

Program Repair Competition

Ridwan Shariffdeen
National University of Singapore
ridwan@comp.nus.edu.sg

Martin Mirchev
National University of Singapore
mmirchev@comp.nus.edu.sg

Abhik Roychoudhury
National University of Singapore
abhik@comp.nus.edu.sg

Abstract—Automated Program Repair (APR) is a rapidly developing technology that can assist developers in fixing software system errors. An essential part of any scientific work is a fair and comprehensive evaluation of the state-of-the-art techniques. At present there is no publicly available standard mechanism to perform comparisons for program repair. In this paper, we propose such a standardization in the form of a competition to foster comparability, advancing the state-of-the-art capabilities and to pave the path for the next-generation repair technology.

Index Terms—automated program repair, competition

I. INTRODUCTION

Improving the productivity of software developers by relieving the burden of manually fixing the ever-increasing amount of software bugs has been the goal of automated program repair technology [1]. The past decade has seen a significant growth in Automated Program Repair (APR) expanding into various tasks in Software Engineering including but not limited to fixing functional errors, repairing security vulnerabilities, assisting educators in teaching programming languages and collateral evolution of software systems [2]. Due to the rapid advancements made by the repair community there are many powerful techniques, making it very difficult to compare.

In general, the repair problem can be formulated as a search problem where; given a buggy program P and a program specification σ (i.e. a test-suite) which P does not satisfy, find a mutation to the program such that it satisfies σ . In practice, σ can be represented in various different forms, focusing on completely different aspects of the program. For example, the program specification σ could be presented as a set of test-cases or as a program assertion that should hold true for all inputs of the program. Additionally, the property being specified by the specification could capture correct functionality of the program or a security property. Repair techniques [2] proposed in the literature assume different types and formats to the inputs albeit attempting the same repair task. In addition each technique differs in its interaction with the repair subjects and in the implementation of its validation, making it difficult to compare and contrast with respect to the quality of the generated patches.

A recent study [3] showed appreciable interest from the software practitioners to incorporate repair technology in their development processes, provided the tools can be integrated in a non-disruptive manner. In order for a repair tool to be used (either making it usable by developers or practitioners simply using research tools) it is important that the existing repair tools improve its usability such that users can easily operate

the tools and is provided with documentation support on how to configure/extend its capabilities. The benefits of developing the tools beyond an academic prototype has little-short term reward but longer-term benefits for wide adoption of repair technology. However, such considerations are generally not included when developing repair techniques, creating an entry barrier for developers to familiarize with the repair technology.

Furthermore, there is no widely adopted set of benchmark programs that remains relevant with the development of the field. Most of the benchmark subjects quickly become obsolete either due to the use of deprecated dependent libraries or simply because the code is no longer maintained. Such programs make it difficult to use to assess the advancement made in the technology and irrelevant for state of the art comparisons. Using out-dated subjects as benchmark programs makes it challenging for the researchers to compare the performance of their techniques, with no reward for the effort. It is also important for the standardized benchmark programs to evolve over-time in order to prevent repair techniques from being over-fitting for a given set of benchmark programs. However, currently there is no mechanism to maintain a standardized benchmark that can remain relevant with the fast-changing landscape of software engineering.

Finally, the recent emergence of large language models [4] has the potential to automatically generate code and will change the traditional processes in software engineering. Traditional program repair technology mainly focused on fixing human-errors in programs by learning fix-patterns from fixes written by developers themselves. In future, software may comprise of code generated by humans and machines, hence it is important as a community to develop techniques that can *co-exist* with AI driven Software Development. Thus, we can explore how both automatic code generation and automated program repair can synergistically work towards of automating software engineering tasks.

We envision a *Program Repair Competition* that provides additional incentive to the community to setup and maintain repair tools and benchmarks relevant for the current needs in practice. A competition can foster the development of standards, tools and benchmarks that can be used to assess the advancement of the technology. It can also provide reward for technical improvements while providing a platform for researchers to compare with state-of-the-art techniques. A competition could also enforce standards for repair tools to agree on. Such a standard can specify how the input and output for the repair problem will be specified, the qualitative

and quantitative measures used to compare techniques and the types of different program specifications that can be repaired. These standards should make it easy to generate and compare new repair techniques. A common platform with tool-agnostic, user-friendly interfaces allows to even add new software programs to be repaired with dynamic workflows.

II. METHODOLOGY

The *program repair competition* will be designed to evaluate state-of-the-art repair tools with respect to effectiveness, efficiency and quality of the results. Effectiveness will be measured using its repair capabilities on different repair tasks while efficiency is measured with respect to usability aspect (i.e. time latency, resource usage) of the tool. Quality of the results (i.e. generated patches) will be evaluated using explainable properties of a patch.

