

Cerberus: a Program Repair Framework

Ridwan Shariffdeen*, Martin Mirchev[†], Yannic Noller[‡], Abhik Roychoudhury[§]

National University of Singapore, Singapore

Email: *ridwan@comp.nus.edu.sg, [†]mmirchev@comp.nus.edu.sg, [‡]yannic.noller@acm.org, [§]abhik@comp.nus.edu.sg

Abstract—Automated Program Repair (APR) represents a suite of emerging technologies which attempt to automatically fix bugs and vulnerabilities in programs. APR is a rapidly growing field with new tools and benchmarks being added frequently. Yet a language agnostic repair framework is not available. We introduce CERBERUS, a program repair framework integrated with 13 program repair tools and 6 repair benchmarks, coexisting in the same framework. CERBERUS is capable of executing diverse set of program repair tasks, using multitude of program repair tools and benchmarks.

Video: <https://www.youtube.com/watch?v=bYtShpsGL68>

Index Terms—automated program repair, repair platform

I. INTRODUCTION

Automated Program Repair (APR) [1] has many applications in software engineering, like supporting software developers in fixing bugs, particularly for security vulnerabilities and concurrency bugs, but also for guiding students to solve programming assignments in an educational context. Furthermore, APR has already been adopted to some initial extent in industry deployments [2], [3], in which they are usually added to the CI/CD pipelines, proposing patches for failing test cases. While we can observe a plethora of APR approaches [4], they greatly vary in required patch ingredients, target languages, and execution environments, including the dependencies and instrumentation requirements. This large diversity poses a challenge for bug and patch reproduction and technique comparisons.

The existing approaches for integrating APR tools into frameworks [5]–[7] fail to provide an environment and architecture that would allow covering multiple application domains and target languages. REPAIRTHEMALL [5] assumes that APR tools are provided as .jar-files and is customized for Java repair. SECURETHEMALL [6] is customized for security vulnerability repair and targets C/C++ programs. The most recent proposed work is MAESTRO [7], a benchmarking framework for automated program repair tools that supports multiple implementation and target languages. However, their design choices prevent the integration of semantic-based techniques like CPR [8] that require complex build infrastructures.

To close this gap, we present CERBERUS, a program repair framework that provides the means to integrate many different APR tools with diverse target languages and application domains, their execution environments, and their experiment data sets, resulting in a unified way of accessing the developed tools. In contrast to the existing works, CERBERUS does not make any assumption about implementation or target language, and is not customized to a specific application

domain. In contrast to MAESTRO, it encapsulates the benchmark and the repair tool in a single container setup, which makes it straightforward to integrate tools with complex build infrastructure. In fact, we have already integrated 13 tools for C/C++ and Java with various repair methodologies covering search-based, semantic-based, and learning-based APR and multiple application domains, including test-based general-purpose repair, security repair, static-based concurrency repair, and student assignment repair.

CERBERUS is useful for software engineering researchers as well as for software developers. Researchers can integrate their new APR tool into our framework and perform evaluations with the already integrated tools. Thereby, they do not need to consider dependency issues or the technical setups of other tools because our framework makes them readily available. Moreover, we have already integrated the corresponding benchmarks and data sets, which makes it straightforward to run additional experiments. Software developers, who may not be familiar with APR and the existing tools, can use CERBERUS to apply different APR tools on their own (private) data set to see which technique is most suited for their needs. CERBERUS makes the existing APR tools accessible beyond their original experimental environment (i.e., as reported in the corresponding research papers), enables the reproducibility of APR studies, allows comparisons between tools, and enables practitioners to get easy access to the state of the art in APR.

To demonstrate the capabilities of CERBERUS in handling a diverse set of APR approaches, we used it to reproduce the experiments of VERIFIX [9] and of SEQUENCER [10]. While SEQUENCER is a learning-based APR approach, VERIFIX applies repair in the educational context, which makes it different from the standard general-purpose repair application. Our results show that CERBERUS can produce similar or the same results as the original works. We observed minor differences in the results for VERIFIX because we used the latest tool version, which the authors had improved since the original publication. We make the following contributions:

- CERBERUS, a fully agnostic repair platform with a layered architecture that allows addition of new tools and benchmarks, including complex build infrastructures,
- the integration of 13 program repair tools and 6 repair benchmarks across multiple target languages and application domains, and
- the demonstration of CERBERUS’s capabilities on executing repair in the educational context with VERIFIX and in general-purpose repair with the learning-based APR technique SEQUENCER.

