# Chatting with GPT-3 for Zero-Shot Human-Like Mobile Automated GUI Testing

Zhe Liu<sup>1,3</sup>,Chunyang Chen<sup>4</sup>, Junjie Wang<sup>1,2,3,\*</sup>, Mengzhuo Chen<sup>1,3</sup>, Boyu Wu<sup>3</sup>, Xing Che<sup>1,3</sup>, Dandan Wang<sup>1,3</sup>, Qing Wang<sup>1,2,3,\*</sup>

<sup>1</sup>Laboratory for Internet Software Technologies, <sup>2</sup>State Key Laboratory of Computer Sciences, Science & Technology on Integrated Information System Laboratory

Institute of Software Chinese Academy of Sciences, Beijing, China; <sup>3</sup>University of Chinese Academy of Sciences, Beijing, China; \*Corresponding author <sup>4</sup>Monash University, Melbourne, Australia;

liuzhe 181@mails.ucas.ac.cn, Chunyang.chen@monash.edu, junjie@iscas.ac.cn, wq@iscas.ac.cn, wqw. ac.cn, wqw. ac.c

## **ABSTRACT**

Mobile apps are indispensable for people's daily life, and automated GUI (Graphical User Interface) testing is widely used for app quality assurance. There is a growing interest in using learning-based techniques for automated GUI testing which aims at generating human-like actions and interactions. However, the limitations such as low testing coverage, weak generalization, and heavy reliance on training data, make an urgent need for a more effective approach to generate human-like actions to thoroughly test mobile apps. Inspired by the success of the Large Language Model (LLM), e.g., GPT-3 and ChatGPT, in natural language understanding and question answering, we formulate the mobile GUI testing problem as a Q&A task. We propose GPTDroid, asking LLM to chat with the mobile apps by passing the GUI page information to LLM to elicit testing scripts, and executing them to keep passing the app feedback to LLM, iterating the whole process. Within it, we extract the static context of the GUI page and the dynamic context of the iterative testing process, design prompts for inputting this information to LLM, and develop a neural matching network to decode the LLM's output into actionable steps to execute the app. We evaluate GPTDroid on 86 apps from Google Play, and its activity coverage is 71%, with 32% higher than the best baseline, and can detect 36% more bugs with faster speed than the best baseline. GPTDroid also detects 48 new bugs on the Google Play with 25 of them being confirmed/fixed. We further summarize the capabilities of GPTDroid behind the superior performance, including semantic text input, compound action, long meaningful test trace, and test case prioritization.

## **KEYWORDS**

Automated GUI testing, Large language model

# 1 INTRODUCTION

Mobile apps have become increasingly popular over the past decade, with millions of apps available for download from app stores like the Apple App Store [3] and Google Play Store [4]. With the rise of app importance in our daily life, it has become increasingly critical for app developers to ensure that their apps are of high quality and perform as expected for users. To avoid time-consuming and labour-extensive manual testing, automated GUI (Graphical User Interface) testing is widely used for quality assurance of mobile apps [62–64, 96, 97, 101] i.e., dynamically exploring mobile apps by

executing different actions such as scrolling, clicking based on the program analysis to verify the app functionality.

However, existing GUI testing tools such as probability-based or model-based ones [50, 78, 85] suffer from low testing coverage when testing practical commercial apps, meaning that they may miss important bugs and issues. This is because of the complex and dynamic nature of modern mobile apps [31, 36, 50, 72, 77, 78], which can have hundreds or even thousands of different screens, each with its own unique set of interactions and possible user actions and logic. In addition, test inputs generated by these methods are significantly different from real users' interaction traces [73], resulting in the low testing coverage. To address these limitations, there has been a growing interest in using deep learning (DL) [35, 54, 98, 100] and reinforcement learning (RL) [29, 61, 71, 76] methods for automated mobile GUI testing. By learning from human testers' behavior, these methods aim to generate human-like actions and interactions that can be used to test the app's GUI more thoroughly and effectively. These approaches are based on the idea that the more closely the actions performed by the testing algorithm mimic those of a human user, the more comprehensive and effective the testing will be.

Nevertheless, there are still some limitations with these DL and RL based GUI testing methods. First, learning algorithms require large amounts of data which is difficult to collect real-world users' interactions. Second, learning algorithms are designed to learn and predict from training data, so they may not generalize well to the new, unseen situations, as apps are constantly evolving and updating. Third, mobile apps can be non-deterministic, meaning that the outcome of an action may not be the same every time, it is performed (e.g., clicking the "delete" button from a list with the last content would produce an empty list for which the delete button no longer works) which specifically makes it difficult for RL algorithms to learn and make accurate predictions. Therefore, another more effective approach to generate human-like actions is highly needed to thoroughly test mobile apps.

The emerging Large Language Model (LLM) [19, 27, 87, 102] trained on ultra-large-scale corpus, which shows promising performance in natural language understanding, logical reasoning and question answering in recent years. For example, GPT-3 [19] (Generative Pre-trained Transformer-3) is one LLM from OpenAI with 175 billion parameters trained on a massive dataset including existing test scripts and bug reports, which makes it capable of understanding and generating text across a wide range of topics and domains.

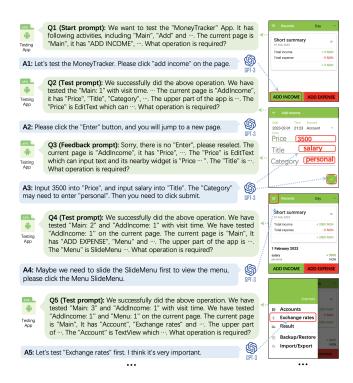


Figure 1: Demonstrated example of how GPTDroid works.

The success of ChatGPT<sup>1</sup> based on GPT-3 demonstrates that LLM can understand human knowledge and interact with humans as a knowledgeable expert. Inspired by ChatGPT<sup>2</sup>, we formulate GUI testing problem as a questions & answering (Q&A) task, i.e., asking the LLM to play a role as a human tester to test the target app.

In detail, we propose a new approach, GPTDroid, asking LLM to chat with mobile apps by passing the GUI page information to LLM to elicit testing scripts, and executing them to keep passing the app feedback to LLM, iterating the whole process. To convert the visual information of the app GUI into the corresponding natural language description, we first extract the semantic information of the app and GUI page by decompiling the target app and view hierarchy files, as well as the dynamic context information of the iterative testing process. Similar to the screenreader through which the blind interact with mobile apps [10, 13], we design linguistic patterns to generate prompts for describing the current GUI page as the input to LLM, which provides ways for LLM to interact with the app. Given the natural language described answer from LLM, we decode it into actionable steps to execute the target app by developing a neural matching network model.

Compared with conventional learning-based algorithms mentioned above, our approach is based on LLM as a zero-shot tester that does not require any training data or corresponding computational resources for training the model. One example chat log can be seen in Figure 1. LLM can understand the app GUI, and provide detailed actions to navigate the app (e.g., A1-A5 at Figure 1). To compensate for its wrong prediction (A2 at Figure 1), the real-time

feedback by our approach guides it to regenerate the input until triggering a valid page transition. It remains clear testing logic even after a long testing trace to make complex reasoning of actions (A3, A4 at Figure 1), and it can prioritize to test important functions earlier (e.g., A5 at Figure 1). A more detailed analysis of our approach's capability and reasons behind it is in Section 6.

To evaluate the effectiveness of GPTDroid, we carry out an experiment on 86 popular Android apps in Google Play with 129 bugs. Compared with 9 common-used and state-of-the-art baselines, GPTDroid can achieve more than 32% boost in activity coverage than the best baseline, resulting in 71% activity coverage. As GPTDroid can cover more activities, the method can detect 36% more bugs with faster speed than the best baseline. Apart from the accuracy of our GPTDroid, we also evaluate the usefulness of our GPTDroid by detecting unseen crash bugs in real-world apps from Google Play. Among 216 apps, we obtain 48 crash bugs with 25 of them being confirmed and fixed by developers, while the remaining are still pending. To reveal reasons behind the promising performance of our approach, we further investigate the experiment results qualitatively and summarize 4 findings in discussion including semantic text input, compound action, long meaningful test trace, and test case prioritization.

The contributions of this paper are as follows:

- Vision. The first work to formulate the automated GUI testing problem to a question & answering task by bringing LLM into GUI testing domain, as far as we know.
- **Technique.** A novel approach GPTDroid based on "pre-train, prompt and predict" paradigm of the LLM by understanding the GUI semantic information and dynamic context of the iterative testing process, for automatically inferring possible operation steps.
- Evaluation. Effectiveness and usefulness evaluation of GPTDroid in the real-world apps with practical bugs detected.
- Insight. Detailed qualitative analysis in discussion revealing the reasons why LLM can generate human-like actions for app testing.

