to be a theoretical comparison as the approaches and goals of existing benchmarks are diverse (cf. Table 1). The generation approach of *CamBench* differs from existing benchmarks (cf. Section 2, Table 1) and is intended to serve as a reference benchmark for the community. Therefore, we expect *CamBench* to differ from the other custom benchmarks in its composition and evaluation of the analyses. Hence, we will investigate these differences, also to understand how our generation process for the test cases compares against the instantiated benchmarks. Additionally, we will evaluate our heuristic approach for measuring the coverage of the JCA.

6 THREATS TO VALIDITY

For the first version of *CamBench*, we introduced some limitations (cf. Section 1) to have a clear target domain and develop a tangible benchmark. These limitations are choosing a specific programming language in Java, a specific API in the JCA, and finally focusing on static analysis tools. However, our methodology (cf. Section 4) is designed for semi-automated benchmark generation and extensibility. Hence, extending *CamBench* to support other APIs or dynamic analyses has very much been taken into consideration during the initial development as a reasonable future step and is as a result not inhibited whatsoever by the design. Furthermore, instantiating benchmarks for other programming languages by applying our methodology is also possible. However, these extensions would each require a lot of additional effort and leave a benchmark that is significantly changed in shape and size as the new real-world applications need to be mined with new criteria, and thus, the benchmarks for analysis capabilities and API coverage need adaptation, too.

Besides this, the generation of the benchmarks for *analysis capabilities* and *API coverage* are subject to manual work by the authors. Since the authors are experts regarding static analysis and crypto API misuse detection, the chosen domain limitations make them an appropriate choice to develop and label the test cases. Furthermore, community feedback is part of the generation process and will therefore help to mitigate this threat to validity.

Lastly, employing the heuristic for the *API coverage* benchmark is a new approach that we will evaluate. Employing the heuristic seems to be a reasonable and practical approach to create the

API coverage feature. However, this approach is untested and needs evaluation.

7 IMPLICATIONS

Benchmarks generally drive development and research in their fields as they provide measurable optimization goals [12]. Therefore, establishing a reference benchmark for the domain of benchmarking crypto API misuse detection tools in the form of *CamBench* with its open and transparent generation process will help the community manifold: *CamBench* will allow to evaluate and compare current crypto API misuse detection tools on the same benchmark which has not been done, yet. Having this comparison, future development can reference the comparison as well as *CamBench* in new contributions.

Furthermore, by deriving, extending, and defining the process for the creation of *CamBench*, we hope to provide a template for creating and maintaining benchmarks for specific domains in need of a reference benchmark and thereby steering current developments towards more open and reproducible research.

ACKNOWLEDGMENTS

Funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – SFB 1119 – 236615297 and by the German Federal Ministry of Education and Research and the Hessen State Ministry for Higher Education, Research and the Arts within their joint support of the National Research Center for Applied Cybersecurity ATHENE.