The main **goals** of the competition are:

- 1) to provide visibility for the state-of-the-art repair tools to the software practitioners;
- 2) bridge the gap between the research tools and current trends and practices in the industry;
- 3) establish an evolving community standard to shape the next-generation of repair technology;
- 4) provide incentives for researchers to improve their techniques beyond academic prototypes.

Organization: The competition is designed following the successful model used in SV-COMP [5], the International Competition on Software Verification. The organizers first announce the competition within the community. The organizing committee will comprise of a ‘core’ team and a ‘jury’ committee. The ‘core’ team will be responsible for conducting the evaluations in an automated, reproducible and fair manner. The jury committee is formed with a representative of each of the participating teams. The jury committee will serve as an advisory board for decision making on qualification of tools, benchmarks and setting up standards.

Timeline: The competition will be announced at the International Workshop on Automated Program Repair co-hosted with International Conference on Software Engineering (ICSE). The tool developers will submit their intention to participate in the competition and provide an initial version of their tool which can be used in pre-runs by the organizing committee and report preliminary results to the jury committee and the tool developers. Once the jury committee approves the final list of benchmarks for evaluation, the tool developers will submit the final versions which will be independently evaluated by the organizing committee. Final results will be reported to the team for inspection and approval, and the results will be announced on the competition website. The competition organizers write the final competition report with reviews/approval from the jury committee.

Standards: The organizing committee will announce the rules and regulations for the competition along with the metrics used to evaluate the repair tools. The rule book will define the minimum qualification for a repair tool to participate in the competition, the standard interfaces it should implement

and the expected inputs for the repair task. The benchmark programs will be selected using a predefined set of criteria that captures various aspects of the program under repair; programming languages and program complexity. In addition, a clear definition of the categories and sub-categories used to determine different capabilities of repair will be announced to the participants. Finally, the scoring criteria will be approved by the running jury committee before the evaluation.

Next-Generation: Incorporating a forward looking view for the program repair technology, a dedicated track will be presented as part of the competition that captures different challenges for next-generation repair techniques. The recent growth of automated code generation from large language models [4] will be considered for constructing this track. This track will be different compared to the traditional repair tasks, where the aim is to refine the program repair problem with emerging technologies and challenges.

III. CONCLUSION

Automated program repair can improve the productivity of developers by alleviating the burden of fixing software errors. Due to the rapid advancements in repair technology and the fast-changing landscape in software engineering there is a need to setup a continuously improving set of standards. Inspired by the success in the fields of software verification and testing, we propose a *Program Repair Competition*. At the end of the evaluation all repair tasks, results and tools will be published for reproducibility purposes. Infrastructure, data and all components will be made available for public access.

The potential benefits of participating includes high-visibility of the technology among software practitioners, acknowledgement for technical improvements and advancing the state of the art research in program repair. We invite the community to participate in this program repair competition, which will hopefully play a significant role in shaping the next-generation of repair technology.

ACKNOWLEDGMENT

The work is partially supported by a Singapore Ministry of Education (MoE) Tier 3 grant “Automated Program Repair”, MOE-MOET32021-0001.

REFERENCES

- [1] C. Le Goues, M. Pradel, and A. Roychoudhury, “Automated program repair,” *Commun. ACM*, vol. 62, no. 12, p. 56–65, nov 2019. [Online]. Available: <https://doi.org/10.1145/3318162>
- [2] M. Monperrus, “The Living Review on Automated Program Repair,” HAL Archives Ouvertes, Technical Report hal-01956501, 2018. [Online]. Available: <https://hal.archives-ouvertes.fr/hal-01956501>
- [3] Y. Noller, R. Shariffdeen, X. Gao, and A. Roychoudhury, “Trust enhancement issues in program repair,” in *Proceedings of the 44th International Conference on Software Engineering*, ser. ICSE ’22. New York, NY, USA: Association for Computing Machinery, 2022, p. 2228–2240. [Online]. Available: <https://doi.org/10.1145/3510003.3510040>
- [4] T. B. Brown *et al.*, “Language models are few-shot learners,” *arxiv:2005.14165*, 2020.
- [5] D. Beyer, “Software verification: 10th comparative evaluation (sv-comp 2021),” in *Tools and Algorithms for the Construction and Analysis of Systems*, J. F. Groote and K. G. Larsen, Eds. Cham: Springer International Publishing, 2021, pp. 401–422.