Action-Based Test Carving for Android Apps

Alessio Gambi*, Hemant Gouni[†], Daniel Berreiter[‡], Vsevolod Tymofyeyev[‡], Mattia Fazzini[†]

*IMC University of Applied Science Krems, Austria; alessio.gambi@fh-krems.ac.at

[†]University of Minnesota, MN, USA; gouni008@umn.edu, mfazzini@umn.edu

[‡]University of Passau, Germany; daniel.berreiter@gmail.com, timvseffolod@gmail.com

Abstract-The test suites of an Android app should take advantage of different types of tests including end-to-end tests, which validate user flows, and unit tests, which provide focused executions for debugging. App developers have two main options when creating unit tests: create unit tests that run on a device (either physical or emulated) or create unit tests that run on a development machine's Java Virtual Machine (JVM). Unit tests that run on a device are not really focused, as they use the full implementation of the Android framework. Moreover, they are fairly slow to execute, requiring the Android system as the runtime. Unit tests that run on the JVM, instead, are more focused and run more efficiently but require developers to suitably handle the coupling between the app under test and the Android framework. To help developers in creating focused unit tests that run on the JVM, we propose a novel technique called ARTISAN based on the idea of test carving. The technique (i) traces the app execution during end-to-end testing on Android devices, (ii) identifies focal methods to test, (iii) carves the necessary preconditions for testing those methods, (iv) creates suitable test doubles for the Android framework, and (v) synthesizes executable unit tests that can run on the JVM. We evaluated ARTISAN using 152 end-to-end tests from five apps and observed that ARTISAN can generate unit tests that cover a significant portion of the code exercised by the end-to-end tests (i.e., 45% of the starting statement coverage on average) and does so in a few minutes.

I. Introduction and Motivation

Testing is an essential aspect of the software development life cycle and is crucial in improving software quality. For Android applications (or apps in short), testing is critical to avoid failures that can lead to the disruption of mission-critical activities, loss of reputation, and customer loss.

A good testing strategy needs to find an appropriate balance between the fidelity of the tests, testing speed, and testing reliability [1]. To achieve this balance, app developers can create tests at different granularity levels [1]. End-to-end tests exercise large parts of an app through its Graphical User Interface (GUI) and help developers check user flows. Unit tests focus on a small portion of an app and help developers debug and test regression. In the Android realm, developers also need to decide on which platform tests shall run [1]. Instrumented tests execute on an Android device, either physical or emulated, whereas Local tests execute on a JVM.

Although it is possible to create tests by combining granularity levels and execution environments (e.g., unit tests that run in the Android device), related work on app testing [2], [3] observed that instrumented end-to-end GUI tests and local unit tests are the most frequently used.

Researchers and practitioners proposed several automatic and semi-automatic techniques to help app developers create end-to-end instrumented tests (e.g., [4]–[8]). However, only a few existing techniques for automatically creating local unit tests have been proposed (see [9], [10]) even though developers would benefit from having those tests, which run fast and simplify debugging activities [2]. Existing techniques have limited applicability as they require substantial manual effort [9] or are not based on the tester's intent [10] (i.e., they are based on an input generation strategy).

Test carving [11] has been proposed for generating unit tests from end-to-end tests by either checkpointing portions of an application state and reloading it during unit testing to set the test preconditions (state-based carving) or by extracting units of execution from end-to-end executions and replaying them as part of the generated unit tests (action-based carving). The potential benefits of test carving include improving regression testing effectiveness by enabling standard regression test selection techniques, reducing the sensitivity to interference bugs by isolating tests better, and enabling developers to perform more focused debugging activities. Consequently, various techniques for performing test carving have been proposed (e.g., [11]-[19]). However, none of the existing techniques carves unit tests across different platforms, as it is needed in Android for extracting focused and efficient unit tests that run on the JVM from end-to-end tests that run on a device. Additionally, carving unit tests in Android also necessitates handling the challenges of testing typical Android apps that operate through event-based and inversion-of-control paradigms (i.e., everything, including the instantiation of some components required as a precondition by the unit tests, is orchestrated by the Android framework), and are tightly coupled with the Android framework, which requires replacing Android components with test doubles [3], [20], [21] during testing.