II. DESIGN AND USAGE

In this section, we describe the design architecture and the components of our platform. The default mode of execution in CERBERUS is using containers, which allows to create isolated, modular, and easily reproducible experiments for empirical studies. However, CERBERUS is designed to cater repair tools in both containerized and non-containerized environments. This is because not all program repair tools are available as a Docker container, but also can be made available as a virtual machine (e.g., SENX [11]). The platform is built with the aim of providing flexibility in executing repair tools with less assumptions about the environment and the dependencies required. For brevity, the rest of the paper discusses the containerized mode of CERBERUS. Instructions for specifying the virtualization can be found in our repository.

Experiments in program repair require two main components to be configured and set up. One component is the APR tool itself, with all dependencies available during runtime. The second component is the benchmark which provides information on the bug that needs to be repaired and the mechanism to reproduce the bug in a new environment. CERBERUS abstracts the nuances of different experiments and alleviates the repetitive, tedious efforts required to set up an experiment by providing a single interface to the user.

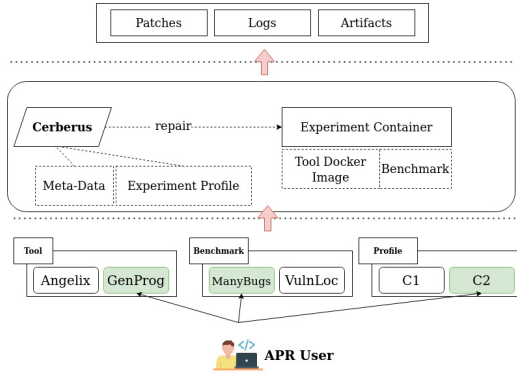


Fig. 1. Repair workflow of CERBERUS

Figure 1 illustrates the workflow for CERBERUS. The user selects a repair tool and a repair benchmark from a pre-configured list. Additionally, we provide experiment profiles that can be configured to control the experiment for different parameters (i.e., time duration for repair, number of test cases provided, etc.). Once the user makes the selection, CERBERUS will extract the relevant meta-data of the bug(s) that needs to be repaired. First, it will load the Docker image for the repair tool as the baseline image and extend the container by setting up the benchmark. Once the container is instantiated, CERBERUS will load the experiment profile to adjust the necessary parameters and invoke the repair module inside the container. Finally, CERBERUS extracts necessary artifacts, e.g., generated patches, debug logs, and other artifacts like repair constraints and additional generated test cases.

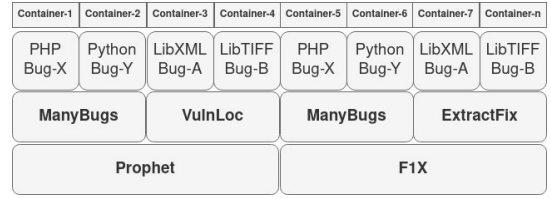


Fig. 2. Layered architecture for containers in CERBERUS

Figure 2 depicts the layered architecture CERBERUS follows to create containers for each experiment. As indicated previously, the baseline image is the repair tool that encapsulates all necessary dependencies to run the repair tool. CERBERUS would then extend the container by setting up the benchmark subject. For each experiment, a container will be spawned and can be used for a controlled experiment. This layered architecture provides several advantages, such as efficient space management, low latency for repeated experiments, and the re-usability of shared layers. Docker uses a union file system that uses a copy-on-write strategy to provide efficient space management. This means only the files modified by the write-layer (the top-most layer) are changed in the container. This strategy allows sharing of common files across different containers. Additionally, since docker provides in-built caching of the layers, the consecutive re-run of experiments can be executed without rebuilding the complete container. For instance, consider a subject in the MANYBUGS benchmark (e.g., PHP-Bug-X) that we will run for one hour on PROPHET. Re-running the same experiment with a 2-hour timeout would be easily instantiated with the minor change to the repair profile, saving time and space for the new experiment.

A. Extendibility

CERBERUS provides necessary abstractions for tool developers and benchmark providers to easily integrate and extend the platform. To integrate a new repair tool, the tool developer needs to create a new driver for the tool. To invoke the repair process, the driver should provide the basic functionalities for CERBERUS. In particular, the driver should transform the meta-data provided by a benchmark into the expected form by the repair tool (i.e., creating configuration files). Once a repair driver is configured, CERBERUS can create the experiment container and execute the experiment with the configuration parameters defined in the repair profile and the information provided by the selected benchmark.

