BURT: A Chatbot for Interactive Bug Reporting

Yang Song*, Junayed Mahmud[†], Nadeeshan De Silva*, Ying Zhou[‡], Oscar Chaparro*, Kevin Moran[†], Andrian Marcus[‡], Denys Poshyvanyk*
*William & Mary (USA), [†]George Mason University (USA), [‡]The University of Texas at Dallas (USA)

Abstract—This paper introduces BURT, a web-based chatbot for interactive reporting of Android app bugs. BURT is designed to assist Android app end-users in reporting high-quality defect information using an interactive interface. BURT guides the users in reporting essential bug report elements, i.e., the observed behavior, expected behavior, and the steps to reproduce the bug. It verifies the quality of the text written by the user and provides instant feedback. In addition, BURT provides graphical suggestions that the users can choose as alternatives to textual descriptions.

We empirically evaluated BURT, asking end-users to report bugs from six Android apps. The reporters found that BURT's guidance and automated suggestions and clarifications are useful and BURT is easy to use. BURT is an open-source tool, available at github.com/sea-lab-wm/burt/tree/tool-demo.

A video showing the full capabilities of BURT can be found at https://youtu.be/SyfOXpHYGRo.

I. INTRODUCTION

Bug reports are essential for successful software maintenance and evolution. These reports are expected to provide clear and detailed information related to a defect, including the system's observed behavior (**OB**), the expected behavior (**EB**), and the steps to reproduce the bug (**S2Rs**) [1]. Unfortunately, bug reports are often unclear, incomplete, and/or ambiguous – often causing delays during the bug resolution process [2], [3].

One main challenge that contributes to low-quality bug reports is the *knowledge gap* between what reporters *know* and what developers *need* [1], [4]. This is especially true when the reporters are software end-users, who are often unfamiliar with the system internals and do not know the information elements important for developers (*e.g.*, the OB, EB, and S2Rs) and how to express them. Most bug reporting systems (*e.g.*, GitHub Issues, Jira, or Bugzilla) are not designed to address this knowledge gap, since they are static web forms and templates that: (1) offer limited guidance on *what* information needs to be reported, and *how* it needs to be reported; and (2) do not provide *feedback* to end-users on whether the information they provide is clear and complete. In consequence, the burden of providing high-quality bug information rests mainly on the reporters.

We propose a web-based chatbot for interactive BUg repoRTing (or BURT). The *software engineering challenge* addressed by BURT is ensuring high-quality bug reporting by end-users, considering the above-mentioned limitations of existing bug reporting systems. The *envisioned users* for BURT are *end-users* who report problems with their app. BURT guides the users during reporting essential bug report elements (*i.e.*, OB, EB, and S2Rs), offering instant quality verification, corrections, and graphical suggestions. BURT's *usage methodology* is described in detail in section II-B.

BURT implements techniques based on natural language processing, dynamic software analysis, and automated bug report quality assessment. We designed and instantiated BURT as a stand-alone web system for Android apps that focuses on bugs manifesting in the app's UI.

We *empirically evaluated* BURT, asking 18 end-users to report 12 bugs from six Android apps using BURT, and assess their experience. The reporters found BURT's guidance and automated suggestions/clarifications to be useful, accurate, and easy to use. Moreover, the bug reports collected by BURT are of higher quality than reports collected via a traditional template-based bug reporting system.

BURT is an open-source tool hosted on GitHub [5] that can be used for any Android app project. More details about BURT's algorithms and evaluation can be found in its original research paper [6].

II. THE BURT INTERACTIVE BUG REPORTING TOOL

BURT is a standalone web-based chatbot, currently tailored for Android apps, that aims to collect high-quality information from the reporter through a guiding dialog. It generates a bug report with textual bug descriptions and app screen captures.

BURT (i) guides the user in reporting essential bug report elements (the OB, EB, and S2Rs), (ii) checks the quality of these elements and offers instant feedback about issues, and (iii) provides graphical suggestions such as the next S2Rs.