REFERENCES

- [1] 2022. OWASP ZAP. <https://github.com/zaproxy/zaproxy> accessed: 2015-06-03.
- [2] 2022. PMD - source code analyzer. <https://github.com/pmd/pmd> accessed: 2012-07-11.
- [3] 2022. spotbugs/spotbugs. <https://github.com/spotbugs/spotbugs> accessed: 2016-11-04.
- [4] 2022. VisualCodeGrepper. <https://github.com/nccgroup/VCG> accessed: 2015-01-07.
- [5] Aijn Abraham, Magaofei, Matan Dobrushin, and Vincent Nadal. 2022. Mobile Security Framework (MobSF). <https://github.com/MobSF/Mobile-Security-Framework-MobSF> accessed: 2015-01-31.
- [6] Sharmin Afrose, Sazzadur Rahaman, and Danfeng Yao. 2019. Cryptoapi-bench: A comprehensive benchmark on java cryptographic api misuses. In *IEEE Cybersecurity Development (SecDev)*. IEEE.
- [7] Booz Allen. 2022. AppCritique. <https://appcritique.boozallen.com/> When we accessed (2022-01-26) AppCritique, it seemed that the project is deprecated.
- [8] Sven Amann, Sarah Nadi, Hoan A Nguyen, Tien N Nguyen, and Mira Mezini. 2016. MUBench: A benchmark for API-misuse detectors. In *2016 ACM 13th International Conference on Mining Software Repositories (MSR)*. ACM.
- [9] Sven Amann, Hoan Anh Nguyen, Sarah Nadi, Tien N Nguyen, and Mira Mezini. 2018. A systematic evaluation of static api-misuse detectors. In *IEEE Transactions on Software Engineering (TSE)*, Vol. 45. IEEE.
- [10] Amazon Web Services (AWS). 2022. How can I check the integrity of an object uploaded to Amazon S3? <https://aws.amazon.com/de/premiumsupport/knowledge-center/data-integrity-s3/>, accessed on 2022-01-21.
- [11] Philippe Arteau. 2020. Find Security Bugs 1.11.0. <https://find-sec-bugs.github.io/> accessed: 2022-01-26.
- [12] S. Blackburn, R. Garner, C. Hoffmann, A. Khang, K. McKinley, R. Bentzur, A. Diwan, D. Feinberg, D. Frampton, S. Guyer, M. Hirzel, A. Hosking, and et. al. 2006. The DaCapo Benchmarks: Java Benchmarking Development and Analysis. In *21st ACM SIGPLAN conference on Object-oriented programming systems, languages, and applications (OOPSLA)*. ACM.
- [13] Alexandre Braga, Ricardo Dahab, Nuno Antunes, Nuno Laranjeiro, and Marco Vieira. 2017. Practical evaluation of static analysis tools for cryptography: Benchmarking method and case study. In *IEEE 28th International Symposium on Software Reliability Engineering (ISSRE)*. IEEE.
- [14] Elisa Burato, Pietro Ferrara, and Fausto Spoto. 2017. Security analysis of the OWASP benchmark with Julia. *First Italian Conference on Cybersecurity (ITASEC)*.
- [15] Maria Christakis and Christian Bird. 2016. What Developers Want and Need from Program Analysis: An Empirical Study. In *31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*. ACM.
- [16] Lisa Nguyen Quang Do, Michael Eichberg, and Eric Bodden. 2016. Toward an Automated Benchmark Management System. In *5th ACM SIGPLAN International Workshop on State Of the Art in Program Analysis (SOAP)*. ACM.
- [17] Robert Dyer, Hoan Anh Nguyen, Hridesh Rajan, and Tien N Nguyen. 2013. Boa: A language and infrastructure for analyzing ultra-large-scale software repositories. In *ACM 35th International Conference on Software Engineering (ICSE)*. ACM.
- [18] Manuel Egele, David Brumley, Yanick Fratantonio, and Christopher Kruegel. 2013. An empirical study of cryptographic misuse in android applications. In *2013 ACM SIGSAC conference on Computer and Communications security (CCS)*. ACM.
- [19] German Federal Office for Information Security (BSI). 2022. BSI TR-02102-1: "Cryptographic Mechanisms: Recommendations and Key Lengths". <https://www.bsi.bund.de/SharedDocs/Downloads/EN/BSI/Publications/TechGuidelines/TG02102/EN> accessed on 2022-01-21.
- [20] RIGS IT GmbH. 2022. Xanitizer. <https://www.xanitizer.com> When we accessed (2022-01-26) the webpage, the Xanitizer page was under maintenance.
- [21] Peter Leo Gorski, Luigi Lo Iacono, Dominik Wermke, Christian Stransky, Sebastian Möller, Yasemin Acar, and Sascha Fahl. 2018. Developers Deserve Security Warnings, Too: On the Effect of Integrated Security Advice on Cryptographic API Misuse. In *14th Symposium on Usable Privacy and Security (SOUPS)*. USENIX Association.
- [22] Synopsys Inc. 2022. Coverity Scan - Static Analysis. <https://scan.coverity.com/> accessed: 2022-01-26.
- [23] B. Johnson, Y. Song, E. Murphy-Hill, and R. Bowdidge. 2013. Why don't software developers use static analysis tools to find bugs?. In *ACM 35th International Conference on Software Engineering (ICSE)*. ACM.
- [24] Alok Kumar Kasgar, Mukesh Kumar Dhariwal, Neeraj Tantubay, and Hina Malviya. 2013. A review paper of message digest 5 (MD5). In *International Journal of Modern Engineering & Management Research*, Vol. 1.
- [25] Barbara Kitchenham. 2004. Procedures for performing systematic reviews. *Keele, UK, Keele University* 33.
- [26] Stefan Krüger, Johannes Späth, Karim Ali, Eric Bodden, and Mira Mezini. 2019. Crysl: An extensible approach to validating the correct usage of cryptographic apis. In *IEEE Transactions on Software Engineering (TSE)*. IEEE.
- [27] XYSEC Labs. 2022. Devknox - Security Plugin for Android Studio. <https://devknox.io/> accessed: 2022-01-26.
- [28] David Lazar, Haogang Chen, Xi Wang, and Nickolai Zeldovich. 2014. Why does cryptographic software fail? A case study and open problems. In *5th Asia-Pacific Workshop on Systems (APSys)*. ACM.
- [29] Li Li, Tegawendé F. Bissyandé, Mike Papadakis, Siegfried Rasthofer, Alexandre Bartel, Damien Octeau, Jacques Klein, and Le Traon. 2017. Static analysis of android apps: A systematic literature review. In *Information and Software Technology*, Vol. 88.
- [30] Benjamin Livshits, Manu Sridharan, Yannis Smaragdakis, Ondřej Lhoták, J Nelson Amaral, Bor-Yuh Evan Chang, Samuel Z Guyer, Uday P Khedker, Anders Møller, and Dimitrios Vardoulakis. 2015. In defense of soundness: A manifesto. *Commun. ACM* 58, 2.
- [31] Joydeep Mitra and Venkatesh-Prasad Ranganath. 2017. Ghera: A repository of android app vulnerability benchmarks. In *13th International Conference on Predictive Models and Data Analytics in Software Engineering (PROMISE)*. ACM.
- [32] Ahmad Mustafa, Wan MN Wan-Kadir, and I Ibrahim. 2017. Comparative evaluation of the state-of-art requirements-based test case generation approaches. In *International Journal on Advanced Science, Engineering and Information Technology*, Vol. 7.
- [33] Sarah Nadi, Stefan Krüger, Mira Mezini, and Eric Bodden. 2016. Jumping through hoops: Why do Java developers struggle with cryptography APIs?. In *38th IEEE/ACM International Conference on Software Engineering (ICSE)*. ACM.
- [34] National Institute of Standards and Technology (NIST). 2022. SP 800-175B Rev. 1 Guideline for Using Cryptographic Standards in the Federal Government: Cryptographic Mechanisms. <https://csrc.nist.gov/publications/detail/sp/800-175b/rev-1/final> accessed: 2022-01-21.
- [35] Duc Cuong Nguyen, Dominik Wermke, Yasemin Acar, Michael Backes, Charles Weir, and Sascha Fahl. 2017. A stitch in time: Supporting android developers in writingsecure code. In *2017 ACM SIGSAC Conference on Computer and Communications Security (CCS)*. ACM.
- [36] Lina Qiu, Yingying Wang, and Julia Rubin. 2018. Analyzing the analyzers: Flow-droid/iccta, amandroid, and droidsafe. In *27th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*. ACM.
- [37] Sazzadur Rahaman, Ya Xiao, Sharmin Afrose, Fahad Shaon, Ke Tian, Miles Frantz, Murat Kantarcioglu, and Danfeng Yao. 2019. Cryptoguard: High precision detection of cryptographic vulnerabilities in massive-sized java projects. In