#### 2 BACKGROUND

## 2.1 Android GUI and GUI event

For a mobile app, the user interface (UI) is the place where interactions between humans and machines occur. App developers design UI to help users understand the features of their apps, and users interact with the apps through the UI. To help developers manipulate views flexibly, Android Software Development Kit (SDK) [2] allows developers to build UI in Android source code using the View and ViewGroup objects. The View objects are usually called "widget" (e.g., ImageView, TextView), while the ViewGroup objects are usually called "layout" which provides various layout structures (e.g., LinearLayout, RelativeLayout). Every widget and layout in the Android source code has its specific attributes, which are used to set its boundaries, clickable and referenced external resources, etc. During the running process, developers can obtain the view hierarchy file corresponding to the current UI page (screenshot) through "uiautomator dump" of the Android Debug Bridge (ADB) command [2]. The view hierarchy file includes the widget information (coordinate information, ID, widget type, text description, etc.),

<sup>&</sup>lt;sup>1</sup>https://openai.com/blog/chatgpt/

<sup>&</sup>lt;sup>2</sup>ChatGPT cannot be directly applied for mobile GUI testing, so we adopt GPT-3 which provides official API and easy-to-control results.

and the layout information on the current UI page [11], which can be used by automated GUI testing tools to obtain the information about widgets.

The graphical user interface (GUI) is the most important type of UI for most mobile apps, where apps present content and actionable widgets on the screen and users interact with the widgets using actions (GUI event) such as clicks, swipes, and text inputs. In the process of using apps, users often need complex GUI events to jump from one page to another, such as filling in multiple text input widgets, sliding widgets to the left to operate, and long pressing widgets to delete.

# 2.2 Large Language Model & Prompt

Pre-trained Large Language Models (LLMs) have been shown effective in many natural language processing (NLP) tasks. It is trained on ultra-large-scale corpus and can understand the input prompts (sentences with prepending instructions or a few examples) and generate reasonable text. When pre-trained on billions of samples from the Internet, recent LLMs (like GPT-3 [19], PaLM [27] and OPT [102]) encode enough information to support many NLP tasks [60, 80, 99]. GPT-3 [19] is one of the most popular and powerful LLM which has great performance in many text generation tasks. It is based on the transformer model [87] including input embedding layers, masked multi-self attention, normalizaiton layers, and feed-forward in Fig 2. Given a sentence, the input embedding layer encodes it through the word embedding. The multi-self attention layer is used to divide a whole high-dimensional space into several different subspaces to calculate the similarity. The normalization layer is implemented through a normalization step that fixes the mean and variance of each layer's inputs. The feed-forward layer compiles the data extracted by previous layers to form the final output.

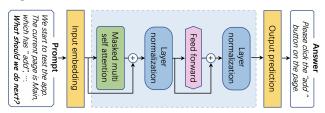


Figure 2: The model structure of GPT-3.

Recently, prompt engineering has been proposed to close the gap between pre-training and downstream tasks. Instead of designing a new training objective for each downstream task, prompt engineering rewrites the input by adding a natural language instruction such as "This is XX app, On its xxx page, it has xxx. What to do next?" to reuse the masking objective for downstream tasks. Formally, a popular prompt engineering employs a prompt template  $T_{prompt}(.)$  to convert the input X to prompt input  $X_{prompt} = T_{prompt}(X)$ . The prompt template is a textual string with unfilled slots to fill the input X and a question.

For automated GUI testing test script generation, the filled input X is the GUI information of the page, and LLM attempts to generate the operation steps to be executed. Although achieving promising results in various NLP tasks, the standard prompt is ineffective for automated GUI testing, because LLM cannot understand the

visual information of the GUI page or the source code of the app. Therefore, this paper focuses on how to describe the GUI page information to let LLM better understand, and how to decode the LLM's feedback to actionable operation steps to execute the app.

## 3 APPROACH

We model the GUI testing as a Question & Answering (Q&A) problem, i.e., asking the LLM to play a role as a human tester, and enabling the interactions between the LLM and the app under testing. To realize this, we propose GPTDroid, as demonstrated in Figure 3. It extracts the static and dynamic context of the current GUI page, encodes them into prompt questions for LLM, decodes LLM's feedback answer into actionable operation scripts to execute the app; and iterates the whole process. Armed with the knowledge learned from large-scale training corpus, GPTDroid would have the potential to guide the testing in exploring more diversified pages, conducting more complex operational actions, and covering more meaningful operational sequences.

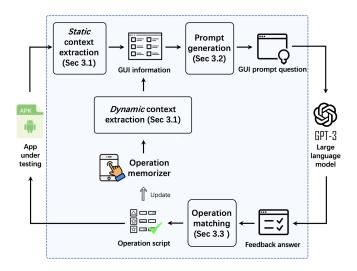


Figure 3: Overview of GPTDroid.

Specifically, in each iteration of the testing, GPTDroid first obtains the view hierarchy file of the mobile app and extracts the static context including the app information, information of the current GUI page, and details of each widget in the page. It also maintains a testing operation memorizer, and extracts the latest dynamic information from it to indicate current testing progress. Based on this contextual information, we design linguistic patterns for generating the prompts as input of LLM, and the LLM would output the operational steps described in natural language. We then design a natural matching network to match the operational steps with the GUI events (i.e., widgets) of the app to enable it to be automatically executed.

## 3.1 Context Extraction

Despite of its excellence on various tasks, the performance of LLM can be significantly influenced by the quality of its input, i.e., whether the input can precisely describe what to ask [26, 55, 104]. In the scenario of this interactive mobile GUI testing, we need to

Id	Attribute	Description	Examples							
	Static context - App information									
1 2	AppName Activities	Name of the app under testing List of names for all activities of the app, obtained from <i>AndroidManifest.xml</i> file	AppName = "Money Tracker"   Activities = ["Main", "AddAccount", "Import", "Income",]							
	Static context - page GUI information									
3	ActivityName	Activity name of the current GUI page	ActivityName = "AddPersonalInformation"							
4	Widgets	List of all widgets in current page, represented with text/id	Widgets = ["Edit Account", "btn_income",]							
5	Position	Relative position of widgets, obtained through their coordinates	Upper = ["Welcome",], Lower = ["Add Income",]							
		Static context - widget information								
6	WidgetText	Widget text, obtained by field 'text' or 'hint-text'	WidgetText = "Welcome to the Money Tracker!"							
7	WidgetID	Widget ID, obtained by field 'resource-id'.	WidgetID = "add_account"							
8	WidgetCategory	Category: TextView, EditText, ImageView, etc, obtained by field 'class'	WidgetCategory = "TextView"							
9	WidgetAction	Widget action, obtained by field 'clickable', such as click, input, etc.	WidgetAction = "Click"							
10	NearbyWidget	Nearby widgets, obtained by the text of parent widgets and sibling widgets	NearbyWidget = "your income: [SEP] \$ "							
		Dynamic context								
11	PageVisits	Set of tested GUI pages with page visits number	PageVisits = [{"Main": "5" . "Account": "3" . "Setting": "1"}]							

Table 1: Extracted information and examples.

accurately depict the GUI page currently under test, as well as its contained widgets information from a more micro perspective, and the app information from a more macro perspective. Furthermore, to act like a human tester, GPTDroid should also capture current testing progress so as to recommend the testing operations from a more global viewpoint to potentially cover more activities and avoid duplicate explorations. This section describes which information will be extracted, organized into static context and dynamic context to facilitate reading. And Section 3.2 will describe how we organize this information into the style that LLM can better understand.

Set of tested widgets of current GUI page with widget visits number

WidgetVisits

3.1.1 Static Context Extraction. Static context relates to the information of the app, the GUI page currently tested, and all the widgets on the page. The app information is extracted from the AndroidMaincast.xml file, while the other two types of information are extracted from the view hierarchy file, which is obtainable with UIAutomator [86]. Table 1 presents the summarized view of them.

**App information** provides the macro-level semantics of the app under testing, which facilitates the LLM to gain a general perspective about the functions of the app. The extracted information includes the name of the app and the name of all its activities.

**Page GUI information** provides the semantics of the current page under testing during the interactive process, which facilitates the LLM to capture the current snapshot. We extract the activity name of the page, all the widgets represented by the "text" field or "resource-id" field (the first non-empty one in order), and the widget position of the page. For the position, inspired by the screen reader [83, 88, 103], we first obtain the coordinates of each widget in order from top to bottom and from left to right, and the widgets whose ordinate is below the middle of the page is marked as lower, and the rest is marked as upper.

