



HAL
open science

Evaluating the Reusability of Android Static Analysis Tools

Jean-Marie Mineau, Jean-François Lalande

► **To cite this version:**

Jean-Marie Mineau, Jean-François Lalande. Evaluating the Reusability of Android Static Analysis Tools. ICSR 2024 - 21st International Conference on Software and Systems Reuse, Jun 2024, Limassol, Cyprus. pp.153-170, 10.1007/978-3-031-66459-5_10 . hal-04557993

HAL Id: hal-04557993

<https://centralesupelec.hal.science/hal-04557993>

Submitted on 24 Apr 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Evaluating the Reusability of Android Static Analysis Tools

Jean-Marie Mineau^[0009-0001-5233-3056] and Jean-Francois Lalande^[0000-0003-4984-2199]

CentraleSupélec, Inria, Univ Rennes, CNRS IRISA, Rennes, France
`firstname.lastname@inria.fr`

Abstract. Reproducibility and reusability in computer science experiments become a requirement for research works. Reproducibility ensures that results can be confirmed by using the same dataset and software of previous papers. Reusability helps other researchers to build new approaches with distributed software artifacts. For researchers in the field of security of mobile platforms, ensuring reproducibility and reusability is difficult to implement. In particular for reusability, datasets of Android applications may contain recent applications that past analysis software cannot process. As a consequence, past software produced by researchers may be difficult to reuse, which endangers the reproducibility of research. This paper intends to explore the reusability of past software dedicated to static analysis of Android applications. We pursue the community effort that identified publications between 2011 and 2017 that perform static analysis of mobile applications and we propose a method for evaluating the reusability of the associated tools. We extensively evaluate the success or failure of these tools on a dataset containing Android applications that can have up to six years of distance from the original publication. We also measure the influence of some important characteristics of the application such as being a goodware or a malware or the application size. Our results show that 54.5% of the evaluated tools are no longer usable and that the size of the bytecode and the min SDK version have the greatest influence on the reusability of tested tools.

1 Introduction

Android is the most used mobile operating system since 2014, and since 2017, it even surpasses Windows all platforms combined¹. The public adoption of Android is confirmed by application developers, with 1.3 millions apps available in the Google Play Store in 2014, and 3.5 millions apps available in 2017². Its popularity makes Android a prime target for malware developers. Consequently, Android has also been an important subject for security research. In the past fifteen years, the research community released many tools to detect or analyze malicious behaviors in applications. Two main approaches can be distinguished: static and

¹ <https://gs.statcounter.com/os-market-share#monthly-200901-202304>

² <https://www.statista.com/statistics/266210>

dynamic analysis [18]. Dynamic analysis requires to run the application in a controlled environment to observe runtime values and/or interactions with the operating system. For example, an Android emulator with a patched kernel can capture these interactions but the modifications to apply are not a trivial task. As a consequence, a lot of efforts have been put in static approaches, which is the focus of this paper.

The usual goal of a static analysis is to compute data flows to detect potential information leaks [33, 31, 5, 15, 12, 23, 16] by analyzing the bytecode of an Android application. The associated developed tools should support the Dalvik bytecode format, the multiplicity of entry points, the event driven architecture of Android applications, the interleaving of native code and bytecode, possibly loaded dynamically, the use of reflection, to name a few. All these obstacles threaten the research efforts. When using a more recent version of Android or a recent set of applications, the results previously obtained may become outdated and the developed tools may not work correctly anymore.

In this paper, we study the reusability of open source static analysis tools that appeared between 2011 and 2017, on a recent Android dataset. The scope of our study is **not** to quantify if the output results are accurate for ensuring reproducibility, because all the studied static analysis tools have different goals in the end. On the contrary, we take as hypothesis that the provided tools compute the intended result but may crash or fail to compute a result due to the evolution of the internals of an Android application, raising unexpected bugs during an analysis. This paper intends to show that sharing the software artifacts of a paper may not be sufficient to ensure that the provided software would be reusable.

Thus, our contributions are the following. We carefully retrieved static analysis tools for Android applications that were selected by Li *et al.* [18] between 2011 and 2017. We contacted the authors, whenever possible, for selecting the best candidate versions and to confirm the good usage of the tools. We rebuild the tools in their original environment and we plan to share our Docker images with this paper. We evaluated the reusability of the tools by measuring the number of successful analysis of applications taken in a custom dataset that contains more recent applications (62 525 in total). The observation of the success or failure of these analysis enables us to answer the following research questions:

- RQ1:** What Android static analysis tools that are more than 5 years old are still available and can be reused without crashing with a reasonable effort?
- RQ2:** How the reusability of tools evolved over time, especially when analyzing applications that are more than 5 years far from the publication of the tool?
- RQ3:** Does the reusability of tools change when analyzing goodware compared to malware?

The paper is structured as follows. Section 2 presents a summary of previous works dedicated to Android static analysis tools. Section 3 presents the methodology employed to build our evaluation process and Section 4 gives the associated experimental results. Section 5 discusses the limitations of this work and gives some takeaways for future contributions. Section 6 concludes the paper.

2 Related Work

We review in this section the past existing datasets provided by the community and the papers related to static analysis tools reusability.

2.1 Application Datasets

Computing if an application contains a possible information flow is an example of a static analysis goal. Some datasets have been built especially for evaluating tools that are computing information flows inside Android applications. One of the first well known dataset is DroidBench, that was released with the tool Flowdroid [2]. Later, the dataset ICC-Bench was introduced with the tool Amandroid [33] to complement DroidBench by introducing applications using Inter-Component data flows. These datasets contain carefully crafted applications containing flows that the tools should be able to detect. These hand-crafted applications can also be used for testing purposes or to detect any regression when the software code evolves. Contrary to real world applications, the behavior of these hand-crafted applications is known in advance, thus providing the ground truth that the tools try to compute. However, these datasets are not representative of real-world applications [26] and the obtained results can be misleading.

Contrary to DroidBench and ICC-Bench, some approaches use real-world applications. Bosu *et al.* [5] use DIALDroid to perform a threat analysis of Inter-Application communication and published DIALDroid-Bench, an associated dataset. Similarly, Luo *et al.* released TaintBench [22] a real-world dataset and the associated recommendations to build such a dataset. These datasets confirmed that some tools such as Amandroid [33] and Flowdroid [2] are less efficient on real-world applications. These datasets are useful for carefully spotting missing taint flows, but contain only a few dozen of applications.