Integrating a new benchmark requires a separate driver and a schema file specifying the necessary meta-data for the defects in the benchmark. Different benchmarks provide different sets of information based on the complexity of the defects and the artifacts required for repair (i.e., test cases). For instance, VULNLOC only provides one failing test case since the benchmark is for vulnerability repair, while MANYBUGS includes a list of passing and failing test cases for each bug. The driver should provide the functionality to

TABLE I
DETAILS OF THE PROGRAM REPAIR TOOLS INTEGRATED WITH CERBERUS

#	Tool	Language	Methodology	Target Defect Type
1	Angelix	C/C++	Semantic	Test Failure
2	Prophet	C/C++	Learning	Test Failure
3	Darjeeling	C/C++	Search	Test Failure
4	CPR	C/C++	Semantic	Test Failure
5	VulnFix	C/C++	Semantic	Security Vulnerabilities
6	F1X	C/C++	Search	Test Failure
7	Fix2Fit	C/C++	Search	Test Failure
8	SenX	C/C++	Search	Security Vulnerabilities
9	GenProg	C/C++	Search	Test Failure
10	ExtractFix	C/C++	Semantic	Security Vulnerabilities
11	Verifix	C	Search	Student Assignments
12	Hippodrome	Java	Search	Concurrency Bugs
13	SequenceR	Java	Learning	Test Failure

TABLE II
DETAILS OF THE BENCHMARK SUBJECTS INTEGRATED WITH CERBERUS

#	Benchmark	Language	Type	# Projects	# Bugs
1	ManyBugs	C/C++	Test Failure	6	60
2	VulnLoc	C/C++	Vulnerabilities	11	43
3	ExtractFix	C/C++	Vulnerabilities	7	30
4	ITSP	C	Student Assignments	10	661
5	Hippodrome	Java	Concurrency Bugs	16	25
6	Defects4J	Java	Test Failure	6	75
Total					894

(configure/build/test/validate) the defects in the benchmark, capturing different stages in the repair process.

We provide extensive material with documentation, tutorials, and examples in our repository¹ to support the integration of new repair tools and new defect benchmarks.

III. IMPLEMENTATION

We have implemented CERBERUS with 13 program repair tools and 6 repair benchmarks consisting of real-world applications and student assignments. The repair tools represent different repair methodologies, including learning-based, semantic-based, and search-based techniques. Table I details the program repair tools integrated with CERBERUS representing each methodology.

The benchmark programs consist of different classes of repair tasks, including but not limited to fixing concurrent bugs, generating feedback for student assignments, and repairing test suite failures. Table II describes the details of the benchmark programs integrated with CERBERUS. MANYBUGS [12] and DEFECTS4J [13] are benchmarks consisting of functionality errors reported as test case failures for C/C++ and Java programs, respectively. EXTRACTFIX [14] and VULNLOC [15] are benchmarks for C/C++ programs capturing a security vulnerability with a proof of concept exploit. ITSP [16] is a benchmark consisting of incorrect solutions for student assignments with a correct solution provided as reference. HIPPODROME [17] is a benchmark for concurrency bugs, which does not include any test case.

IV. EVALUATION

We demonstrate the capabilities of CERBERUS by analyzing the performance improvement introduced for repair tasks and

reproducing reported experimented values on literature. First, we analyze the performance improvement gained by using CERBERUS with respect to space and time to set up. For this purpose, we selected VULNLOC benchmark [15], which consists of real-world applications with a single failing test. Table III shows the comparison of running F1X [18] and VULNFIX [19] on Binutils in VULNLOC benchmark. Columns *s* and *t* indicate the space consumed by preparing the subject for repair and the time duration for setup, respectively.

TABLE III
ANALYSIS ON TIME AND SPACE IMPROVEMENT

Bug-ID	Original				Cerberus			
	F1X		VULNFIX		F1X		VULNFIX	
	<i>s</i>	<i>t</i>	<i>s</i>	<i>t</i>	<i>s</i>	<i>t</i>	<i>s</i>	<i>t</i>
CVE-2017-14745	849MB	23s	1.4GB	46s	849MB	27s	597MB	21s
CVE-2017-15020	840MB	30s	1.1GB	37s	840MB	23s	252MB	17s
CVE-2017-15025	851MB	27s	1.1GB	37s	851MB	24s	252MB	17s
CVE-2017-6965	840MB	29s	1.4GB	44s	840MB	24s	582MB	21s