**Widget information** denotes the micro-level semantics of the GUI page, i.e., the inherent meaning of all its widgets, which facilitates the LLM in providing actionable operational steps related to these widgets. The extracted information includes "text", "hint-text", and "resource-id" field (the first non-empty one in order), "class" field, and "clickable" field. We also extract the information

from nearby widgets to provide a more thorough perspective, which includes the "text" of parent node widgets and sibling node widgets.

WidgetVisits = [{"Income": "2", "Add": "2", "Delete": "3", ...}]

3.1.2 **Dynamic Context Extraction**. Dynamic context relates to the detailed testing progress, which facilitates the LLM being well aware of the process context and making informed decisions. We design an **operation memorizer** to keep the record of this information, i.e., whether and how many times a GUI page has been explored or a widget has been operated, as shown in Table 1.

Specifically, during the iteration, when an operation is conducted, we can obtain the widget information of the operation, and the GUI pages information after the operation, and then the operation memorizer is updated accordingly. In detail, the visit number of the widget is updated through finding the same widget in the operation memorizer by the "text" field and "resource-id" field of the widget. The visit number of the GUI page is updated through finding the same activity in the memorizer with the "ActivityName" field of the page.

## 3.2 Prompt Generation

With the extracted information, we design linguistic patterns to generate prompts for inputting into the LLM.

We first conduct preprocessing for the extracted information, to facilitate the follow-up design. For each static attribute in Table 1, we tokenize it by the underscore and Camel Case [6] (e.g. Capital letters of each word) considering the naming convention in app development and remove the stop words to reduce noise. We then conduct the part-of-speech (POS) tagging with Standford NLP parser [30], and only retain the noun, verb and prepositions for the linguistic patterns.

3.2.1 Linguistic Patterns of Prompt. To design the patterns, each of the five authors is asked to write the prompt sentence following regular prompt template [21, 26, 42], and questions the LLM for generating the operation steps. He/she then checks to what extent the recommended operation is reasonable considering the whole testing process. Each author can access a random-chosen

Table 2: The example of linguistic patterns of prompts and prompt generation rules.

Id	Sample of linguistic patterns/rules	Instantiation							
	Static context patterns: (StaticContext)								
1 2 3	We want to test the 'AppName' App. It has the following activities, including 'Activities'.  The current page is 'ActivityName', it has 'Widgets'. The upper part of the app is 'Position', the lower part is 'Position'.  The widgets which can operation are 'WidgetText / WidgetID'. 'WidgetText / WidgetID' is 'WidgetCategory' which can 'WidgetAction' and its nearby widget is 'NearbyWidget'.	We want to test "Money tracker" App. It has the following activities, including "Main", "AddAccount", "Import", "Setting",  The current page is "Main", it has "Income", "Add", The upper part of the app is "Welcome to, delete,", the lower part of the app is "Income, add,".  The widgets which can operation are "add", "delete", "add" is Button which can click and its nearby widget is "Add account,", "delete" is TextView which can click and its nearby widget is							
	Dynamic context patterns: ⟨DynamicContext⟩								
4	We have tested 'PageVisits' with visit time. We have tested the 'WidgetVisits' on the current page.	We have tested "Main: 7", "About: 2", "Account: 5", with visit time. We have tested the "Add: 1", "Delete: 2", "Edit the order: 1" on the current page.							
	Operation & feedback question patterns: <i>(OperationQuestion)</i>								
5 6	What operation is required?   There is no <i>'WidgetText / WidgetID'</i> on the current page, please reselect.	What operation is required? There is no "Input" on the current page, please reselect.							
	Prompt generation rules								
		ion rules							
1	Start Prompt: $\langle StaticContext \rangle$ [1,2,3] + $\langle OperationQuestion \rangle$ [5]	We want to test "Money tracker" App, It has the following activities, including "Main", "AddAccount", "Exchange", The current page is "Main", it has "Income", "Add" The upper of the app is "Welcome to, delete,", the lower of the app is "Income, add,". The "Income" is Button which can click and its nearby widget is "Please input", What operation is required?							
2	Start Prompt: $\langle StaticContext \rangle [1,2,3] + \langle OperationQuestion \rangle [5]$ Test Prompt: We successfully did the above operation. $\langle DynamicContext \rangle [4] + \langle StaticContext \rangle [2,3] + \langle OperationQuestion \rangle [5]$	We want to test "Money tracker" App, It has the following activities, including "Main", "AddAccount", "Exchange", The current page is "Main", it has "Income", "Add" The upper of the app is "Welcome to, delete,", the lower of the app is "Income, add,". The "Income" is Button which can click and its							

Notes: "[1,2, ..., 5]" means the id of each pattern.

100 apps from Google play, and he/she can obtain the preprocessed static context information and the dynamic context information. After 10 hours of trial, he/she is required to provide the most promising and diversified 20 prompt sentences, which are served as the seeds for designing patterns. With the prompt sentences, the five authors then conduct card sorting [82] and discussion to derive the linguistic patterns. As shown in Table 2, this process comes out with 6 linguistic patterns corresponding with the four sub-types of information in Table 1 and two operation & feedback patterns.

**Pattern related to static context:** We design three patterns to describe the overview of the GUI page currently under testing, respectively corresponding to the app information, page GUI information, and widget information in Table 1.

**Pattern related to dynamic context:** We design one pattern to describe the testing progress with the dynamic context as shown in Table 1.

Pattern related to operation & feedback question: We design two patterns to describe the operational and feedback question. For the operational question, we ask the LLM what operation is required. And for the feedback question, after deciding the previous operation is not applicable (as described in Section 3.3), we inform the LLM that there is no such widget on the current page, and let it re-try.

3.2.2 **Prompt Generation Rules:** Since the designed patterns describe information from different points of view, we combine the patterns from different viewpoints and generate the prompt rules as shown in Table 2. We design three kinds of prompts respectively for starting the test, routine inquiry, and the feedback in case of error occurred. Note that, due to the robustness of the LLM, the generated prompt sentence does not need to follow the grammar completely.

**Test prompt** is the most commonly used prompt for informing the LLM of the current status and query for the next operation. Specifically, we first tell the LLM the dynamic context, i.e., how many times each GUI page and widget has been explored; followed by the static context, i.e., the information about the current GUI page and detailed widget information; then ask the LLM which operation is required.

**Feedback prompt** is used for informing the LLM error occurred and re-try for querying the next operation. Specifically, we first tell LLM its generation operation cannot correspond to the widget on the page; re-provide it the detailed widget information of the page and let the LLM recommend the operation again.

Besides the above two kinds of prompts, we additionally design **start prompt** to start the testing of the app. Different from the test prompt, it provides the LLM with the information about the app with all activities for a global overview; and it does not have dynamic context since the testing just begins. This prompt is only

used once since the LLM can somehow remember this global app information during the testing process [54, 84].

# 3.3 Operation Matching

After inputting the generated prompt, LLM would output a natural language sentence describing the operation steps for the testing, e.g., click the save button. We need to convert the natural language described operation steps to the GUI events (i.e., widgets) of the app to enable it to be automatically executed. This is non-trivial considering the natural language description can be arbitrary, and inherently imprecise. We design a neural matching network to predict which widget can be most likely to be mapped to the operation step. Since training the neural network usually requires a large amount of labeled data, we develop a heuristic-based automated training data generation method to facilitate the model training in Section 3.3.2.

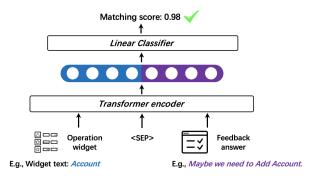


Figure 4: The matching network architecture.

3.3.1 Neural Matching Network. As shown in Figure 4, one input of the neural matching network is the LLM's feedback answer, which is the natural language described operation step  $C_{step}$ , while the other input is the textual information  $C_{text}$  of app's widget. We choose the first non-empty "text", "ID", and "description" fields of the widget in turn. The operation step  $C_{step}$  and the textual information of the widget  $C_{text}$  are concatenated with symbol  $\langle SEP \rangle$ , then input into the pre-trained transformer encoder to generate the hidden state of the text. Finally, the textual hidden state is input into a fully connected layer for obtaining the matching score of the feedback answer and the widget.