Pauck *et al.* [25] used those three datasets to compare Amandroid [33], DIALDroid [5], DidFail [15], DroidSafe [12], FlowDroid [2] and IccTA [16] – all these tools will be also compared in this paper. To perform their comparison, they introduced the AQL (Android App Analysis Query Language) format. AQL can be used as a common language to describe the computed taint flow as well as the expected result for the datasets. It is interesting to notice that all the tested tools timed out at least once on real-world applications, and that Amandroid [33], DidFail [15], DroidSafe [12], IccTA [16] and ApkCombiner [17] (a tool used to combine applications) all failed to run on applications built for Android API 26. These results suggest that a more thorough study of the link between application characteristics (e.g. date, size) should be conducted. Luo *et al.* [22] used the framework introduced by Pauck *et al.* to compare Amandroid [33] and Flowdroid [2] on DroidBench and their own dataset TaintBench, composed of real-world android malware. They found out that those tools have a low recall on real-world malware, and are thus over adapted to micro-datasets. Unfortunately, because AQL is only focused on taint flows, we cannot use it to evaluate tools performing more generic analysis.

2.2 Static Analysis Tools Reusability

Several papers have reviewed Android analysis tools produced by researchers. Li *et al.* [18] published a systematic literature review for Android static analysis before May 2015. They analyzed 92 publications and classified them by goal, method used to solve the problem and underlying technical solution for handling the bytecode when performing the static analysis. In particular, they listed 27 approaches with an open-source implementation available. Nevertheless, experiments to evaluate the reusability of the pointed out software were not performed. We believe that the effort of reviewing the literature for making a comprehensive overview of available approaches should be pushed further: an existing published approach with a software that cannot be used for technical reasons endanger both the reproducibility and reusability of research.

A first work about quantifying the reusability of static analysis tools was proposed by Reaves *et al.* [28]. Seven Android analysis tools (Amandroid [33], AppAudit [35], DroidSafe [12], Epicc [24], FlowDroid [2], MalloDroid [9] and TaintDroid [8]) were selected to check if they were still readily usable. For each tool, both the usability and results of the tool were evaluated by asking auditors to install and use it on DroidBench and 16 real world applications. The auditors reported that most of the tools require a significant amount of time to setup, often due to dependencies issues and operating system incompatibilities. Reaves *et al.* propose to solve these issues by distributing a Virtual Machine with a functional build of the tool in addition to the source code. Regrettably, these Virtual Machines were not made available, preventing future researchers to take advantage of the work done by the auditors. Reaves *et al.* also report that real world applications are more challenging to analyze, with tools having lower results, taking more time and memory to run, sometimes to the point of not being able to run the analysis. We will confirm and expand this result in this paper with a larger dataset than only 16 real-world applications.

3 Methodology

3.1 Collecting Tools

We collected the static analysis tools from [18], plus one additional paper encountered during our review of the state-of-the-art (DidFail [15]). They are listed in Table 1, with the original release date and associated paper. We intentionally limited the collected tools to the ones selected by Li *et al.* [18] for several reasons. First, not using recent tools enables to have a gap of at least 5 years between the publication and the more recent APK files, which enables to measure the reusability of previous contributions with a reasonable gap of time. Second, collecting new tools would require to describe these tools in depth, similarly to what have been performed by Li *et al.* [18], which is not the primary goal of this paper. Additionally, selection criteria such as the publication venue or number of citations would be necessary to select a subset of tools, which would require an additional methodology. These possible contributions are left for future work.

Table 1. Considered tools [18]: availability and usage reliability

Tool	Availability			Repo type	Decision	Comments
	Bin	Src	Doc			
A3E [3] (2013)	–	●	●	github	×	Hybrid tool (static/dynamic)
A5 [32] (2014)	–	●	×	github	×	Hybrid tool (static/dynamic)
Adagio [10] (2013)	–	●	●	github	●	
Amandroid [33] (2014)	●	●	●	github	●	
Anadroid [19] (2013)	×	●	●	github	●	
Androguard [7] (2011)	–	●	●●	github	●	
Android-app-analysis [11] (2015)	×	●	●●	google	×	Hybrid tool (static/dynamic)
Apparecium [31] (2015)	×	●	×	github	●	
BlueSeal [30] (2014)	×	●	○	github	●	
Choi et al. [6] (2014)	×	●	○	github	×	Works on source files only
DIALDroid [5] (2017)	●	●	●	github	●	
DidFail [15] (2014)	●	●	○	bitbucket	●	
DroidSafe [12] (2015)	×	●	●	github	●	
Flowdroid [2] (2014)	×	●	●●	github	●	
Gator [29, 36] (2014), (2015)	×	●	●●	edu	●	
IC3 [23] (2015)	●	●	○	github	●	
IccTA [16] (2015)	●	●	●	github	●	
Lotrack [20] (2014)	×	×	×	github	○	Authors ack. a partial doc.
MalloDroid [9] (2012)	–	●	●	github	●	
PerfChecker [21] (2014)	×	×	○	request	●	Binary obtained from authors
Poeplau et al. [27] (2014)	×	○	×	github	×	Related to Android hardening
Redexer [14] (2012)	×	●	●	github	●	
SAAF [13] (2013)	×	●	●	github	●	
StadynA [37] (2015)	×	●	●	request	×	Hybrid tool (static/dynamic)
Thresher [4] (2013)	×	●	●	github	○	Not built with author’s help
Wognsen et al. [34] (2014)	–	●	×	bitbucket	●	

binaries, sources: –: not relevant, ●: available, ○: partially available, ×: not provided
documentation: ●●: excellent, MWE, ●: few inconsistencies, ○: bad quality, ×: not available
decision: ●: considered; ○: considered but not built; ×: out of scope of the study

Some tools use hybrid analysis (both static and dynamic): A3E [3], A5 [32], Android-app-analysis [11], StaDynA [37]. They have been excluded from this paper. We manually searched the tool repository when the website mentioned in the paper is no longer available (e.g. when the repository have been migrated from Google code to GitHub for example) and for each tool we searched for:

- an optional binary version of the tool that would be usable as a fall back (if the sources cannot be compiled for any reason);
- the source code of the tool;
- the documentation for building and using the tool with a MWE (Minimum Working Example).