A. Burt's Graphical User Interface (GUI)

BURT's GUI is composed of a standard chatbot interface and various panels for interactive bug reporting (see figure 1). The *Chat Box* 1 allows the reporter to provide textual descriptions of the OB, EB, and S2Rs and select BURT's graphical suggestions (e.g., the next S2Rs via screenshots). The Reported Steps Panel 2 enumerates and displays the S2Rs that the user has reported, allowing them to edit the steps to correct mistakes. The Screen Capture Panel 3 displays screen captures of the last three S2Rs. The Quick Action Panel 4 provides buttons to finish reporting the bug, restart the reporting session, and (pre)view the bug report. The *Tips Panel* **5** displays suggestions to reporters on how to use BURT and how to better express the OB, EB, and S2Rs. The tips change depending on the current stage of the conversation. BURT also provides a Developer Panel that allows developers to add new apps to BURT (via the *Settings* icon next to the *Help* button).

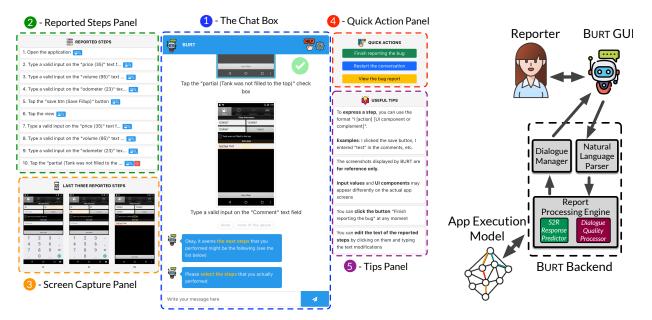


Fig. 1 BURT's graphical user interface

Fig. 2 Overview of BURT's architecture

B. Reporting a Bug with BURT

To report a bug, the user first selects the target app exhibiting the defect by clicking on its icon—BURT lists the apps that it supports. Then, BURT guides the reporter through three phases: OB, EB, and S2R reporting. In each phase, BURT prompts the user to provide individual descriptions of the OB, EB, and S2Rs, respectively. BURT automatically parses the descriptions and verifies their quality by matching them to states of a GUI-level execution model for the app (see Sec. III-C).

If the OB/EB/S2R is *matched* to an app screen from BURT's execution model, BURT asks the user to confirm the matched screen. If the user confirms, BURT proceeds to the next phase of the conversation (*e.g.*, asking for the EB or next S2Rs); otherwise, BURT asks the user to rephrase the bug element.

If there are *no* app screen matches, BURT suggests the user to revise their description and asks them to rephrase the OB/EB/S2R. With a new description, the quality verification is re-executed. If there are *multiple* matches, BURT provides a list of up to five app screenshots (derived from the app execution model) that match the description. The user can then inspect the app screens and select the one that she believes best matches her description of the bug element. If none are selected, BURT suggests additional app screens, if any. If the user selects one app screen, BURT saves the bug element description and screen, and proceeds to prompt the user for the next bug element. After three unsuccessful attempts to provide a high-quality description, BURT records the (last) provided description for bug report generation. This process proceeds for each bug element starting with the OB.

BURT includes an additional feature to help users save time writing S2Rs: it suggests the probable next S2Rs that the user may have performed during actual app usage. BURT suggests the first five S2Rs from the most likely path from the current state to the OB state in the app execution model. This dialogue flow uses a predictive algorithm that uses BURT's execution

model (see section III-C). The suggestions are displayed as a list of generated app screens, each screen representing a S2R. The generated screen is visually annotated with a yellow oval highlighting the GUI component (*e.g.*, a button) executed by the step. The user can select none, one, or multiple of the suggested S2Rs. When a S2R is selected, BURT suggests additional S2Rs, if any. When none are selected and BURT has more suggestions, BURT asks the user if they want additional suggestions. If so, BURT displays them. Otherwise, BURT prompts the user to describe the next S2R.