- 2019 ACM SIGSAC Conference on Computer and Communications Security (CCS). ACM.
- [38] Venkatesh-Prasad Ranganath and Joydeep Mitra. 2020. Are free android app security analysis tools effective in detecting known vulnerabilities?. In *Empirical Software Engineering*, Vol. 25. Springer.
- [39] Michael Reif, Michael Eichberg, Ben Hermann, and Mira Mezini. 2017. Hermes: assessment and creation of effective test corpora. In *6th ACM SIGPLAN International Workshop on State Of the Art in Program Analysis (SOAP)*. ACM.
- [40] Thomas Reps. 2000. Undecidability of Context-Sensitive Data-Dependence Analysis. *ACM Trans. Program. Lang. Syst.* 22, 1.
- [41] Joaquín Rinaudo and Juan Heguiabehere. 2022. Marvin Static Analyzer. <https://github.com/programa-stic/Marvin-static-Analyzer> accessed: 2022-01-26.
- [42] Michael Scovetta and Daniel Ruf. 2022. Yasca. <https://github.com/scovetta/yasca> when we accessed (2022-01-26) the GitHub page, the project was deprecated.
- [43] Senior Officials Group Information Systems Security (SOG-IS). 2022. SOG-IS Crypto Evaluation Scheme Agreed Cryptographic Mechanisms. https://www.sogis.eu/uk/supporting_doc_en.html, accessed: 2022-01-21.
- [44] SonarSource. 2022. SonarQube. <https://github.com/SonarSource/sonarqube> accessed: 2011-01-05.
- [45] Fausto Spoto. 2016. The Julia Static Analyzer for Java. In *Static Analysis*, Xavier Rival (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg.
- [46] Ewan Tempero, Craig Anslow, Jens Dietrich, Ted Han, Jing Li, Markus Lumpe, Hayden Melton, and James Noble. 2010. The Qualitas Corpus: A curated collection of Java code for empirical studies. In *Asia Pacific Software Engineering Conference (APSEC)*. IEEE.
- [47] Andrew van der Stock, Brian Glas, Neil Smithline, and Torsten Giger. 2021. OWASP Top 10:2021. <https://owasp.org/Top10/> accessed: 2022-01-23.
- [48] Fengguo Wei, Sankardas Roy, and Xinming Ou. 2014. Amandroid: A precise and general inter-component data flow analysis framework for security vetting of android apps. In *2014 ACM SIGSAC conference on Computer and Communications security (CCS)*. ACM.
- [49] Dave Wichers. 2022. OWASP Benchmark | OWASP Foundation. <https://owasp.org/www-project-benchmark/> accessed: 2022-01-23.
- [50] Anna-Katharina Wickert, Lars Baumgärtner, Florian Breitfelder, and Mira Mezini. 2021. Python Crypto Misuses in the Wild. In *15th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*. ACM.
- [51] Anna-Katharina Wickert, Michael Reif, Michael Eichberg, Anam Dodhy, and Mira Mezini. 2019. A dataset of parametric cryptographic misuses. In *16th International Conference on Mining Software Repositories (MSR)*. IEEE/ACM.
- [52] Li Zhang, Jiongyi Chen, Wenrui Diao, Shanqing Guo, Jian Weng, and Kehuan Zhang. 2019. CryptoREX: Large-scale Analysis of Cryptographic Misuse in IoT Devices. In *22nd International Symposium on Research in Attacks, Intrusions and Defenses (RAID)*. USENIX Association.