Although both F1X and VULNFIX are repairing the same subjects, they both require different methods for preparing the subjects as observed in the space difference in Table III. This is because F1X does not require the binary files to be built; however, VULNFIX expects the binary executable to be provided as input to the repair. In a traditional environment where the user runs both VULNFIX and F1X in parallel, they would need to keep two copies of the setup, one for each tool. However, using the layered architecture in CERBERUS, we only need one copy of the source code, which saves significant space, as shown in Table III. For example, CVE-2017-14745 setup would require 849MB amount of space, on top of which F1X will run the repair. VULNFIX for the same defect would require an additional 597MB space totaling to 1.4GB, before attempting to repair. In the native setup, the two experiments using F1X and VULNFIX would require 2.3GB. For CERBERUS, the total space required is 1.4GB saving 597MB for the two experiments. Similarly, using caching to re-use previous setup CERBERUS saves time for consecutive repair on the same defect. Note that the space and time saving reported are for two consecutive repairs. For each additional repair, the savings would be a factorial of the initial saving.

Next, we successfully reproduce experimental results for two selected repair tasks. To demonstrate the diversity of repair tools CERBERUS can cater, we select two repair tasks. In contrast to traditional test failure repair, we select fixing student assignments. For this purpose, we executed the education repair tool VERIFIX [9] using CERBERUS and were able to generate similar results as reported in [9]. Table IV shows the results for the latest version of the tool VERIFIX. The repair percentages observed in our experiments are slightly better since we are using the latest version of the tool, which the authors have improved since the publication.

We also demonstrate that CERBERUS can support learning-based repair tools. Especially since learning-based repair tools require different infrastructures to run the repair task. For this purpose, we selected SEQUENCER, a popular learning-based

¹<https://github.com/nus-apr/cerberus/blob/main/doc/Extending.md>

TABLE IV
EXPERIMENTAL RESULTS FOR VERIFIX IN ITSP BENCHMARK USING

Lab-ID	# Assignments	# Programs	Repair %	Average Time (s)
Lab-3	4	63	91.66	9.95
Lab-4	8	117	82.24	15.80
Lab-5	8	82	71.95	3.03
Lab-6	8	79	49.36	6.72

tool, and reproduced the results reported in [10]. Table V shows the results of SEQUENCER on a subset of DEFECTS4J benchmark. We are able to reproduce identical results in terms of candidate patches, compilable patches, plausible patches, and correct patches as reported in [10].

TABLE V
EXPERIMENTAL RESULTS FOR SEQUENCER IN DEFECTS4J BENCHMARK

Metric	Observed	Reported
# Bugs candidate patches found	57	57
# Bugs compilable patches found	52	52
# Bugs plausible patches generated	19	19
# Bugs correct patches generated	14	14

V. RELATED WORK

Several works have been proposed in the literature to address the gap of a standard platform for empirical evaluation in automated program repair. MAESTRO [7] is a recently proposed platform to evaluate automated program repair tools across different benchmarks with a low overhead to the user. The platform is designed to work as micro-service containers communicating via RESTful APIs to perform repair tasks for different benchmark programs. The proposed decentralized approach creates separate micro-service containers, each for the benchmark program and for the program repair tool itself. In contrast, CERBERUS creates a single container encapsulating the benchmark program and the repair tool. This design choice allows CERBERUS to generate on-the-fly Docker containers customized for a repair tool and a benchmark program pair. However, a decentralized approach as proposed in MAESTRO [7] is difficult to extend with repair tools that require a large dependent toolchain. For example, semantic-based repair tools such as CPR [8] require LLVM build infrastructure with KLEE [20] runtime support. Therefore, separating the benchmark program and the repair tool prevents running repair using semantic-based repair tools.

Similar to our approach, REPAIRTHEMALL [5] proposed a framework to evaluate multiple program repair tools across multiple sets of benchmarks. REPAIRTHEMALL implements a monolithic architecture targeted for Java. SECURETHEMALL [6] follows a similar design but focuses on security vulnerabilities. SECURETHEMALL is implemented for C programs and specifically designed to cater only security vulnerabilities. In contrast, CERBERUS is not restricted to a specific language or a class of defects. We have shown that CERBERUS is capable of catering to multiple programming languages, multiple classes of defects, and multiple types of repair techniques (i.e., learning-based, semantic-based, and search-based).