C. Adding New Apps to BURT

Developers can add new apps to BURT through the *Developer Panel*, which requires the developer to upload the app icon and ZIP files containing app execution data required by BURT to build the app execution model. BURT uploads/extracts the files, parses the data to identify the app name and version, verifies the data format, and stores the data for later use. BURT provides feedback if there are any detected issues with the data.

III. BURT'S ARCHITECTURE & IMPLEMENTATION

BURT has three main components (see figure 2). The *Natural Language Parser (NL)* parses users' bug descriptions. The *Dialogue Manager (DM)* implements the conversation flows for the reporting process. The *Report Processing Engine (RP)* matches the parsed bug descriptions to the app execution model, to assess bug element quality and provide suggestions.

A. Natural Language Parser (NL)

BURT parses the textual OB/EB/S2R descriptions using dependency parsing via the Stanford CoreNLP toolkit [7], which produces a tree of grammatical dependencies between words. BURT first utilizes a heuristic-based approach from our prior work [8] to identify the sentence type of each user message (*e.g.*, conditional, imperative, or passive voice). Then, BURT executes one of its 16 parsing algorithms (one

for each sentence type) that traverse the tree to extract the relevant words from the sentences. This parsing is based on our prior work on S2R quality assessment [9] and extracts a phrase using the following format: [subject] [action] [object] [preposition] [object2]. For example, for the Mileage app [10], the OB sentence "The fuel economy shows a NaN value on page", written in present tense, is parsed as [fuel economy] [shows] [NaN value] [on] [page]. The S2R sentence "Save the car fillup", written imperatively, is parsed as [user] [saves] [car fillup].

B. Dialogue Manager (DM)

BURT's dialogue flow guides users to report the OB, EB, and S2Rs. BURT's dialogue is multi-modal and capable of suggesting both natural language and graphical elements, *e.g.*, screenshots, to assist the user through the reporting process. The DM relies upon the RP engine to assess the quality of the bug elements reported by end users. There are two main dialogue flows that BURT navigates: (i) performing quality checks on written bug report elements (applies to all bug elements), and (ii) automated suggestion of S2Rs (for S2Rs only).

C. Report Processing Engine (RP)

BURT's RP Engine consists of three sub-components:

- (1) the App Execution Model is a directed graph where nodes represent app screens and the edges are transitions between screens, triggered by GUI interactions (e.g., taps or type events) on GUI components (e.g., buttons or text fields). The graph screens contain the hierarchy of GUI components with their metadata (e.g., labels) and a screen capture. The graph edges contain the action type, the interacted GUI component, and a screen capture highlighting the component. Burt builds the execution model from app usage data collected automatically, via automated app exploration, or manually (see section III-D).
- (2) the *Dialogue Quality Processor* performs quality verification of the parsed OB/EB/S2R descriptions, by mapping them to app states/interactions from the execution model. A textual description is *high-quality* if it can be matched to the model, otherwise, it is *low-quality*. This definition and BURT's dialogue features that prompt users to improve low-quality descriptions aim to reduce the knowledge gap between reporters and developers. To perform quality verification, BURT extends the bug description resolution/matching algorithm from our prior work [9] and performs exploration of the execution model, driven by the matching of the reported OB/EB/S2Rs and user confirmations during the bug reporting process.
- (3) the S2R Response Predictor determines and suggests to the reporter the next S2Rs that she may have performed in practice. BURT implements a shortest-path approach to predict the next S2Rs [6]. BURT determines the paths between the current graph state and the corresponding OB state, and then, it computes the likelihood score based on the execution model edge weights, which are higher for manual app usage [6].

D. Collecting App Execution Data

BURT requires app execution data for building the execution model. This data encodes sequential interactions made on the app features and comes from two sources: systematic app exploration via CRASHSCOPE [11], [12] and crowdsourced app usage (*e.g.*, from app users or developers).