To facilitate unit testing of Android apps, we propose ARTISAN, an *action-based* test carving technique that takes as input the app under test (AUT) and its end-to-end GUI tests and automatically produces a set of locally executable unit tests as output. As illustrated in Figure 2, ARTISAN (i) instruments the AUT to enable tracing of method invocations at runtime during the execution of end-to-end GUI tests, (ii) it identifies *focal* methods to test, (iii) carves the relevant method invocations for testing a unit based on the collected traces, (iv) creates suitable test doubles for the Android framework needed by the unit tests, and (v) synthesizes executable unit tests that can run on the JVM.

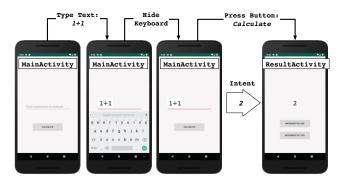


Fig. 1: Main user flow of BASICCALCULATOR.

To assess the usefulness of ARTISAN, we performed an empirical evaluation based on 152 end-to-end GUI tests from five apps and assessed the effectiveness and efficiency of our technique. As far as effectiveness is concerned, ARTISAN can generate unit tests that cover a significant portion of the code exercised by the end-to-end tests (i.e., 45% of the starting statement coverage on average). As for efficiency, ARTISAN can generate unit tests in just a few minutes. Overall, we believe that our results are promising and provide initial evidence of the usefulness of our technique.

In summary, this paper makes the following contributions:

- An automated technique that performs action-based test carving for Android apps.
- A publicly available implementation of the technique (see replication package [22]).
- An empirical evaluation that provides initial evidence of the effectiveness and efficiency of our technique.

II. BACKGROUND

In this section, we report background information that aims to facilitate the understanding of our technique. Android apps are composed of different types of components. *Activities* are one of the main types of components and are one of the primary ways to interface with app users. Only one activity can be active at any time and activities communicate via message passing. The messages used by activities are called *intents* and can contain an arbitrary, serializable payload. Under the hood, Android delivers intents by serializing and deserializing their content. Activities are also event-based: their logic is encapsulated into *listeners* and *callbacks*, and the dispatch of predefined events by the Android framework triggers their execution. Typical examples of such event listeners are lifecycle events and events generated by GUI elements.

We further illustrate background concepts using an example Android app called BASICCALCULATOR, which solves user-provided mathematical expressions through two activities: MainActivity (which solves the expressions) and ResultActivity (which displays the results of the expressions). Figure 1 displays the main user flow of BASICCALCULATOR.

MainActivity sets up the app GUI, i.e., the text input field and the "CALCULATE" button shown in Figure 1, and registers the various callbacks to handle user interactions (e.g., the pressing of the "CALCULATE" button). The MainActivity has a method called sendResult and this method is invoked when the user presses the "CALCULATE" button, which fetches the content of the text input field on the GUI, computes the result of the expression provided by the user, and sends the result to the ResultActivity (which displays the result).

In this process, sendResult retrieves the expression of the text input field using the Android findViewById method [23], which takes as input the ID of the text input field. Starting another activity requires calling the Android startActivity method [24] and passing an intent that contains a reference to the activity to start and an optional payload. In the case of our example, the activity is ResultActivity and the optional payload is the result of the expression.

Listing 1: A GUI test stressing BASICCALCULATOR.

Testing BASICCALCULATOR main user flow can be done using the GUI test in Listing 1; doing so requires developers to (i) build the app, (ii) install it inside an Android device or emulator, and (iii) execute the test interacting with the app's GUI *on* that device.

An alternative for testing BASICCALCULATOR is to use unit tests focusing on specific code units that can run on any standard JVM. Listing 2 reports a local unit test that checks the behavior of the sendResult method under the same conditions observed while running the GUI test in Listing 1.

The test_MainActivity_sendResult unit test implements a complex setup to ensure that objects such as the "CALCULATE" button and the text input field, which are normally provided by the Android framework, are also available during unit testing. Since the Android runtime does not manage those objects during local unit testing, it is necessary to replace them with mocks and stubs. To that end, the example uses Mockito [25] and Robolectric [26] Mockito is a framework for creating mocks and stubs and Robolectric provides basic stubbing capabilities by (partially) simulating the Android framework on the JVM.