During the iterative testing process, each time when the LLM provides the feedback answer, GPTDroid would first separate the answer with "and", "," "to derive the atomic operation step, considering the compound operations can be recommended by the LLM. Then for each atomic operation step, we compute the matching score with all candidate widgets in the current GUI page, and choose the widget with the highest matching score as the target widget. We then use the WidgetAction attribute of the widget in Table 1 to performance the action on the target widget for executing the app. Besides, if the matching score of all widgets is less than 0.5, we determine there is no satisfied widget on the page. It indicates an error occurred in the LLM's feedback answer, and would activate the feedback prompt in the next iteration.

3.3.2 Heuristic-based Training Data Generation. Training such a neural matching network requires a large amount of labeled data

with natural language described operation step and GUI event, e.g., press the back button and corresponding event click back button in the app. However, there is no such open dataset, and collecting it from scratch is time- and effort-consuming. Meanwhile, by examining the operation step of LLM's output, we observe that they tend to follow certain linguistic patterns. This motivates us in developing a heuristic-based automated training data generation method for collecting the satisfied training data.

The primary idea is that for each interactive widget in a GUI page, we heuristically generate the operation step which can operate on it for transitioning to the next state; meanwhile, generate the negative data instances between the operation step and the irrelevant widgets. To derive the heuristic rules, the five authors examine 400 operation steps of LLM's output, summarizing the linguistic patterns for writing the operation steps. We have summarized 31 non-repetitive operation descriptions and their variants, such as click, enter, press, etc. We uploaded the complete table to our website?? For each iterative widget, we randomly generate three different operation descriptions to serve as the positive data instances in the training data. For the negative data instances, we follow the hard negative mining strategy [58] to enhance the discriminability of the model.

# 3.4 Implementation

Our GPTDroid is implemented as a fully automated GUI app testing tool, which uses or extends the following tools: VirtualBox [12] and the Python library pyvbox [8] for running and controlling the Android-x86 OS, Android UIAutomator [86] for extracting the view hierarchy file, and Android Debug Bridge (ADB) [1] for interacting with the app under test (Section 3.1).

For the LLM (Section 3.2), we use the pre-trained GPT-3 model which was released on the OpenAI website<sup>3</sup>. The basic model of GPT-3 is the *text-davinci-003* model which is extremely powerful and good at answering questions. We build our neural matching network model (Section 3.3) based on PyTorch [7] and Sentence Transformers [74]. The text processing module is loaded with Distil-Bert [79], a 12-layer transformer-based pre-training model. We use AdamW as the optimizer, BCEWithLogitsLoss as the loss function, and train the model with the batch size set to 20.

## 4 EFFECTIVENESS EVALUATION

In order to verify the performance of GPTDroid, we evaluate it by investigating the activity coverage (RQ1) and the number of detected bugs (RQ2). In addition, we present the performance of operation matching in RQ3. Note that, this section utilizes the previously-detected bugs in the app's repositories to demonstrate the effectiveness of our approach, and the next section will further evaluate the usefulness of GPTDroid in detecting new bugs.

## 4.1 Experimental Setup

The experimental dataset comes from two sources. The first is from the apps in Themis benchmark [85], which contains 20 open-source apps with 34 bugs in GitHub. Considering the small number of apps in the benchmark, we collect a second dataset following similar procedures as the benchmark.

 $<sup>^3</sup> https://beta.openai.com/docs/models/gpt-3\\$ 

In detail, we crawl the 50 most popular apps of each category from Google Play [4], and we keep the ones with at least one update after May. 2022, resulting in 317 apps in 12 Google Play categories. Then, we use 9 common-used and state-of-the-art automated GUI testing tools (details are in Section 4.2) to run these apps in turn to ensure that they work properly. We then filter out the unusable apps by the following criteria: (1) UIAutomator [86] can't obtain the view hierarchy file; (2) they would constantly crash on the emulator; (3) one or more tools can't run on them; (4) The registration and login functions cannot be skipped with scripts [35, 61, 85]; (5) They don't have issue records or pull requests on GitHub.

There are 66 apps (with 93 bugs) remaining for this effectiveness evaluation. Note that, same as the benchmark, all bugs are crash bugs. Specifically, for each app, we select the version in which the bugs are confirmed by developers (merged GitHub pull requests) as our experimental data, following the practice of the benchmark. The details of all 86 experimental apps (20 + 66) and related bugs are shown in Table 3.

Table 3: Dataset of effectiveness evaluation.

Statistics	#Activities	#Bugs	#Download	#Update
Min	7	1	50K+	05/22
Max	21	9	50M+	11/22
Median	10	3	1M+	-
Average	9	1.5	10M+	-
All	790	129	-	-

Note that, there are 101 apps that are filtered out for effectiveness evaluation, yet can successfully run with our proposed approach. We apply them to the manual prompt generation in Section 3.2.1 and heuristic training data generation in Section 3.3.2. And this ensures there is no overlapping between the apps in approach design and evaluation.

We employ activity coverage and the number of detected bugs, which are widely used metrics for evaluating GUI testing [15, 43, 57]. We also present the number of covered activities and widgets which are also commonly-used metrics [54, 71, 85] in Table 4.

# 4.2 Baselines

To demonstrate the advantage of GPTDroid, we compare it with 9 common-used and state-of-the-art baselines. We roughly divide the GUI testing approaches into random-/rule-based methods, model-based methods, and learning-based methods, to facilitate understanding.

For random-/rule-based methods, we use Monkey [34] and Droid-bot [53]. For model-based methods, we use Stoat [84], Ape [41], Fastbot [20], ComboDroid [90], TimeMachine [35]. For learning-based methods, we use Humanoid [54] and Q-testing [71].

We deploy the baselines and our approach on a 64-bit Ubuntu 18.04 machine (64 cores, AMD 2990WX CPU, and 128GB RAM) and evaluate them on Google Android 7.1 emulators (API level 25). Each emulator is configured with 2GB RAM, 1GB SDCard, 1GB internal storage, and X86 ABI image. Different types of external files (including PNGs / MP3s / PDFs / TXTs / DOCXs) are stored on the SDCard to facilitate file access from apps. Following common practice [41, 53], we registered separate accounts for each bug that requires login and wrote the login scripts, and during testing reset

the account data before each run to avoid possible interference. In order to ensure fair and reasonable use of resources, we set up the running time of each tool in one app to 30 minutes, which is widely used in other GUI testing studies [36, 41, 53, 85]. We run each tool three times and obtain the highest performance to mitigate potential bias.

# 4.3 Results and Analysis

4.3.1 **Performance of Activity Coverage (RQ1)**. Table 4 shows the number of covered widgets, number of covered activities, and average activity coverage of GPTDroid and the baselines. We can see that GPTDroid covers far more widgets and activities than the baselines, and the average activity coverage achieves 71% across the 86 apps. It is 32% (0.71 vs. 0.54) higher even compared with the best baseline (TimeMachine). This indicates the effectiveness of GPTDroid in covering more activities and widgets, thus bring higher confidence to the app quality and potentially uncovering more bugs.

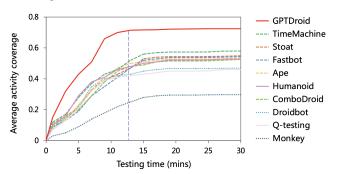


Figure 5: Activity coverage with varying time (RQ1).

Figure 5 additionally demonstrates the average activity coverage with varying time. We can see that, in every time point, GPTDroid achieves higher activity coverage than the baselines, and it achieves high coverage within about 13 minutes. This again indicates the effectiveness and efficiency of GPTDroid in covering more activities with less time, which is valuable considering the testing budget.

Among the baselines, the model-based and learning-based approaches have relatively higher performance. Yet the model-based approaches cannot capture the GUI semantic information and the exploration could not well understand the inherent business logic of the app. In addition, existing learning-based approaches only use few context information for guiding the exploration, and the learners only have limited intelligence restricted by the model architecture and amount of labeled training data.

We further analyze the potential reasons for the uncovered cases. First, some widgets or inputs do not have meaningful "text" or "resource-id", which hinders the approach of effectively understanding the GUI page. Second, some app requires specific operations, e.g., database connection, long press and drag widgets to a fixed location, which is difficult if not impossible to be automatically achieved.

4.3.2 **Performance of Bug Detection (RQ2).** Figure 6 shows the overall number of detected bugs of GPTDroid and baselines with varying times. GPTDroid detects 72 bugs for the 86 apps, 36%

Metric	Random-/r	ule-based		Model-based			Learning-based			
Metric	Monkey	Droidbot	Stoat	Ape	Fastbot	ComboDroid	TimeMachine	Humanoid	Q-testing	GPTDroid
#Widgets	351	893	1337	1582	1437	1388	1701	1453	1398	1989
#Activities	104	269	333	370	391	383	401	340	323	477
Ave activity coverage	0.16	0.34	0.44	0.51	0.56	0.53	0.57	0.40	0.45	0.71

Table 4: Performance of activity coverage (RQ1).