In Table 1 we rated the quality of these artifacts with ”●” when available but may have inconsistencies, a ”○” when too much inconsistencies (inaccurate remarks about the sources, dead links or missing parts) have been found, a ”×

We finally excluded Choi *et al.* [6] as their tool works on the sources of Android applications, and Poeplau *et al.* [27] that focus on Android hardening.

Table 2. Selected tools, forks, selected commits and running environment

Tool	Origin	Alive Forks		Last commit	Authors	Environment
	Stars Alive	Nb	Usable	Date	Reached	Language - OS
Adagio [10]	74 ●	0	×	2022-11-17	●	Python - U20.04
Amandroid [33]	161 ×	2	×	2021-11-10	●	Scala - U22.04
Anadroid [19]	10 ×	0	×	2014-06-18	×	Scala/Java/Py. - U22.04
Androguard [7]	4430 ●	3	×	2023-02-01	×	Python - Python 3.11 slim
Apparecium [31]	0 ×	1	×	2014-11-07	×	Python - U22.04
BlueSeal [30]	0 ×	0	×	2018-07-04	●	Java - U14.04
DIALDroid [5]	16 ×	1	×	2018-04-17	×	Java - U18.04
DidFail [15]	4 ×			2015-06-17	●	Java/Python - U12.04
DroidSafe [12]	92 ×	3	×	2017-04-17	●	Java/Python - U14.04
Flowdroid [2]	868 ●	1	×	2023-05-07	●	Java - U22.04
Gator [29, 36]				2019-09-09	●	Java/Python - U22.04
IC3 [23]	32 ×	3	●	2022-12-06	×	Java - U12.04 / 22.04
IccTA [16]	83 ×	0	×	2016-02-21	●	Java - U22.04
Lotrack [20]	5 ×	2	×	2017-05-11	●	Java - ?
MalloDroid [9]	64 ×	10	×	2013-12-30	×	Python - U16.04
PerfChecker [21]				-	●	Java - U14.04
Redexer [14]	153 ×	0	×	2021-05-20	●	Ocaml/Ruby - U22.04
SAAF [13]	35 ×	5	×	2015-09-01	●	Java - U14.04
Thresher [4]	31 ×	1	×	2014-10-25	●	Java - U14.04
Wognsen <i>et al.</i> [34]			×	2022-06-27	×	Python/Prolog - U22.04

●: yes, ×: no, UX.04: Ubuntu X.04

As a summary, in the end we have 20 tools to compare. Some specificities should be noted. The IC3 tool will be duplicated in our experiments because two versions are available: the original version of the authors and a fork used by other tools like IccTa. For Androguard, the default task consists of unpacking the bytecode, the resources, and the Manifest. Cross-references are also built between methods and classes. Because such a task is relatively simple to perform, we decided to duplicate this tool and ask to Androguard to decompile an APK and create a control flow graph of the code using its decompiler: DAD. We refer to this variant of usage as `androguard_dad`. For Thresher and Lotrack, because these tools cannot be built, we excluded them from experiments.

Finally, starting with 26 tools of Table 1, with the two variations of IC3 and Androguard, we have in total 22 static analysis tools to evaluate in which two tools cannot be built and will be considered as always failing.

3.2 Source Code Selection and Building Process

In a second step, we explored the best sources to be selected among the possible forks of a tool. We reported some indicators about the explored forks and our decision about the selected one in Table 2. For each source code repository called "Origin", we reported in Table 2 the number of GitHub stars attributed by users and we mentioned if the project is still alive (● in column Alive when a commit exist in the last two years). Then, we analyzed the fork tree of the project. We searched recursively if any forked repository contains a more recent commit than the last one of the branch mentioned in the documentation of the original repository. If such a commit is found (number of such commits are reported in

column Alive Forks Nb), we manually looked at the reasons behind this commit and considered if we should prefer this more up-to-date repository instead of the original one (column "Alive Forks Usable"). As reported in Table 2, we excluded all forks, except IC3 for which we selected the fork JordanSamhi/ic3, because they always contain experimental code with no guarantee of stability. For example, a fork of Aparecium contains a port for Windows 7 which does not suggest an improvement of the stability of the tool. For IC3, the fork seems promising: it has been updated to be usable on a recent operating system (Ubuntu 22.04 instead of Ubuntu 12.04 for the original version) and is used as a dependency by IccTa. We decided to keep these two versions of the tool (IC3 and IC3_fork) to compare their results.

Then, we self-allocated a maximum of four days for each tool to successfully read and follow the documentation, compile the tool and obtain the expected result when executing an analysis of a MWE. We sent an email to the authors of each tool to confirm that we used the more suitable version of the code, that the command line we used to analyze an application is the most suitable one and, in some cases, requested some help to solve issues in the building process. We reported in Table 2 the authors that answered our request and confirmed our decisions.

From this building phase, several observations can be made. Using a recent operating system, it is almost impossible in a reasonable amount of time to rebuild a tool released years ago. Too many dependencies, even for Java based programs, trigger compilation or execution problems. Thus, if the documentation mentions a specific operating system, we use a Docker image of this OS. Most of the time, tools require additional external components to be fully functional. It could be resources such as the android.jar file for each version of the SDK, a database, additional libraries or tools. Depending of the quality of the documentation, setting up those components can take hours to days. This is why we automatized in a Dockerfile the setup of the environment in which the tool is built and run³.

3.3 Runtime Conditions

As shown in Figure 1, before benchmarking the tools, we built and installed them in a Docker containers for facilitating any reuse of other researchers. We converted them into Singularity containers because we had access to such a cluster and because this technology is often used by the HPC community for ensuring the reproducibility of experiments. We performed manual tests using these Singularity images to check:

³ To guarantee reproducibility we published the results, datasets, Dockerfiles and containers. Source code is located at: <https://github.com/histausse/rasta>, datasets and experiments results are available at: <https://zenodo.org/records/10144014>, Singularity containers are available at: <https://zenodo.org/records/10980349>, Docker images (`histausse/rasta-<toolname>:icsr2024`) can be downloaded from [Docker Hub](#).