VI. CONCLUSION

We presented a platform that is capable of executing a diverse set of program repair tasks using a multitude of program repair tools and benchmarks. We implemented our solution in CERBERUS and demonstrated the capability to reproduce previously reported results in the literature. Our experiments also showed there is a significant cost saving in terms of setup time. Our vision is that CERBERUS will be extended to be used as a repair service constituting a variety of repair capabilities. As future work CERBERUS will be integrated with testing and analysis tools such as fuzzers, static analyzers, and symbolic execution engines. This will combine bug detection and repair in a single framework.

CERBERUS is open-source and available for use via:
<https://github.com/nus-apr/cerberus>

REFERENCES

- [1] C. L. Goues *et al.*, “Automated program repair,” *Commun. ACM*, vol. 62, no. 12, p. 56–65, nov 2019.
- [2] A. Marginean *et al.*, “Sapfix: Automated end-to-end repair at scale,” in *2019 IEEE/ACM 41st ICSE: SEIP*, 2019.
- [3] S. Kirbas *et al.*, “On The Introduction of Automatic Program Repair in Bloomberg,” *IEEE Software*, vol. 38, pp. 43–51, 2021.
- [4] M. Monperrus, “Automatic software repair: A bibliography,” *ACM Comput. Surv.*, vol. 51, no. 1, jan 2018.
- [5] T. Durieux *et al.*, “Empirical Review of Java Program Repair Tools: A Large-Scale Experiment on 2,141 Bugs and 23,551 Repair Attempts,” in *Proceedings of the 27th ACM ESEC/FSE ’19*, 2019.
- [6] E. Pinconschi *et al.*, “A comparative study of automatic program repair techniques for security vulnerabilities,” in *2021 IEEE 32nd ISSRE*, 2021.
- [7] —, “Maestro: A platform for benchmarking automatic program repair tools on software vulnerabilities,” in *Proceedings of the 31st ACM SIGSOFT ISSTA*, 2022, p. 789–792.
- [8] R. Shariffdeen *et al.*, “Concolic program repair,” in *Proceedings of the 42nd ACM SIGPLAN PLDI*, 2021, p. 390–405.
- [9] U. Z. Ahmed *et al.*, “Verifix: Verified repair of programming assignments,” *ACM Trans. Softw. Eng. Methodol.*, vol. 31, no. 4, jul 2022.
- [10] Z. Chen *et al.*, “Sequencer: Sequence-to-sequence learning for end-to-end program repair,” *IEEE Trans. on Soft. Eng.*, vol. 47, no. 9, pp. 1943–1959, 2021.
- [11] Z. Huang *et al.*, “Using safety properties to generate vulnerability patches,” in *2019 IEEE SoSP*, 2019, pp. 539–554.
- [12] C. Le Goues *et al.*, “The manybugs and introclass benchmarks for automated repair of c programs,” *IEEE Trans. on Soft. Eng.*, vol. 41, no. 12, pp. 1236–1256, 2015.
- [13] R. Just *et al.*, “Defects4j: A database of existing faults to enable controlled testing studies for java programs,” in *Proceedings of the 2014 ISSTA*, 2014, p. 437–440.
- [14] X. Gao *et al.*, “Beyond tests: Program vulnerability repair via crash constraint extraction,” *ACM Trans. Softw. Eng. Methodol.*, vol. 30, no. 2, 2021.
- [15] S. Shen *et al.*, “Localizing vulnerabilities statistically from one exploit,” in *Proceedings of the 2021 ACM Asia CCS*, 2021, p. 537–549.
- [16] J. Yi *et al.*, “A feasibility study of using automated program repair for introductory programming assignments,” in *Proceedings of the 2017 11th ESEC/FSE*, 2017, p. 740–751.
- [17] A. Costea *et al.*, “Hippodrome: Data race repair using static analysis summaries,” *ACM Trans. Softw. Eng. Methodol.*, jul 2022, just Accepted.
- [18] S. Mechtaev *et al.*, “Test-equivalence analysis for automatic patch generation,” *ACM Trans. on Software Engg. and Meth. (TOSEM)*, vol. 27, no. 4, 2018.
- [19] Y. Zhang *et al.*, “Program vulnerability repair via inductive inference,” in *ISSTA*, 2022.
- [20] C. Cadar *et al.*, “Klee: Unassisted and automatic generation of high-coverage tests for complex systems programs,” in *Proceedings of the 8th USENIX OSDI*, 2008, p. 209–224.