In the example, the test creates a simple mock for the "CAL-CULATE" button (Lines 8–11), nested mocks that simulate accessing the text input field (Lines 18–25), and injects these mocks into Robolectric to enable their execution (Lines 29–30).

```
1 @RunWith (RobolectricTestRunner.class)
2 @Config(shadows = { EditTextShadowForSendResult.class })
3 public class Test028 {
    @Test(timeout = 4000)
    public void test_MainActivity_sendResult() {
      // Mock the Calculate button
      Button button2 = Mockito.mock(Button.class);
      Stubber stubber3 = Mockito.doReturn(R.id.
       calculateButton);
10
      Button button3 = stubber3.when(button2);
      button3.getId();
      // Instantiate MainActivity using Robolectric
      ActivityController controller = Robolectric.
13
       buildActivity (MainActivity.class);
14
      MainActivity mainactivity = controller.get();
      // Simulate triggering "create" event
      controller.create();
16
      // Mock the text input field
18
      SpannableStringBuilder stringbuilder3 = Mockito.mock(
       SpannableStringBuilder.class);
      Stubber stubber4 = Mockito.doReturn("1+1");
19
      SpannableStringBuilder stringbuilder4 = stubber4.when
20
       (spannablestringbuilder3);
21
      spannablestringbuilder4.toString();
      EditText edittext3 = Mockito.mock(EditText.class);
      Stubber stubber5 = Mockito.doReturn(
       spannablestringbuilder3);
24
      EditText edittext4 = stubber5.when(edittext3);
      edittext4.getText();
       // Inject the mocks in Robolectric
      EditText edittext2 = mainactivity.findViewById(R.id.
       input);
      EditTextShadowForSendResult edittextshadow = Shadow.
28
       extract (edittext2);
      edittextshadow.setMockFor("android.widget.EditText:
29
       android.text.Editable getText()", edittext3);
      edittextshadow.setStrictShadow();
       // Invoke the Method Under Test
      mainactivity.sendResult(button3);
33
34 }
```

Listing 2: A unit test carved by ARTISAN for BASICCALCULATOR

In summary, this motivating example shows how complex Android apps' unit tests can be and illustrates some technical challenges in their automated generation.

III. TECHNIQUE

ARTISAN is an end-to-end approach composed of several steps (Figure 2). It starts by instrumenting the (non-obfuscated) original AUT to enable tracing method invocations. Next, it executes the Instrumented AUT on an Android device against GUI Tests to collect Execution Traces. Then, it parses the traces into a form amenable to automatic analysis (i.e., graphs) and carves the original executions. Finally, it augments the Carved Executions with code that mocks dependencies provided by the Android framework and synthesizes the source code of the Carved Unit Tests from the Extended Carved Executions.

Notably, ARTISAN generates carved unit tests in the static single-assignment (SSA) form [27]. Although programs written in SSA form are generally longer than programs written in other forms, which might affect their readability, we decided

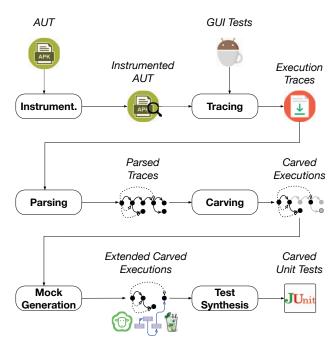


Fig. 2: An overview of ARTISAN's end-to-end approach to carving local unit tests from instrumented GUI tests.

to generate test code in SSA form for two main reasons: on the one hand, it is easy to translate the carved executions in this form; on the other hand, SSA enables the application of standard program analyses (e.g., live variable analysis) and optimization techniques (e.g., test minimization [28]) that can improve the generated tests' quality.

A. Instrumentation and Tracing

Action-based test carving is a dynamic analysis technique that requires execution traces to identify the units of execution, i.e., method invocations, that are relevant for testing code units. ARTISAN instruments the original AUT using the byte-code modification library Soot [29].

ARTISAN adopts a light-weighted approach to trace Android apps and generates traces in plain text, making it possible for developers to easily inspect them. Specifically, ARTISAN injects code that logs for each method invocation, the method signature, the actual parameters, and any returned values or thrown exceptions. Tracing focuses only on the AUT's operations, as we want unit tests that focus on the AUT. To this end, ARTISAN instruments only the method bodies of the methods that belong to the AUT and not other components such as third-party libraries, standard language libraries, and the Android runtime. In the instrumentation, ARTISAN differentiates between calls to instance and static methods, traces method calls at different visibility levels and distinguishes whether methods return normally or exceptionally. In the latter case, the trace contains also whether the exceptional behavior was caused by a checked or an unchecked exception.