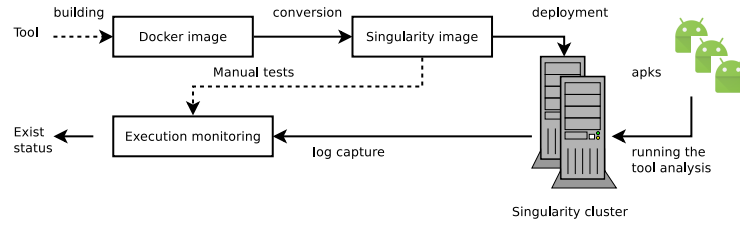


Fig. 1. Methodology overview

- the location where the tool is writing on the disk. For the best performances, we expect the tools to write on a mount point backed by an SSD. Some tools may write data at unexpected locations which required small patches from us.
- the amount of memory allocated to the tool. We checked that the tool could run a MWE with a 64 GB limit of RAM.
- the network connection opened by the tool, if any. We expect the tool not to perform any network operation such as the download of Android SDKs. Thus, we prepared the required files and cached them in the images during the building phase. In a few cases, we patched the tool to disable the download of resources.

A campaign of tests consists in executing the 20 selected tools on all APKs of a dataset. The constraints applied on the clusters are:

- No network connection is authorized in order to limit any execution of malicious software.
- The allocated RAM for a task is 64 GB.
- The allocated maximum time is 1 hour.
- The allocated object space / stack space is 64 GB / 16 GB if the tool is a Java based program.

For the disk files, we use a mount point that is stored on a SSD disk, with no particular limit of size. Note that, because the allocation of 64 GB could be insufficient for some tool, we evaluated the results of the tools on 20% of our dataset (described later in Section 3.4) with 128 GB of RAM and 64 GB of RAM and checked that the results were similar. With this confirmation, we continued our evaluations with 64 GB of RAM only.

3.4 Dataset

We built a dataset named **Rasta** to cover all dates between 2010 to 2023. This dataset is a random extract of Androzoo [1], for which we balanced applications between years and size. For each year and inter-decile range of size in Androzoo, 500 applications have been extracted with an arbitrary proportion of 7% of malware. This ratio has been chosen because it is the ratio of goodware/malware

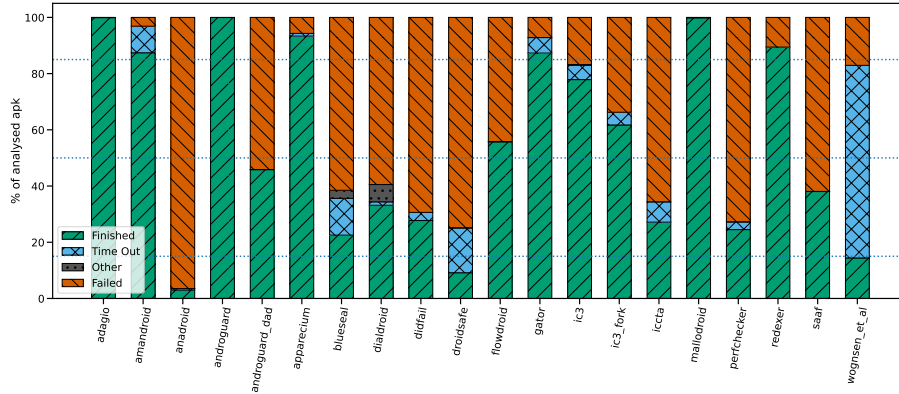


Fig. 2. Exit status for the Rasta dataset

that we observed when performing a raw extract of Androzoo. For checking the maliciousness of an Android application we rely on the VirusTotal detection indicators. If more than 5 antiviruses have flagged the application as malicious, we consider it as a malware. If no antivirus has reported the application as malicious, we consider it as a goodware. Applications in between are dropped.

For computing the release date of an application, we contacted the authors of Androzoo to compute the minimum date between the submission to Androzoo and the first upload to VirusTotal. Such a computation is more reliable than using the DEX date that is often obfuscated when packaging the application.

4 Experiments

4.1 RQ1: Re-Usability Evaluation

Figure 2 represents the success/failure rate (green/orange) of the tools. We distinguished failure to compute a result from timeout (blue) and crashes of our evaluation framework (in grey, probably due to out of memory kills of the container itself). Because it may be caused by a bug in our own analysis stack, exit status represented in grey (Other) are considered as unknown errors and not as failure of the tool.

We observe a global increase of the failed status: 12 tools (54.5%) have a finishing rate below 50%. Three tools (androguard_dad, blueseal, saaf) reach the bar of 50% of failure. 7 tools keep a high success rate: Adagio, Amandroid, Androguard, Apparecium, Gator, Mallodroid, Redexer. Regarding IC3, the fork with a simpler build process and support for modern OS has a lower success rate than the original tool.

Two tools should be discussed in particular. Androguard has a high success rate which is not surprising: it used by a lot of tools, including for analyzing

application uploaded to the Androzo repository. Nevertheless, when using Androguard decompiler (DAD) to decompile an APK, it fails more than 50% of the time. This example shows that even a tool that is frequently used can still run into critical failures. Concerning Flowdroid, our results show a very low timeout rate (0.06%) which was unexpected: in our exchanges, Flowdroid’s author were expecting a higher rate of timeout and fewer crashes.

As a summary, the final ratio of successful analysis for the tools that we could run is 54.9%. When including the two defective tools, this ratio drops to 49.9%.

RQ1 answer: On a recent dataset we consider that 54.5% of the tools are unusable. For the tools that we could run, 54.9% of analysis are finishing successfully.

4.2 RQ2: Size, SDK and Date Influence

To measure the influence of the date, SDK version and size of applications, we fixed one parameter while varying an other. For the sake of clarity, we separated Java based / non Java based tools.

Fixed application year. (5000 APKs) We selected the year 2022 which has a good amount of representatives for each decile of size in our application dataset. Figure 3 (resp. 4) shows the finishing rate of the tools in function of the size of the bytecode for Java based tools (resp. non Java based tools) analyzing applications of 2022. We can observe that all Java based tools have a finishing rate decreasing over years. 50% of non Java based tools have the same behavior.