(72 vs. 53) higher than the best baseline (Stoat). We also compare the similarities and differences of the bugs between Stoat and our approach, and the results show that its detected bugs are a subset of our detected bugs. This indicates the effectiveness of GPTDroid in detecting bugs and helps to ensure app quality.

We can also see that, in every time point, GPTDroid detects more bugs than the baselines, and reaches the highest value in about 17 minutes, saving 35% (17 vs. 26) of the testing time compared with the best baseline (also with more detected bugs). This again proves its effectiveness and efficiency of GPTDroid, which is valuable for saving more time for the follow-up bug fixing. We will conduct a further discussion about the reason behind the superior performance in Section 4.3.3.

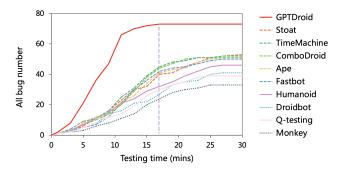


Figure 6: Bug detection with varying time (RQ2).

4.3.3 **Performance of Operation Matching (RQ3).** When testing the experimental apps in RQ1, we randomly choose 1,000 LLM's feedback answer, i.e., natural language described operations outputted by LLM, and evaluate whether the operations can successfully match the correct GUI widgets. Specifically, the two authors follow the principle of open-coding, analyze the feedback answer and find the matching widget on the current page. For the inconsistent labeling, the third author will judge until all the authors reach an agreement.

Results show that the operation matching can achieve 0.96 accuracy, indicating that most of the LLM's feedback answers can be accurately matched to the GUI widget. This lays a solid foundation for the high activity coverage of our GPTDroid.

#### 5 USEFULNESS EVALUATION

## 5.1 Experimental Setup

This section further evaluates the usefulness of GPTDroid in detecting new crash bugs. We employ a similar experimental setup to the previous section. To make it brief, we only compare the best baselines for bug detection in the last section, i.e., Droidbot, Stoat

and Humanoid, the best one from each type of methods (random/rule-based methods, model-based methods and learning-based methods).

We begin with the most popular and recently updated 317 apps from 12 categories as in the previous section. Then we reuse the five criteria in the previous section for filtering the unusable apps. Differently, we loosen the criteria 5, which only requires the app to have ways for bug reporting, since the issue records or pull requests are not mandatory in this section. This results in 216 apps for our usefulness evaluation. Note that, this section aims at evaluating whether GPTDroid can detect new bugs on these apps, thus the overlap between the apps of this section and the previous section is allowed.

We use the same hardware and the software configurations as the previous evaluation section. When the crash bugs are detected, we report them to the app development team through online issue reports or email.

Table 5: confirmeded or fixed bugs

Id	APP Name	Category	Download	Status
1	PerfectPia	Music	50M+	confirmed
2	MusicPlayer	Music	50M+	confirmed
3	NoxSecu	Tool	10M+	fixed
4	Degoo	Tool	10M+	fixed
5	Proxy	Tool	10M+	confirmed
6	Secure	Tool	10M+	confirmed
7	Thunder	Tool	10M+	confirmed
8	ApowerMir	Tool	5M+	confirmed
9	MediaFire	Product	5M+	confirmed
10	Postegro	Commun	500K+	fixed
11	Deezer MP	Music	500K+	fixed
12	MTG	Utilities	500K+ 500K+	fixed
13	OFF	Health	500K+ 500K+	confirmed
13 14	Yucata	Tool	500K+ 500K+	confirmed
		1		
15	ClassySha	Tool	500K+	confirmed
16	Linphone	Commun	500K+	confirmed
17	Paytm	Finance	100K+	confirmed
18	Transdroid	Tool	100K+	confirmed
19	Transistor	Music	10K+	fixed
20	Onkyo	Music	10K+	fixed
21	Democracy	News	10K+	confirmed
22	NewPipe	Media	10K+	confirmed
23	LessPass	Product	10K+	confirmed
24	CEToolbox	Medical	10K+	confirmed
25	OSM	Health	10K+	fixed

# 5.2 Results and Analysis

For the 216 apps, GPTDroid detects 135 bugs involving 115 apps, of which 48 bugs involving 39 apps are newly-detected bugs. Furthermore, these new bugs are not detected by the three baselines. We submit these 48 bugs to developers, and 25 of them have been fixed/confirmed so far (8 fixed and 17 confirmed), while the remaining are still pending (none of them is rejected). This further indicates

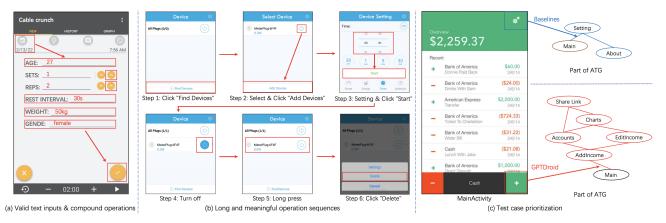


Figure 7: Examples of our finding.

the effectiveness of our proposed GPTDroid in bug detection. Due to space limit, Table 5 presents these fixed/confirmed bugs, and the full lists can be found on our website??

We further analyze the details of our found bugs, and 17 of them involve multiple text inputs or compound operations. Besides, we also observe that there are 11 bugs with more than 20 operation steps before triggering the bug, counting from the *MainActivity* page, which indicates the ability of GPTDroid in testing deeper features. Furthermore, we find at least 28 bugs related to the main business logic of the app, for example, a bug about the health data statistics is revealed for a digital health app.

## 6 DISCUSSION

Despite the superior performance of GPTDroid in the last section, it is still unclear the reason behind it. To fully understand the testing capability of the LLM, we carry out a qualitative study by investigating the cases in which our model outperforms baselines. We summarize four kinds of capabilities including low-level (i.e., valid text input, and compound actions) and high-level ones (i.e., long meaningful test trace, and test case prioritization). These findings pave the way for further research in this area.

Valid text inputs. Our approach can automatically fill in valid text content to the input widget which is essentially the key for passing the page as seen in Figure 7 (a). More importantly, our model can generate semantic text input (e.g., income, date, ID number, searching item, etc) accordingly. Besides single text input, it can also successfully fill in multiple input widgets at the same time which are correlated to each other like the departure and arrival cities and dates in the flight booking app. It may be due to GPT-3's language generation capabilities [19] which was well learned during the training phase.

Compound actions. GPTDroid can conduct complex compound operations guided by the LLM. As shown in the above example (Figure 7 (a)), to add the "Cable crunch" information, it first inputs the text, selects the date, sets the "SETS" and "REPS" by clicking the upper or lower button, then click the submit button in the lower right corner. Since there are many tutorials or bug reports on the web including natural-language descriptions of how certain actions can lead to specific outcomes which may be adopted into GPT-3's training corpus. It may help with a better understanding of the

causal relationships between actions and outcomes, resulting in LLM's compound action ability in GUI testing.

Long meaningful test trace. GPTDroid can automatically generate the test cases with a long sequence of operations which together accomplish a business logic of the app, and this is quite important for covering the app features and ensuring its quality. As shown in Figure 7 (b), in SmartMeter app [9], to test a commonlyused app feature "delete equipment", the automatic tool first needs to click "find devices" in the device page, then select a device (Bluetooth is turned on and there are candidate devices) and click "add device" for adding it in the device page; input the related information and click "start" to start the device; then turn off this device in the device page, long press it and click "delete" from the pop-up menu. Only with this long sequence of operations which touches the "deleting equipment" feature, a crash can be revealed. It may be because of GPT-3's exposure to tutorial or bug reports which contain step-by-step instructions or descriptions of how to trigger a certain feature or reproduce a certain bug [85, 89] in the training corpus. Therefore, provided with low-level semantic information (i.e., current GUI page) and high-level testing history, GPTDroid can capture long-term dependencies among GUI pages for generating long meaningful exploration sequences.

Test case prioritization. We also observe that GPTDroid usually prioritizes testing the "important" widgets, which is valuable for reaching a higher activity coverage and covering more key activities with relatively less time. As shown in Figure 7(c), in the Main page of the Moni app [5], the baseline tools tend to first click the "Setting" button following the exploration order from upper to lower, which leads the testing easily trapped into the setting page cycle. Our GPTDroid chooses to first test the "AddIncome" activity, i.e., click the "add" button, which is based on the semantics of the GUI page and app information, and can quickly explore the activities related to the key features of the app. This may be because that GPT-3's training corpus contains a wide variety of software-related information, including user manuals, release notes, and app/software descriptions where developers often highlight the most important features at the front.