Differently than state-based carving, ARTISAN does not check-point the application state nor serialize the objects used as parameters or return values. Instead, it manages all the object instances, including exceptions, by-reference and stores in the trace only their object id. ARTISAN obtains the object id of non-null instances by calling System.getObjectIdentity() to capture the actual object types along with their hash codes. To avoid cluttering the execution traces, ARTISAN manages primitive types, boxed primitive types, and "string" types (e.g., String) by-value and saves only their string representation.

Action-based carving implicitly assumes that all interactions between objects happen exclusively via method invocations; unfortunately, units of executions like array stores, array accesses, and field assignments are not implemented as method calls in Java and would be missed if not properly handled. To avoid missing such fundamental units of execution, ARTI-SAN implements a custom instrumentation code that traces them as *synthetic* methods. For example, it traces an array store like array[0] = 10 as the generic method invocation abc.ArrayOperation.set (array, 0, 10).

As discussed in Section II, activities communicate by passing Intent objects that the Android framework serializes in a stream of bytes and deserializes into actual objects. Thus, the object ids of serialized and deserialized objects differ, which effectively breaks the (logical) connection between them. To avoid losing this connection, ARTISAN leverages application level taint tracking [30], [31]. Instead of modifying the Intent's bytecode to accommodate the tainting value, i.e., the object id of the intents to be sent, ARTISAN stores the tainting values directly in the intents as a regular payload using a special key before Android sends them. To read the tainting value, instead, ARTISAN injects custom code that is invoked before the receiving activity accesses the payload. This custom code extracts from the payload the tainted value by invoking a standard Intent method using the special key as a parameter. Doing so, ARTISAN is able to trace the logical connection between sent and received intents.

B. Trace Parsing

After tracing the execution of the GUI tests, ARTISAN parses the generated execution traces into graph data structures that capture chronological, data, and call dependencies of method invocations and object instances. For each trace, the technique creates three graph data structures.

Execution Flow Graph (EFG). This graph captures the chronological dependencies of the method invocations appearing in the trace; hence, EFGs are useful to identify (past) method invocations that may set test preconditions. The graph is a (doubly) linked list whose nodes represent method invocations and edges the (strict) precedence/follow relations. Data Dependency Graph (DDG). This graph is a directed graph that links method invocations to data nodes and data nodes to method invocations. Data nodes can be either object instances or primitive values. Object instances can be linked to multiple method invocations, whereas primitive values are

always linked to one and only one method invocation. In this graph, two nodes are linked when (i) an object instance OWNS a method invocation; (ii) a data node is used as a PARAMETER of a method invocation; (iii) a STATIC data node is used inside a method invocation; (iv-a) a method invocation RETURNS a data node or (iv-b) THROWS an exception. DDG is useful to identify test preconditions.

Call Dependency Graph (CDG). This graph captures the nesting relations among the method invocations; hence, it is useful to ensure that carved method invocations are executed the right amount of times (e.g., no duplicate executions). The graph is a directed and acyclic graph whose nodes represent method invocations and edges the INVOKE relation. The graph is a forest because apps can have multiple entry points.

After parsing is completed, ARTISAN decorates the graphs by including additional information that will be used later during carving and test synthesis. This step includes (i) identifying method invocation nodes that are owned by Android components, like Activities, and tagging the nodes corresponding to life cycle events callbacks (e.g., onCreate); (ii) aliasing data nodes that correspond to sources and sinks of tainted intents (i.e., logically related intents) and their payload; and, (iii) injecting static dependencies.