Fixed application bytecode size. (6252 APKs) We selected the sixth decile (between 4.08 and 5.20 MB), which is well represented in a wide number of years. Figure 5 and 6 represent the finishing rate depending of the year at a fixed bytecode size. We observe that 9 tools over 12 have a finishing rate dropping below 20% for Java based tools, which is not the case for non Java based tools.

We performed similar experiments by variating the min SDK and target SDK versions, still with a fixed bytecode size between 4.08 and 5.2 MB, as shown in Figure 7 and 8. We found that contrary to the target SDK, the min SDK version has an impact on the finishing rate of Java based tools: 8 tools over 12 are below 50% after SDK 16. It is not surprising, as the min SDK is highly correlated to the year.

RQ2 answer: The success rate varies based on the size of bytecode and SDK version. The date is correlated with the success rate for Java based tools only.

4.3 RQ3: Malware vs Goodware

We compared the finishing rate of malware and goodware applications for evaluated tools. Because, the size of applications impacts this finishing rate, it is interesting to compare the success rate for each decile of bytecode size. Table 3

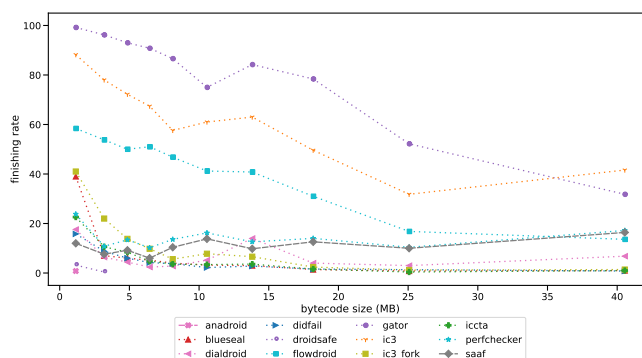


Fig. 3. Finishing rate by bytecode size for APK detected in 2022 – Java based tools

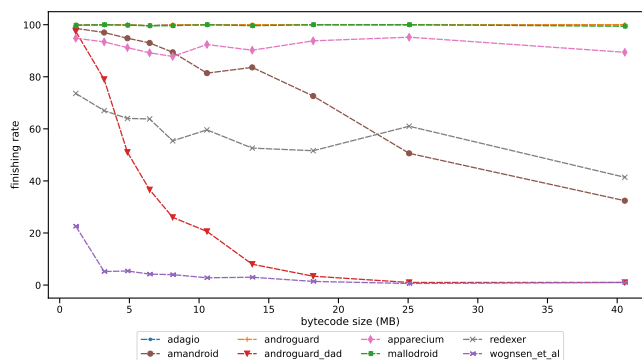


Fig. 4. Finishing rate by bytecode size for APK detected in 2022 – Non Java based tools

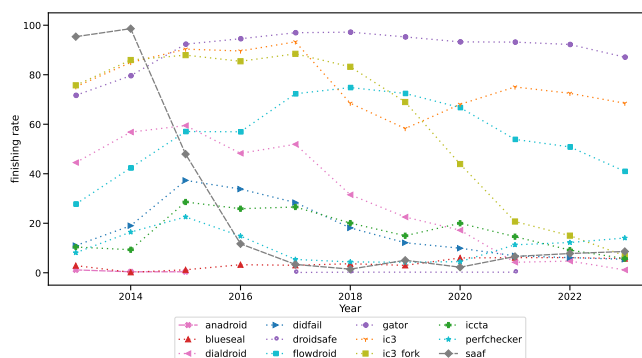


Fig. 5. Finishing rate by discovery year with a bytecode size $\in [4.08, 5.2]$ MB – Java based tools

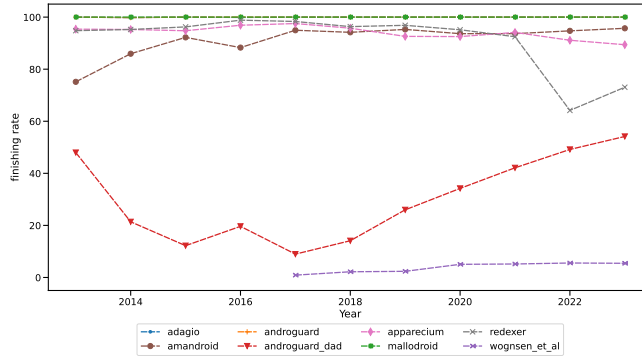


Fig. 6. Finishing rate by discovery year with a bytecode size $\in [4.08, 5.2]$ MB – Non Java based tools

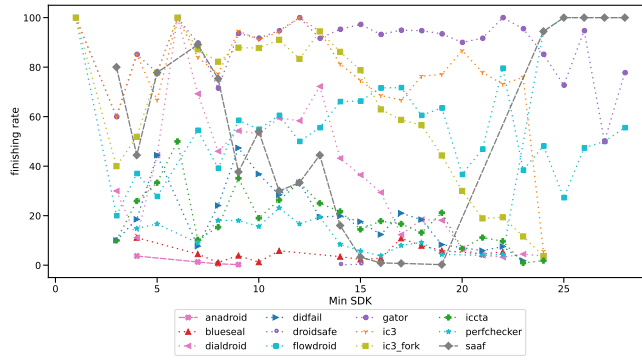


Fig. 7. Finishing rate by min SDK with a bytecode size $\in [4.08, 5.2]$ MB – Java based tools

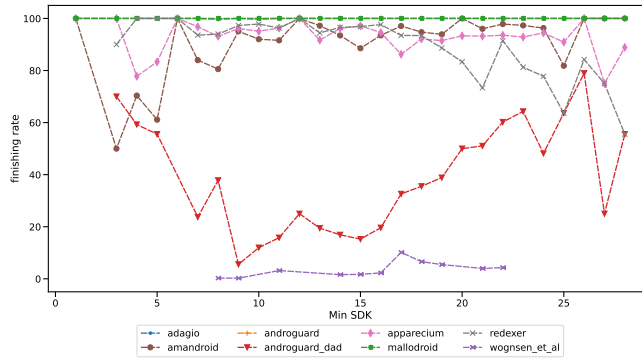


Fig. 8. Finishing rate by min SDK with a bytecode size $\in [4.08, 5.2]$ MB – Non Java based tools