CRASHSCOPE's *GUI-ripping engine* automatically generates app execution data (app interactions such as *taps* or *type* events) by running an app on a mobile device/emulator and a set of systematic exploration strategies (*e.g.*, top-down or bottom-up) on the app [11], to uncover as many app screens as possible and interact with most screen GUI components.

Crowdsourced app usage data *complements* the automated app usage data [13]. This data can be collected through built-in trace recording app features that capture app usages "in the wild", or during in-house GUI-level app testing performed by developers. BURT provides two tools to assist trace capture: AVT and TRACEREPLAYER. AVT is a desktop app that allows humans to collect video recordings and getevent traces from a mobile device/emulator while humans are using an app. TRACEREPLAYER parses the collected AVT traces (*i.e.*, from humans) and converts them into the CRASHSCOPE data format. In this way, BURT can read/use this data to augment the app execution model (see details in BURT's repository [5]).

E. Burt Implementation

BURT is currently implemented as a web application, using the React Chatbot Kit [14] and Spring Boot [15]. BURT also provides command-line tools for CRASHSCOPE and TRACEREPLAYER, and the desktop AVT app, along with detailed documentation on how to use them. BURT's implementation is tailored for Android applications, however, its underlying techniques are generic enough to be easily adapted for other types of software—the App Execution Model Data Collection is the only platform-specific part.

IV. BURT'S EVALUATION

We conducted an empirical study to evaluate: (1) BURT's perceived usefulness/usability ($\mathbf{RQ}_1/\mathbf{RQ}_2$); (2) BURT's intrinsic accuracy in performing bug report element quality verification and prediction (\mathbf{RQ}_3); and (3) the quality of the bug reports collected with BURT (\mathbf{RQ}_4).

A. Methodology

We selected 12 Android bugs (seven crashes, one handled error, and four non-crashes) from the dataset of our prior work [16]. The bugs come from six Android apps of different domains: AntennaPod, Time Tracker, GnuCash, GrowTracker, and Droid Weight. To collect app execution data, we executed CRASHSCOPE on these apps and asked two computer science (CS) students to use their main features—see our original paper for more details [6]. We recruited 18 participants to report these bugs (each reported three bugs) using BURT (most had little or no bug-reporting experience) and asked them to evaluate their experience via a questionnaire (with Likert-scale and open-ended questions). We analyzed the conversations the reporters had with BURT and measured how accurate BURT was during the reporting process. We asked additional 18 participants to report the same bugs with

a template-based bug-reporting system (a.k.a. ITRAC) and analyzed the collected bug reports to measure their quality based on the bug element correctness framework from our prior work [9]. ITRAC resembles existing issue trackers as it implements a web form with text fields and templates that explicitly ask for the bug summary/title and the OB/EB/S2Rs.

B. Results

1) RQ_1/RQ_2 : BURT's User Experience: The participants evaluated the **usefulness** of BURT's main features.

<u>Screen Suggestions</u>: Half of the 18 reporters agreed that the screen suggestions were *useful*, and another half (nine reporters) agreed that they were *sometimes useful*.

OB/EB/S2R Quality Verification: The reporters had a positive impression on how often BURT correctly verified the quality of their bug descriptions. BURT was able to always (sometimes) verify the OB/EB/S2R descriptions of 9/10/11 (9/6/6) reporters (out of 18). Only two/one participant(s) felt that BURT rarely recognized and verified their EB/S2Rs.

<u>BURT Messages & Questions</u>: 11 of 18 users often understood BURT's messages/questions, while 6 reporters understood them sometimes. Only one reporter rarely understood them.

<u>Panel of Reported S2Rs</u>: BURT's panel of reported S2Rs was deemed to be useful by nine participants and somewhat useful by six participants – one reporter found it somewhat useless.

The reporters also assessed Burt's overall **ease of use**. 12 of 18 reporters indicated Burt was either easy or somewhat easy to use. Four reporters were neutral, while two reporters expressed that Burt was somewhat difficult to use.