## 7 RELATED WORK

# 7.1 Mobile App GUI and GUI Analysis.

GUI provides a visual bridge between apps and users, and is drawing increasing attention [46, 69, 81]. The graphical user interface (GUI) is the most important type of UI for most mobile apps, where apps present content and actionable widgets on the screen and users interact with the widgets using actions such as clicks, swipes, and text inputs. GUI is an indispensable part of software on most major platforms including Android. Analyzing the app's GUI is of great interest to many researchers and practitioners. Related studies included automatic GUI search [17, 22, 75], GUI code generation [23, 65, 70], GUI changes detection and summarization [66, 67], GUI design [25, 95], etc. Gao et al. [39] and Li et al. [52] analyzed the possible problems in UI rendering, and developed automatic approaches to detect them. Nayebi et al. [68] and Holzinger et al. [45] found that different resolutions of mobile devices have brought challenges in Android app design and implementation. Huang et al. [47] and Rubin et al. [77] proposed to detect stealthy behaviors in Android apps by comparing the actual behaviors with the UI. Chen et al. [28] introduced a machine learning-based method to extract UI skeletons from UI images, in order to facilitate GUI development. In human-computer interaction research, software GUI is mainly used to mine UI design practices [14, 51] and interaction patterns [32] at scale. The mined knowledge can further be used to guide UI and UX (user experience) design. Due to these challenges, how to reasonably summarize the app is a problem to be solved. Our focus is to use the reasonable natural language to describe the GUI information of the app page, so that the large language model can understand the GUI information.

## 7.2 Automated GUI testing

To ensure the quality of mobile apps, many researchers study the automatic generation of large-scale test scripts to test apps [93]. Monkey [34] is the popular random-based automated GUI testing tool, which emits pseudo-random streams of UI events and some system events. It is easy to use and compatible with different Android versions. However, the random-based testing strategy cannot formulate a reasonable testing path according to the characteristics of the app, resulting in low test coverage. To improve the test coverage, researchers propose model-based [35, 64, 96, 97, 101] automated GUI testing methods, design corresponding models through the research and analysis of large-scale apps. Sapienz [63] used genetic algorithms as the model and Stoat [84] used the stochastic model learned from an app to optimize test suite generation. Ape [41] used the runtime information to dynamically evolve its abstraction criterion via a decision tree and generated UI events via a random and greedy depth-first state exploration strategy. ComboDroid [90] obtained such use cases either from humans or automatically generates from a GUI model constructed by GUI exploration and analyzed the data flow between obtained use cases, and combined them to generate final tests. Although model-based automated GUI testing tools can improve test coverage, the coverage is still low because it does not consider the semantic information of the app's GUI and Page. Researchers further proposed human-like testing strategies and designed learning-based automated GUI testing methods. Humanoid [54] used a deep neural network model that predicts which

UI elements on the current GUI page are more likely to be interacted with by users and how to interact with them. Q-testing [71] used a reinforcement learning-based method to compare GUI pages and give rewards. These rewards are used and iteratively updated to guide the testing to cover more functionalities of apps. Although the learning-based approach can improve the test coverage by learning a large number of interactive processes or using the idea of reinforcement learning. However, it is still unable to better understand the semantic information of the page and plan the path according to the actual situation of the app, and is greatly affected by the training data. This study aims at proposing a more effective approach to generate human-like actions for thoroughly and more effectively testing the app, and accomplishing it with LLM.

# 7.3 Large Language Model

Recently, there has been a great success of pre-trained Large Language Models (e.g., RoBERTa [59], GPT-3 [19], PaLM [27], OPT [102]) in a variety of NLP tasks. Due to the large amounts of available pre-training data from the internet, research shows that LLMs can already be used for very specific downstream tasks through the new paradigm "pre-train, prompt and prediction" [56] without any finetuning of special data sets. This paradigm for LLM was widely used in many works and achieved state-of-the-art performance on downstream tasks [33, 48]. The core of this paradigm is to use prompt engineering [18, 21, 40, 56], where a natural language description of the task is provided to the LLM. Considering the powerful performance of LLM, researchers try to use LLM to solve relevant tasks in the field of software engineering. Supported by code naturalness [44], researchers applied the LLMs to code writing in different programming languages [24, 37, 38, 94]. Huang et al. [48] used LLM for the type inference in statically-typed partial Code. In testing, LLMFuzz [33] used LLMs to generate input programs for fuzzing Deep Learning libraries. Xia et al. [91, 92] applied LLM to automatic program repair to improve the accuracy of the generated repair patches. Austin et al. [16] and Jain et al. [49] used LLM for program synthesis in general-purpose programming languages. This paper opens a new dimension for automated GUI testing by formulating it as a Q&A task with LLM.

## 8 CONCLUSION

As one of the most important quality assurance activities for mobile apps, automated GUI testing has made much progress, yet still suffers from low activity coverage and may miss critical bugs. This paper aims at generating human-like actions to facilitate app testing more thoroughly and effectively. Inspired by the success of LLM like ChatGPT, we formulate the GUI testing problem as a Q&A task and propose GPTDroid. It extracts the static and dynamic context of the current GUI page, encodes them into prompt question to ask the LLM, decodes the LLM's feedback answer into actionable operations to execute the app, and iterates the whole process. Results on 86 popular apps demonstrate that GPTDroid can achieve 71% activity coverage, with 32% higher than the best baseline, and can detect 36% more bugs with faster speed than the best baseline. GPTDroid also detects 48 new bugs on Google Play with 25 of them being confirmed/fixed, with the remaining pending. The capability of GPTDroid in generating semantic text input and compound actions, guiding to explore the long meaningful test trace, and prioritizing test cases, further proves the effectiveness and human-like aspects of our proposed GPTDroid.

In the future, we will fine-tune the LLM to improve the performance and conduct a systematic study to understand reasons why LLM can help the GUI testing.

## REFERENCES

- 2023. Android Debug Bridge (adb). https://developer.android.com/studio/ command-line/adb.html#forwardports.
- [2] 2023. Android development. http://developer.android.com/reference/android.
- [3] 2023. App Store. https://www.apple.com.cn/app-store/.
- [4] 2023. Google play. https://play.google.com/store/apps/.
- [5] 2023. Moni. https://play.google.com/store/apps/details?id=Moni.
- [6] 2023. pascal case. https://en.wikipedia.org/wiki/Camel\_case.
- [7] 2023. Pytorch. https://pytorch.org/.
- [8] 2023. pyvbox. https://pypi.org/project/pyvbox/.
- [9] 2023. SmartMeter. https://play.google.com/store/apps/details?id=SmartMeter.
- [10] 2023. Talkback. https://play.google.com/store/apps/details?id=talkback.
- [11] 2023. View hierachy. https://developer.android.google.cn/topic/performance/.
- [12] 2023. virtualbox. https://www.virtualbox.org/.
- [13] 2023. VoiceOver. https://play.google.com/store/apps/details?id=voiceover.
- [14] Khalid Alharbi and Tom Yeh. 2015. Collect, decompile, extract, stats, and diff: Mining design pattern changes in Android apps. In Proceedings of the 17th international conference on human-computer interaction with mobile devices and services. 515–524.
- [15] Yauhen Leanidavich Arnatovich, Lipo Wang, Ngoc Minh Ngo, and Charlie Soh. 2018. Mobolic: An automated approach to exercising mobile application GUIs using symbiosis of online testing technique and customated input generation. Software: Practice and Experience 48, 5 (2018), 1107–1142.
- [16] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. arXiv preprint arXiv:2108.07732 (2021).
- [17] Farnaz Behrang, Steven P Reiss, and Alessandro Orso. 2018. GUIfetch: supporting app design and development through GUI search. In Proceedings of the 5th International Conference on Mobile Software Engineering and Systems. ACM, 236–246
- [18] Gwern Branwen. 2020. Gpt-3 creative fiction. https://www.gwern.net/GPT-3.
- [19] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. Advances in neural information processing systems 33 (2020), 1877–1901.
- [20] Tianqin Cai, Zhao Zhang, and Ping Yang. 2020. Fastbot: A Multi-Agent Model-Based Test Generation System Beijing Bytedance Network Technology Co., Ltd.. In Proceedings of the IEEE/ACM 1st International Conference on Automation of Software Test. 93–96.
- [21] Andrew Cantino. 2016. Prompt Engineering Tips and Tricks with GPT-3. https://blog.andrewcantino.com/blog/2021/04/21/prompt-engineering-tipsand-tricks/.
- [22] Chunyang Chen, Sidong Feng, Zhenchang Xing, Linda Liu, Shengdong Zhao, and Jinshui Wang. 2019. Gallery DC: Design Search and Knowledge Discovery through Auto-created GUI Component Gallery. CSCW (2019).
- [23] Chunyang Chen, Ting Su, Guozhu Meng, Zhenchang Xing, and Yang Liu. 2018. From ui design image to gui skeleton: a neural machine translator to bootstrap mobile gui implementation. In ICSE. ACM.
- [24] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374 (2021).
- [25] Qiuyuan Chen, Chunyang Chen, Safwat Hassan, Zhengchang Xing, Xin Xia, and Ahmed E Hassan. 2021. How Should I Improve the UI of My App? A Study of User Reviews of Popular Apps in the Google Play. TOSEM 30, 3 (2021), 1–38.
- [26] Xiang Chen, Ningyu Zhang, Xin Xie, Shumin Deng, Yunzhi Yao, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2022. Knowprompt: Knowledge-aware prompttuning with synergistic optimization for relation extraction. In *Proceedings of the ACM Web Conference* 2022. 2778–2788.
- [27] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. arXiv preprint arXiv:2204.02311 (2022).
- [28] Chen Chunyang, Feng Sidong, Liu Zhengyang, Xing Zhenchang, Zhao Sheng-dong, and Liu Linda. 2020. From Lost to Found: Discover Missing UI Design Semantics through Recovering Missing Tags. In CSCW.