Table 3. DEX size and Finishing Rate (FR) per decile

Decile	Avg DEX size MB		Finishing Rate: FR		Ratio Size	Ratio FR
	Good	Mal	Good	Mal	Good/Mal	Good/Mal
1	0.13	0.11	0.85	0.82	1.17	1.04
2	0.54	0.55	0.74	0.72	0.97	1.03
3	1.37	1.25	0.63	0.66	1.09	0.97
4	2.41	2.34	0.57	0.62	1.03	0.92
5	3.56	3.55	0.53	0.59	1.00	0.90
6	4.61	4.56	0.50	0.61	1.01	0.82
7	5.87	5.91	0.47	0.57	0.99	0.83
8	7.64	7.63	0.43	0.56	1.00	0.76
9	11.39	11.26	0.39	0.58	1.01	0.67
10	24.24	21.36	0.33	0.46	1.13	0.73

reports the bytecode size and the finishing rate of goodware and malware in each decile of size. We also computed the ratio of the bytecode size and finishing rate for the two populations. We observe that the ratio for the finishing rate decreases from 1.04 to 0.73, while the ratio of the bytecode size is around 1. We conclude from this table that analyzing malware triggers less errors than for goodware.

RQ3 answer: Analyzing malware applications triggers less errors for static analysis tools than analyzing goodware for comparable bytecode size.

5 Discussion

5.1 State-of-the-art comparison

Our findings are consistent with the numerical results of Pauck *et al.* that showed that 58.9% of DIALDroid-Bench [5] real-world applications are analyzed successfully with the 6 evaluated tools [25]. Six years after the release of DIALDroid-Bench, we obtain a lower ratio of 40.1% for the same set of 6 tools but using the Rasta dataset of 62 525 applications. We extended this result to a set of 20 tools and obtained a global success rate of 54.9%. We confirmed that most tools require a significant amount of work to get them running [28].

Investigating the reason behind tools' errors is a difficult task and will be investigated in a future work. For now, our manual investigations show that the nature of errors varies from one analysis to another, without any easy solution for the end user for fixing it.

5.2 Recommendations

Finally, we summarize some takeaways that developers should follow to improve the success of reusing their developed software.

For improving the reliability of their software, developers should use classical development best practices, for example continuous integration, testing, code review. For improving the reusability developers should write a documentation

about the tool usage and provide a minimal working example and describe the expected results. Interactions with the running environment should be minimized, for example by using a docker container, a virtual environment or even a virtual machine. Additionally, a small dataset should be provided for a more extensive test campaign and the publishing of the expected result on this dataset would ensure to be able to evaluate the reproducibility of experiments.

Finally, an important remark concerns the libraries used by a tool. We have seen two types of libraries: a) internal libraries manipulating internal data of the tool; b) external libraries that are used to manipulate the input data (APKs, bytecode, resources). We observed by our manual investigations that external libraries are the ones leading to crashes because of variations in recent APKs (file format, unknown bytecode instructions, multi-DEX files). We believe that the developer should provide enough documentation to make possible a later upgrade of these external libraries.

5.3 Threats to validity

Our application dataset is biased in favor of Androguard, because Androzoo have already used Androguard internally when collecting applications and discarded any application that cannot be processed with this tool.

Despite our best efforts, it is possible that we made mistakes when building or using the tools. It is also possible that we wrongly classified a result as a failure. To mitigate this possible problem we contacted the authors of the tools to confirm that we used the right parameters and chose a valid failure criterion.

The timeout value, amount of memory are arbitrarily fixed. For mitigating their effect, a small extract of our dataset has been analyzed with more memory/time for measuring any difference.

Finally, the use of VirusTotal for determining if an application is a malware or not may be wrong. For limiting this impact, we used a threshold of at most 5 antiviruses (resp. no more than 0) reporting an application as being a malware (resp. goodware) for taking a decision about maliciousness (resp. benignness).

6 Conclusion

This paper has assessed the suggested results of the literature [22, 25, 28] about the reliability of static analysis tools for Android applications. With a dataset of 62 525 applications we established that 54.5% of 22 tools are not reusable, when considering that a tool that has more than 50% of time a failure is unusable. In total, the analysis success rate of the tools that we could run for the entire dataset is 54.9%. The characteristics that have the most influence on the success rate is the bytecode size and min SDK version. Finally, we showed that malware APKs have a better finishing rate than goodware.

In future works, we plan to investigate deeper the reported errors of the tools in order to analyze the most common types of errors, in particular for Java based tools. We also plan to extend this work with a selection of more recent tools performing static analysis.

Acknowledgment

This work was supported by the ANR Research under the Plan France 2030 bearing the reference ANR-22-PECY-0007.