- RQ₁/RQ₂ Summary: Overall, reporters found BURT's screen suggestions and S2R panel useful and BURT is easy to use. Reporters also suggested improvements to BURT to support additional wording of bug report elements and provide more accurate suggestions. Improvements are planned for future work to improve BURT 's ability to recognize additional vocabulary and ways of phrasing the OB/EB/S2Rs.
- 2) RQ_3 : BURT's Intrinsic Accuracy: We analyzed the 54 conversations with BURT to assess how often BURT was able to 1) match OB/EB/S2R descriptions to the execution model, and 2) suggest relevant OB/S2R app screens to the reporters.
- <u>OB reporting</u>: BURT matched the OB description to the correct screen in 3 of 54 (5.5%) conversations and multiple screens in 35 of 54 (64.8%) conversations. In 29 of 35 (80%) conversations, the reporters selected one of the suggested screens. Overall, BURT correctly matched the users' OB descriptions in 32 of 54 (59.3%) conversations.

<u>EB reporting</u>: BURT correctly matched the EB against the OB screen in 17 of 32 (53.1%) cases without having to ask the reporter for confirmation. In 6 of 32 (18.8%) cases, BURT needed to ask the users for confirmation. In the remaining 9 cases, BURT struggled to parse the EB description.

<u>S2R reporting</u>: In the 54 conversations, BURT matched 205 of the written S2R descriptions from the reporters. BURT matched

the correct screen in 157 of 205 (76.6%) cases. As for S2R prediction, among 32 conversations with a matched OB screen, S2R prediction occurred 146 times (mean: 4.6 per conversation). The reporters selected 1.6 of the 3.9 suggested S2Rs (on avg.) in 91 of 146 cases (62.3%).

- \mathbf{RQ}_3 Summary: The results confirm the users' perception on the usefulness of screen suggestions and ability of performing correct bug element verification, which indicates that the techniques used in designing BURT's components are adequate.
- 3) RQ_4 : Bug Report Quality: We compared the quality of the $54 \times 2 = 108$ bug reports, collected with ITRAC and BURT. <u>S2R Quality</u>: On average, BURT's reports contain fewer incorrect S2Rs than the ITRAC reports (8.3% vs. 20.4%) and fewer missing S2Rs (19.4% vs. 32%).

<u>OB/EB Quality</u>. BURT and ITRAC reports have a comparable number of incorrect EB descriptions (8 vs. 6 out of 54 reports). However, BURT reports have more incorrect OB descriptions compared to ITRAC reports (16 vs. 8 out of 54 reports).

RQ₄ **Summary:** BURT bug reports have higher-quality S2Rs than ITRAC reports, and comparable EB descriptions. BURT improvements are needed to better collect OB descriptions.

V. MOST RELATED EXISTING TOOLS

Moran et al. [4], [17] proposed FUSION, a web system that allows the user to report the S2Rs graphically via dropdown lists of GUI components and actions (taps, swipes, etc.). Song et al. [18] proposed BEE, a GitHub plugin that identifies the type of each sentence in bug reports and alerts reporters of missing OB/EB/S2Rs. Fazzini et al. proposed EBUG [19], a mobile app bug reporting system, similar to FUSION, that suggests possible future S2Rs as they are written. Shi et al. [20] proposed BUGLISTENER, an approach that identifies bug report dialogs in community live chats and synthesizes bug reports from them.

BURT has two main advancements over these prior tools: (1) it supports end-users with little or no bug reporting experience; for example, FUSION was not specifically designed to end-users, as inexperienced users found it *more difficult* to use; and (ii) it offers an interactive interface, supporting automated suggestions, instant quality verification, and prompts for information clarification.

VI. FINAL REMARKS AND FUTURE WORK

BURT is a web-based chatbot for interactive bug reporting. Unlike existing bug reporting systems, BURT can guide end-users in reporting essential bug report elements (*i.e.*, OB, EB, and S2Rs), provide instant feedback about issues, and produce graphical suggestions of the elements that are likely to be reported next. BURT can help end-users easily report bugs and provide higher-quality bug reports. As future work, we plan to integrate BURT with existing issue trackers.