- [29] Eliane Collins, Arilo Neto, Auri Vincenzi, and José Maldonado. 2021. Deep reinforcement learning based Android application GUI testing. In Proceedings of the XXXV Brazilian Symposium on Software Engineering. 186–194.
- [30] Marie-Catherine De Marneffe and Christopher D Manning. 2008. The Stanford typed dependencies representation. In Coling 2008: proceedings of the workshop on cross-framework and cross-domain parser evaluation. 1–8.
- [31] Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibschman, Daniel Afergan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. 2017. Rico: A Mobile App Dataset for Building Data-Driven Design Applications. In UIST.
- [32] Biplab Deka, Zifeng Huang, and Ranjitha Kumar. 2016. ERICA: Interaction mining mobile apps. In Proceedings of the 29th annual symposium on user interface software and technology. 767–776.
- [33] Yinlin Deng, Chunqiu Steven Xia, Haoran Peng, Chenyuan Yang, and Lingming Zhang. 2022. Fuzzing Deep-Learning Libraries via Large Language Models. arXiv preprint arXiv:2212.14834 (2022).
- [34] Android Developers. 2012. Ui/application exerciser monkey.
- [35] Zhen Dong, Marcel Böhme, Lucia Cojocaru, and Abhik Roychoudhury. 2020. Time-travel testing of android apps. In ICSE. IEEE.
- [36] Lingling Fan, Ting Su, Sen Chen, Guozhu Meng, Yang Liu, Lihua Xu, Geguang Pu, and Zhendong Su. 2018. Large-scale analysis of framework-specific exceptions in android apps. In ICSE. IEEE, 408–419.
- [37] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. Codebert: A pre-trained model for programming and natural languages. EMNLP (2020).
- [38] Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wen-tau Yih, Luke Zettlemoyer, and Mike Lewis. 2022. Incoder: A generative model for code infilling and synthesis. arXiv preprint arXiv:2204.05999 (2022).
- [39] Yi Gao, Yang Luo, Daqing Chen, Haocheng Huang, Wei Dong, Mingyuan Xia, Xue Liu, and Jiajun Bu. 2017. Every pixel counts: Fine-grained UI rendering analysis for mobile applications. In *IEEE INFOCOM 2017-IEEE Conference on Computer Communications*. IEEE, 1–9.
- [40] Songwei Ge and Devi Parikh. 2021. Visual Conceptual Blending with Large-scale Language and Vision Models. arXiv preprint arXiv:2106.14127 (2021).
- [41] Tianxiao Gu, Chengnian Sun, Xiaoxing Ma, Chun Cao, Chang Xu, Yuan Yao, Qirun Zhang, Jian Lu, and Zhendong Su. 2019. Practical GUI testing of Android applications via model abstraction and refinement. In ICSE. IEEE.
- [42] Yuxian Gu, Xu Han, Zhiyuan Liu, and Minlie Huang. 2021. Ppt: Pre-trained prompt tuning for few-shot learning. (2021).
- [43] Yuyu He, Lei Zhang, Zhemin Yang, Yinzhi Cao, Keke Lian, Shuai Li, Wei Yang, Zhibo Zhang, Min Yang, Yuan Zhang, et al. 2020. TextExerciser: feedbackdriven text input exercising for android applications. In 2020 IEEE Symposium on Security and Privacy (SP). IEEE, 1071–1087.
- [44] Abram Hindle, Earl T Barr, Mark Gabel, Zhendong Su, and Premkumar Devanbu. 2016. On the naturalness of software. Commun. ACM 59, 5 (2016), 122–131.
- [45] Andreas Holzinger, Peter Treitler, and Wolfgang Slany. 2012. Making apps useable on multiple different mobile platforms: On interoperability for business application development on smartphones. In *International Conference on Availability, Reliability, and Security*. Springer, 176–189.
- [46] Huaxun Huang, Ming Wen, Lili Wei, Yepang Liu, and Shingchi Cheung. 2021. Characterizing and Detecting Configuration Compatibility Issues in Android Apps. In ASE.
- [47] Jianjun Huang, Xiangyu Zhang, Lin Tan, Peng Wang, and Bin Liang. 2014. Asdroid: Detecting stealthy behaviors in android applications by user interface and program behavior contradiction. In ICSE. 1036–1046.
- [48] Qing Huang, Zhiqiang Yuan, Zhenchang Xing, Xiwei Xu, Liming Zhu, and Qinghua Lu. 2022. Prompt-tuned Code Language Model as a Neural Knowledge Base for Type Inference in Statically-Typed Partial Code. In ICSE.
- [49] Naman Jain, Skanda Vaidyanath, Arun Iyer, Nagarajan Natarajan, Suresh Parthasarathy, Sriram Rajamani, and Rahul Sharma. 2022. Jigsaw: Large language models meet program synthesis. In ICSE.
- [50] Pingfan Kong, Li Li, Jun Gao, Kui Liu, Tegawendé F Bissyandé, and Jacques Klein. 2018. Automated testing of android apps: A systematic literature review. IEEE Transactions on Reliability 68, 1 (2018), 45–66.
- [51] Ranjitha Kumar, Arvind Satyanarayan, Cesar Torres, Maxine Lim, Salman Ahmad, Scott R Klemmer, and Jerry O Talton. 2013. Webzeitgeist: design mining the web. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. 3083–3002
- [52] Wenjie Li, Yanyan Jiang, Chang Xu, Yepang Liu, Xiaoxing Ma, and Jian Lü. 2019. Characterizing and detecting inefficient image displaying issues in Android apps. In 2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER). IEEE, 355–365.
- [53] Yuanchun Li, Ziyue Yang, Yao Guo, and Xiangqun Chen. 2017. Droidbot: a lightweight ui-guided test input generator for android. In ICSE. IEEE.
- [54] Yuanchun Li, Ziyue Yang, Yao Guo, and Xiangqun Chen. 2019. Humanoid: a deep learning-based approach to automated black-box Android app testing. In 2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE). IEEE, 1070–1073.