References

1. Allix, K., Bissyandé, T.F., Klein, J., and Traon, Y.L.: AndroZoo: Collecting Millions of Android Apps for the Research Community. In: 13th Working Conference on Mining Software Repositories (MSR), pp. 468–471 (2016)
2. Arzt, S., Rasthofer, S., Fritz, C., Bodden, E., Bartel, A., Klein, J., Le Traon, Y., Octeau, D., and McDaniel, P.: FlowDroid: Precise Context, Flow, Field, Object-sensitive and Lifecycle-aware Taint Analysis for Android Apps. In: ACM SIGPLAN Conference on Programming Language Design and Implementation, pp. 259–269. ACM Press, Edinburgh, UK (2014). DOI: [10.1145/2666356.2594299](https://doi.org/10.1145/2666356.2594299)
3. Azim, T., and Neamtiu, I.: Targeted and Depth-First Exploration for Systematic Testing of Android Apps. In: Hosking, A.L., Eugster, P.T., and Lopes, C.V. (eds.) Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications, OOPSLA 2013, Part of SPLASH 2013, Indianapolis, IN, USA, October 26-31, 2013, pp. 641–660. ACM (2013). DOI: [10.1145/2509136.2509549](https://doi.org/10.1145/2509136.2509549)
4. Blackshear, S., Chang, B.-Y.E., and Sridharan, M.: Thresher: Precise Refutations for Heap Reachability. SIGPLAN Not. 48(6), 275–286 (2013). DOI: [10.1145/2499370.2462186](https://doi.org/10.1145/2499370.2462186)
5. Bosu, A., Liu, F., Yao, D., and Wang, G.: Collusive Data Leak and More: Large-scale Threat Analysis of Inter-app Communications. In: Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, pp. 71–85. ACM, Abu Dhabi United Arab Emirates (2017). DOI: [10.1145/3052973.3053004](https://doi.org/10.1145/3052973.3053004)
6. Choi, K., and Chang, B.-M.: A Type and Effect System for Activation Flow of Components in Android Programs. Information Processing Letters 114(11), 620–627 (2014). DOI: [10.1016/j.ipl.2014.05.011](https://doi.org/10.1016/j.ipl.2014.05.011)
7. Desnos, A., and Gueguen, G.: Android: From Reversing to Decompilation. Black Hat Abu Dhabi (2011)
8. Enck, W., Gilbert, P., Chun, B.-G., Cox, L.P., Jung, J., McDaniel, P., and Sheth, A.N.: TaintDroid: An Information-Flow Tracking System for Realtime Privacy Monitoring on Smartphones. In: 9th USENIX Symposium on Operating Systems Design and Implementation, pp. 393–407. USENIX Association, Vancouver, BC, Canada (2010)
9. Fahl, S., Harbach, M., Muders, T., Baumgärtner, L., Freisleben, B., and Smith, M.: Why Eve and Mallory Love Android: An Analysis of Android SSL (in)Security. In: Proceedings of the 2012 ACM Conference on Computer and Communications Security, pp. 50–61. ACM, Raleigh North Carolina USA (2012). DOI: [10.1145/2382196.2382205](https://doi.org/10.1145/2382196.2382205)
10. Gascon, H., Yamaguchi, F., Arp, D., and Rieck, K.: Structural Detection of Android Malware Using Embedded Call Graphs. In: Proceedings of the 2013 ACM Workshop on Artificial Intelligence and Security, pp. 45–54. ACM, Berlin Germany (2013). DOI: [10.1145/2517312.2517315](https://doi.org/10.1145/2517312.2517315)

11. Geneiatakis, D., Fovino, I.N., Kounelis, I., and Stirparo, P.: A Permission Verification Approach for Android Mobile Applications. *Computers & Security* 49, 192–205 (2015). DOI: [10.1016/j.cose.2014.10.005](https://doi.org/10.1016/j.cose.2014.10.005)
12. Gordon, M.I., Kim, D., Perkins, J.H., Gilham, L., Nguyen, N., and Rinard, M.C.: Information Flow Analysis of Android Applications in DroidSafe. In: 22nd Annual Network and Distributed System Security Symposium, NDSS 2015, San Diego, California, USA, February 8-11, 2015. The Internet Society (2015)
13. Hoffmann, J., Ussath, M., Holz, T., and Spreitzenbarth, M.: Slicing Droids: Program Slicing for Smali Code. In: Proceedings of the 28th Annual ACM Symposium on Applied Computing, SAC '13, pp. 1844–1851. Association for Computing Machinery, New York, NY, USA (2013). DOI: [10.1145/2480362.2480706](https://doi.org/10.1145/2480362.2480706)
14. Jeon, J., Micinski, K.K., Vaughan, J.A., Fogel, A., Reddy, N., Foster, J.S., and Millstein, T.: Dr. Android and Mr. Hide: Fine-Grained Permissions in Android Applications. In: Proceedings of the Second ACM Workshop on Security and Privacy in Smartphones and Mobile Devices, pp. 3–14. ACM, Raleigh North Carolina USA (2012). DOI: [10.1145/2381934.2381938](https://doi.org/10.1145/2381934.2381938)
15. Klieber, W., Flynn, L., Bhosale, A., Jia, L., and Bauer, L.: Android Taint Flow Analysis for App Sets. In: Proceedings of the 3rd ACM SIGPLAN International Workshop on the State of the Art in Java Program Analysis, pp. 1–6. ACM, Edinburgh United Kingdom (2014). DOI: [10.1145/2614628.2614633](https://doi.org/10.1145/2614628.2614633)
16. Li, L., Bartel, A., Bissyande, T.F., Klein, J., Le Traon, Y., Arzt, S., Rasthofer, S., Bodden, E., Octeau, D., and McDaniel, P.: IccTA: Detecting Inter-Component Privacy Leaks in Android Apps. In: 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering, pp. 280–291. IEEE, Florence, Italy (2015). DOI: [10.1109/ICSE.2015.48](https://doi.org/10.1109/ICSE.2015.48)
17. Li, L., Bartel, A., Bissyandé, T.F., Klein, J., and Traon, Y.L.: ApkCombiner: Combining Multiple Android Apps to Support Inter-App Analysis. In: Federrath, H., and Gollmann, D. (eds.) *ICT Systems Security and Privacy Protection*, pp. 513–527. Springer International Publishing, Cham (2015). DOI: [10.1007/978-3-319-18467-8_34](https://doi.org/10.1007/978-3-319-18467-8_34)
18. Li, L., Bissyandé, T.F., Papadakis, M., Rasthofer, S., Bartel, A., Octeau, D., Klein, J., and Le Traon, Y.: Static Analysis of Android Apps: A Systematic Literature Review. *Information and Software Technology* 88, 67–95 (2017). DOI: [10.1016/j.infsof.2017.04.001](https://doi.org/10.1016/j.infsof.2017.04.001)
19. Liang, S., Keep, A.W., Might, M., Lyde, S., Gilray, T., Aldous, P., and Van Horn, D.: Sound and Precise Malware Analysis for Android via Pushdown Reachability and Entry-Point Saturation. In: Proceedings of the Third ACM Workshop on Security and Privacy in Smartphones & Mobile Devices. SPSM '13, pp. 21–32. Association for Computing Machinery, New York, NY, USA (2013). DOI: [10.1145/2516760.2516769](https://doi.org/10.1145/2516760.2516769)
20. Lillack, M., Kästner, C., and Bodden, E.: Tracking Load-Time Configuration Options. In: Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering. ASE '14, pp. 445–456. Association for Computing Machinery, New York, NY, USA (2014). DOI: [10.1145/2642937.2643001](https://doi.org/10.1145/2642937.2643001)
21. Liu, Y., Xu, C., and Cheung, S.-C.: Characterizing and Detecting Performance Bugs for Smartphone Applications. In: Proceedings of the 36th International Conference on Software Engineering, pp. 1013–1024. ACM, Hyderabad India (2014). DOI: [10.1145/2568225.2568229](https://doi.org/10.1145/2568225.2568229)
22. Luo, L., Pauck, F., Piskachev, G., Benz, M., Pashchenko, I., Mory, M., Bodden, E., Hermann, B., and Massacci, F.: TaintBench: Automatic Real-World Malware