ACKNOWLEDGEMENTS

This work is supported in part by the NSF grants: CCF-1955837, CCF-1955853, CCF-2217733, and CCF-2007246. Any opinions, findings, and conclusions expressed herein are the authors and do not necessarily reflect those of the sponsors.

REFERENCES

- N. Bettenburg, S. Just, A. Schröter, C. Weiss, R. Premraj, and T. Zimmermann, "What Makes a Good Bug Report?" in FSE'08.
- [2] https://github.com/dear-github/dear-github, 2019.
- [3] T. Zimmermann, N. Nagappan, P. J. Guo, and B. Murphy, "Characterizing and predicting which bugs get reopened," in *ICSE'12*.
- [4] K. Moran, M. Linares-Vásquez, C. Bernal-Cárdenas, and D. Poshyvanyk, "Auto-completing Bug Reports for Android Applications," in FSE'15.
- [5] "https://github.com/sea-lab-wm/burt/tree/tool-demo," 2023.
- [6] Y. Song, J. Mahmud, Y. Zhou, O. Chaparro, K. Moran, A. Marcus, and D. Poshyvanyk, "Toward interactive bug reporting for (android app) end-users," in ESEC/FSE'22.
- [7] "https://stanfordnlp.github.io/CoreNLP/," 2022.
- [8] O. Chaparro, J. Lu, F. Zampetti, L. Moreno, M. Di Penta, A. Marcus, G. Bavota, and V. Ng, "Detecting missing information in bug descriptions," in ESEC/FSE'17.
- [9] O. Chaparro, C. Bernal-Cárdenas, J. Lu, K. Moran, A. Marcus, M. Di Penta, D. Poshyvanyk, and V. Ng, "Assessing the quality of the steps to reproduce in bug reports," in ESEC/FSE'19.
- [10] "Mileage," https://fossdroid.com/a/mileage.html, 2021.
- [11] K. Moran, M. Linares-Vasquez, C. Bernal-Cardenas, C. Vendome, and D. Poshyvanyk, "Crashscope: A practical tool for automated testing of android applications," in *ICSE'17*.

- [12] K. Moran, M. Linares-Vasquez, C. Bernal-Cardenas, C. Vendome, and D. Poshyvanyk, "Automatically discovering, reporting and reproducing android application crashes," in *ICST'16*, 2016.
- [13] M. Linares-Vasquez, M. White, C. Bernal-Cardenas, K. Moran, and D. Poshyvanyk, "Mining android app usages for generating actionable gui-based execution scenarios," in MSR'15, 2015.
- [14] "React chatbot kit," https://tinyurl.com/yhz3ws6h, 2022
- [15] "Spring boot," https://spring.io/projects/spring-boot, 2022.
- [16] N. Cooper, C. Bernal-Cárdenas, O. Chaparro, K. Moran, and D. Poshyvanyk, "It takes two to tango: Combining visual and textual information for detecting duplicate video-based bug reports," in *ICSE*'21.
- [17] K. Moran, M. Linares-Vasquez, C. Bernal-Cardenas, and D. Poshyvanyk, "Fusion: A tool for facilitating and augmenting android bug reporting," in 2016 IEEE/ACM 38th International Conference on Software Engineering Companion (ICSE-C), 2016, pp. 609–612.
- [18] Y. Song and O. Chaparro, "Bee: a tool for structuring and analyzing bug reports," in *ESEC/FSE'20*.
- [19] M. Fazzini, K. P. Moran, C. Bernal-Cardenas, T. Wendland, A. Orso, and D. Poshyvanyk, "Enhancing mobile app bug reporting via real-time understanding of reproduction steps," TSE, 2022.
- [20] L. Shi, F. Mu, Y. Zhang, Y. Yang, J. Chen, X. Chen, H. Jiang, Z. Jiang, and Q. Wang, "Buglistener: identifying and synthesizing bug reports from collaborative live chats," in *ICSE* 22.