- [55] Jinzhi Liao, Xiang Zhao, Jianming Zheng, Xinyi Li, Fei Cai, and Jiuyang Tang. 2022. PTAU: Prompt Tuning for Attributing Unanswerable Questions. In Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval. 1219–1229.
- [56] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *Comput. Surveys* 55, 9 (2023), 1–35.
- [57] Peng Liu, Xiangyu Zhang, Marco Pistoia, Yunhui Zheng, Manoel Marques, and Lingfei Zeng. 2017. Automatic text input generation for mobile testing. In ICSE. IEFE
- [58] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. 2016. Ssd: Single shot multibox detector. In European conference on computer vision. Springer, 21–37.
- [59] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettle-moyer, and V. Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. (2019). http://arxiv.org/abs/1907.11692
- [60] Li Lucy and David Bamman. 2021. Gender and representation bias in GPT-3 generated stories. In Proceedings of the Third Workshop on Narrative Understanding. 48–55.
- [61] Zhengwei Lv, Chao Peng, Zhao Zhang, Ting Su, Kai Liu, and Ping Yang. 2022. Fastbot2: Reusable Automated Model-based GUI Testing for Android Enhanced by Reinforcement Learning. In ICSE.
- [62] Aravind Machiry, Rohan Tahiliani, and Mayur Naik. 2013. Dynodroid: An input generation system for android apps. In Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering. 224–234.
- [63] Ke Mao, Mark Harman, and Yue Jia. 2016. Sapienz: Multi-objective automated testing for Android applications. In Proceedings of the 25th International Symposium on Software Testing and Analysis. 94–105.
- [64] Nariman Mirzaei, Joshua Garcia, Hamid Bagheri, Alireza Sadeghi, and Sam Malek. 2016. Reducing combinatorics in GUI testing of android applications. In ICSE, IEEE.
- [65] Kevin Moran, Carlos Bernal-Cárdenas, Michael Curcio, Richard Bonett, and Denys Poshyvanyk. 2018. Machine Learning-Based Prototyping of Graphical User Interfaces for Mobile Apps. arXiv preprint arXiv:1802.02312 (2018).
- [66] Kevin Moran, Boyang Li, Carlos Bernal-Cárdenas, Dan Jelf, and Denys Poshyvanyk. 2018. Automated reporting of GUI design violations for mobile apps. In ICSE. ACM.
- [67] Kevin Moran, Cody Watson, John Hoskins, George Purnell, and Denys Poshyvanyk. 2018. Detecting and Summarizing GUI Changes in Evolving Mobile Apps. ASE (2018).
- [68] Fatih Nayebi, Jean-Marc Desharnais, and Alain Abran. 2012. The state of the art of mobile application usability evaluation. In CCECE. IEEE, 1–4.
- [69] Michael Nebeling, Maximilian Speicher, and Moira C. Norrie. 2013. W3touch: metrics-based web page adaptation for touch. In 2013 ACM SIGCHI Conference on Human Factors in Computing Systems, CHI '13, Paris, France, April 27 - May 2, 2013. ACM, 2311–2320. https://doi.org/10.1145/2470654.2481319
- [70] Tuan Anh Nguyen and Christoph Csallner. 2015. Reverse engineering mobile application user interfaces with remaui (t). In ASE. IEEE, 248–259.
- [71] Minxue Pan, An Huang, Guoxin Wang, Tian Zhang, and Xuandong Li. 2020. Reinforcement learning based curiosity-driven testing of Android applications. In Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis. 153–164.
- [72] Fabiano Pecorelli, Gemma Catolino, Filomena Ferrucci, Andrea De Lucia, and Fabio Palomba. 2022. Software testing and Android applications: a large-scale empirical study. Empirical Software Engineering 27, 2 (2022), 1–41.
- [73] Chao Peng, Zhao Zhang, Zhengwei Lv, and Ping Yang. 2022. MUBot: Learning to Test Large-Scale Commercial Android Apps like a Human. In 2022 IEEE International Conference on Software Maintenance and Evolution (ICSME). IEEE, 543–552.
- [74] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics. http://arxiv.org/abs/1908.10084
- [75] Steven P Reiss, Yun Miao, and Qi Xin. 2018. Seeking the user interface. Automated Software Engineering 25, 1 (2018), 157–193.
- [76] Andrea Romdhana, Alessio Merlo, Mariano Ceccato, and Paolo Tonella. 2022. Deep reinforcement learning for black-box testing of android apps. TOSEM (2022).
- [77] Julia Rubin, Michael I Gordon, Nguyen Nguyen, and Martin Rinard. 2015. Covert communication in mobile applications (t). In 2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE). IEEE, 647–657.
- [78] Konstantin Rubinov and Luciano Baresi. 2018. What are we missing when testing our android apps? *Computer* 51, 4 (2018), 60–68.
- [79] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. arXiv preprint arXiv:1910.01108 (2019).

- [80] Mike Sharples. 2022. Automated Essay Writing: An AIED Opinion. International Journal of Artificial Intelligence in Education (2022), 1–8.
- [81] Maximilian Speicher, Andreas Both, and Martin Gaedke. 2015. S.O.S.: Does Your Search Engine Results Page (SERP) Need Help?. In CHI. ACM. https://doi.org/10.1145/2702123.2702568
- [82] Donna Spencer. 2009. Card sorting: Designing usable categories. Rosenfeld Media.
- [83] Abigale Stangl, Nitin Verma, Kenneth R Fleischmann, Meredith Ringel Morris, and Danna Gurari. 2021. Going Beyond One-Size-Fits-All Image Descriptions to Satisfy the Information Wants of People Who are Blind or Have Low Vision. In The 23rd International ACM SIGACCESS Conference on Computers and Accessibility. 1–15.
- [84] Ting Su, Guozhu Meng, Yuting Chen, Ke Wu, Weiming Yang, Yao Yao, Geguang Pu, Yang Liu, and Zhendong Su. 2017. Guided, stochastic model-based GUI testing of Android apps. In Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering. 245–256.
- [85] Ting Su, Jue Wang, and Zhendong Su. 2021. Benchmarking automated GUI testing for Android against real-world bugs. In Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering. 119–130.
- [86] UIAutomator. 2021. Python wrapper of Android uiautomator test tool. https://github.com/xiaocong/uiautomator.
- [87] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. Advances in neural information processing systems (2017).
- [88] Bryan Wang, Gang Li, Xin Zhou, Zhourong Chen, Tovi Grossman, and Yang Li. 2021. Screen2words: Automatic mobile UI summarization with multimodal learning. In The 34th Annual ACM Symposium on User Interface Software and Technology. 498–510.
- [89] Jue Wang, Yanyan Jiang, Ting Su, Shaohua Li, Chang Xu, Jian Lu, and Zhendong Su. 2022. Detecting non-crashing functional bugs in Android apps via deep-state differential analysis. In ESE.
- [90] Jue Wang, Yanyan Jiang, Chang Xu, Chun Cao, Xiaoxing Ma, and Jian Lu. 2020. Combodroid: generating high-quality test inputs for android apps via use case combinations. In ICSE. 469–480.
- [91] Chunqiu Steven Xia, Yuxiang Wei, and Lingming Zhang. 2022. Practical Program Repair in the Era of Large Pre-trained Language Models. arXiv preprint arXiv:2210.14179 (2022).
- [92] Chunqiu Steven Xia and Lingming Zhang. 2022. Less training, more repairing please: revisiting automated program repair via zero-shot learning. In Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering. 959–971.
- [93] Qing Xie and Atif M Memon. 2007. Designing and comparing automated test oracles for GUI-based software applications. TOSEM (2007).
- [94] Frank F Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. 2022. A systematic evaluation of large language models of code. In Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming. 1–10.
- [95] Bo Yang, Zhenchang Xing, Xin Xia, Chunyang Chen, Deheng Ye, and Shanping Li. 2021. Don't Do That! Hunting Down Visual Design Smells in Complex UIs against Design Guidelines. In ICSE. IEEE, 761–772.
- [96] Shengqian Yang, Haowei Wu, Hailong Zhang, Yan Wang, Chandrasekar Swaminathan, Dacong Yan, and Atanas Rountev. 2018. Static window transition graphs for Android. Automated Software Engineering 25, 4 (2018), 833–873.
- [97] Wei Yang, Mukul R Prasad, and Tao Xie. 2013. A grey-box approach for automated GUI-model generation of mobile applications. In *International Conference* on Fundamental Approaches to Software Engineering. Springer, 250–265.
- [98] Yanming Yang, Xin Xia, David Lo, and John Grundy. 2022. A survey on deep learning for software engineering. ACM Computing Surveys (CSUR) 54, 10s (2022). 1-73.
- [99] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Yumao Lu, Zicheng Liu, and Lijuan Wang. 2022. An empirical study of gpt-3 for few-shot knowledge-based vqa. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 36, 3081-3089
- [100] Husam N Yasin, Siti Hafizah Ab Hamid, and Raja Jamilah Raja Yusof. 2021. Droidbotx: Test case generation tool for android applications using Q-learning. Symmetry (2021)
- [101] Xia Zeng, Dengfeng Li, Wujie Zheng, Fan Xia, Yuetang Deng, Wing Lam, Wei Yang, and Tao Xie. 2016. Automated test input generation for android: Are we really there yet in an industrial case?. In Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering. 987–902
- [102] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068 (2022).
- [103] Xiaoyi Zhang, Lilian de Greef, Amanda Swearngin, Samuel White, Kyle Murray, Lisa Yu, Qi Shan, Jeffrey Nichols, Jason Wu, Chris Fleizach, et al. 2021. Screen recognition: Creating accessibility metadata for mobile applications from pixels.

In Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems. 1–15.

[104] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Learning to prompt for vision-language models. *International Journal of Computer Vision* (2022), 1–12.