- Benchmarking of Android Taint Analyses. *Empir Software Eng* 27(1), 16 (2022). DOI: [10.1007/s10664-021-10013-5](https://doi.org/10.1007/s10664-021-10013-5)
23. Oceau, D., Luchaup, D., Dering, M., Jha, S., and McDaniel, P.: Composite Constant Propagation: Application to Android Inter-Component Communication Analysis. In: 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering, pp. 77–88. IEEE, Florence, Italy (2015). DOI: [10.1109/ICSE.2015.30](https://doi.org/10.1109/ICSE.2015.30)
 24. Oceau, D., McDaniel, P., Jha, S., Bartel, A., Bodden, E., Klein, J., and Le Traon, Y.: Effective Inter-Component communication mapping in android: An essential step towards holistic security analysis. In: 22nd USENIX Security Symposium (USENIX Security 13), pp. 543–558 (2013)
 25. Pauck, F., Bodden, E., and Wehrheim, H.: Do Android Taint Analysis Tools Keep Their Promises? In: Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, pp. 331–341. ACM, Lake Buena Vista FL USA (2018). DOI: [10.1145/3236024.3236029](https://doi.org/10.1145/3236024.3236029)
 26. Pendlebury, F., Pierazzi, F., Jordaney, R., Kinder, J., and Cavallaro, L.: TESSERACT: Eliminating Experimental Bias in Malware Classification across Space and Time. (2018)
 27. Poeplau, S., Fratantonio, Y., Bianchi, A., Kruegel, C., and Vigna, G.: Execute This! Analyzing Unsafe and Malicious Dynamic Code Loading in Android Applications. In: 21st Annual Network and Distributed System Security Symposium, NDSS 2014, San Diego, California, USA, February 23-26, 2014. The Internet Society (2014)
 28. Reaves, B., Bowers, J., Gorski III, S.A., Anise, O., Bobhate, R., Cho, R., Das, H., Hussain, S., Karachiwala, H., Scaife, N., Wright, B., Butler, K., Enck, W., and Traynor, P.: *droid: Assessment and Evaluation of Android Application Analysis Tools. *ACM Computing Surveys* 49(3), 55:1–55:30 (2016). DOI: [10.1145/2996358](https://doi.org/10.1145/2996358)
 29. Rountev, A., and Yan, D.: Static Reference Analysis for GUI Objects in Android Software. In: Proceedings of Annual IEEE/ACM International Symposium on Code Generation and Optimization, pp. 143–153. ACM, Orlando FL USA (2014). DOI: [10.1145/2544137.2544159](https://doi.org/10.1145/2544137.2544159)
 30. Shen, F., Vishnubhotla, N., Todarka, C., Arora, M., Dhandapani, B., Lehner, E.J., Ko, S.Y., and Ziarek, L.: Information Flows as a Permission Mechanism. In: Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering, pp. 515–526. ACM, Vasteras Sweden (2014). DOI: [10.1145/2642937.2643018](https://doi.org/10.1145/2642937.2643018)
 31. Titze, D., and Schutte, J.: Apparecium: Revealing Data Flows in Android Applications. In: 2015 IEEE 29th International Conference on Advanced Information Networking and Applications, pp. 579–586. IEEE, Gwangju, South Korea (2015). DOI: [10.1109/AINA.2015.239](https://doi.org/10.1109/AINA.2015.239)
 32. Vidas, T., Tan, J., Nahata, J., Tan, C.L., Christin, N., and Tague, P.: A5: Automated Analysis of Adversarial Android Applications. In: Proceedings of the 4th ACM Workshop on Security and Privacy in Smartphones & Mobile Devices, pp. 39–50. ACM, Scottsdale Arizona USA (2014). DOI: [10.1145/2666620.2666630](https://doi.org/10.1145/2666620.2666630)
 33. Wei, F., Roy, S., Ou, X., and Robby: Amandroid: A Precise and General Inter-component Data Flow Analysis Framework for Security Vetting of Android Apps. In: ACM SIGSAC Conference on Computer and Communications Security, pp. 1329–1341. ACM, Scottsdale Arizona USA (2014). DOI: [10.1145/2660267.2660357](https://doi.org/10.1145/2660267.2660357)
 34. Wognsen, E.R., Karlsen, H.S., Olesen, M.C., and Hansen, R.R.: Formalisation and Analysis of Dalvik Bytecode. *Science of Computer Programming* 92, 25–55 (2014). DOI: [10.1016/j.scico.2013.11.037](https://doi.org/10.1016/j.scico.2013.11.037)

35. Xia, M., Gong, L., Lyu, Y., Qi, Z., and Liu, X.: Effective Real-Time Android Application Auditing. In: 2015 IEEE Symposium on Security and Privacy, pp. 899–914. IEEE, San Jose, CA (2015). DOI: [10.1109/SP.2015.60](https://doi.org/10.1109/SP.2015.60)
36. Yang, S., Yan, D., Wu, H., Wang, Y., and Rountev, A.: Static Control-Flow Analysis of User-Driven Callbacks in Android Applications. In: 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering, pp. 89–99. IEEE, Florence, Italy (2015). DOI: [10.1109/ICSE.2015.31](https://doi.org/10.1109/ICSE.2015.31)
37. Zhauniarovich, Y., Ahmad, M., Gadyatskaya, O., Crispo, B., and Massacci, F.: StaDynA: Addressing the Problem of Dynamic Code Updates in the Security Analysis of Android Applications. In: Proceedings of the 5th ACM Conference on Data and Application Security and Privacy, pp. 37–48. ACM, San Antonio Texas USA (2015). DOI: [10.1145/2699026.2699105](https://doi.org/10.1145/2699026.2699105)