C. Action-based Carving

The carving process begins after parsing the execution traces and considers only one trace at a time. At first, ARTISAN selects "carvable" targets, i.e., the method invocations for which to generate unit tests. In general, there are no restrictions on which method invocations can be carved; however, carving some method invocations, such as private methods and methods that do not belong to the AUT might produce non-compilable or irrelevant unit tests. Therefore, ARTISAN automatically filters out invocations of private methods and invocations whose owner type does not match the AUT's package name (i.e., they do not belong to the AUT). In case an end-to-end test makes multiple calls to the same target method ARTISAN provides a special option to either select one (i.e., the first) or all the invocations of that method as "carvable" targets (i.e., ARTISAN offers different carving strategies to developers). By default, ARTISAN selects the first invocation of a method to be in the set of "carvable" targets. The idea behind using this default strategy is that it allows for covering the behavior of multiple methods while limiting the cost of running the technique. This strategy might carve the same method multiple times across different traces. Still, we think that this is reasonable as traces originating from different GUI tests likely have different objectives.

After selecting the target method invocations, ARTISAN finds all the method invocations that are relevant to them, either directly or indirectly, using a backward slicing algorithm: Starting from a target method invocation, this algorithm identifies the past method invocations that match one of the following three conditions: (1) the method invocation shares the same owner with the target method invocation; (2) the method invocation is owned by a parameter used by the target

method invocation; or, (3) the method invocation belongs to a static class used within the target method body.

Since the selected method invocations might introduce additional dependencies (e.g., parameters to be set), this algorithm iteratively carves each of them. The algorithm converges because, at every iteration, it considers a smaller set of dependencies. Moreover, as the "carvable" targets are considered sequentially and might share dependencies, ARTISAN caches intermediate results and speeds up the algorithm.

Once the algorithm selects all the method invocations relevant to a "carvable" target, ARTISAN uses the CDG to retrieve all the additional method invocations that would be executed because the relevant method invocations are (e.g., a method invocation triggered by an already selected method invocation). Given this "extended" set of invocations, ARTISAN creates carved executions by extrapolating the connected components they form in the EFG, the DDG, and the CDG.

Finally, ARTISAN "cleans" the carved executions by removing those units of computation, such as lambdas, which occur in Android apps but cannot be directly instantiated in the unit tests, and re-carves the target invocations within the carved trace executions to ensure they remain consistent.

D. Mock Generation

Carved executions contain a list of executed method invocations and their data dependencies. However, they might not be executable yet as they might contain method invocations on objects belonging to the Android framework. To be able to execute those invocations, ARTISAN uses the Roboletric [26] framework, which offers simplified implementations for the classes in the Android framework (which the framework calls shadow classes). These simplified implementations are designed for testing purposes. However, Roboletric is not comprehensive [26], and some of the method invocations in the carved executions might not be executable as the framework does not offer an implementation for them. Additionally, there might be objects that, during end-to-end test execution, are created within the Android framework (e.g., GUI elements) or by third-party libraries. Consequentially, ARTISAN does not know how to either define them or should not use their implementation (as ARTISAN wants to minimize the use of external dependencies such as third-party libraries). ARTI-SAN deals with these cases by means of mocking, a standard technique that improves the reliability of unit tests by replacing complex dependencies with pre-programmed test doubles.

ARTISAN analyzes how each of those objects without an "accessible" implementation is used within the carved executions and automatically programs a mock object that can replicate the (observed) behavior within the unit tests through stubbed methods (i.e., methods that return canned data). If a stubbed method invocation on a mocked object returns another object without an accessible implementation, ARTISAN reproduces its behavior utilizing another automatically configured mock object; the process continues until there are no more objects without an accessible implementation left. ARTISAN takes advantage of forward slicing to identify

which mock objects are needed. Since the number of data dependencies to consider is finite and shrinks at every iteration of the algorithm, synthesizing mocks is guaranteed to always terminate. After the mock generation phase, we call carved executions as *extended carved executions*.

E. Synthesis of Carved Unit Tests

At this point, all the test preconditions are either instantiated or mocked, and ARTISAN can finally synthesize the code implementing the unit tests by transforming the extended carved execution's CDG root-level nodes, i.e., the directly "visible" method invocations, into their corresponding source-code method invocations. While doing so, ARTISAN relies on the DDG to generate all the variables needed to host the references or values that correspond to methods' owners, parameters, and return values. Notably, the extended carved executions do not contain complex control flows; thus, the generated tests consist of a number of variable declarations and a sequence of method invocations, as one would expect from unit tests.

To generate the mocking code which stubs the methods of objects without an accessible implementation and contained in the extended carved executions, ARTISAN uses a predefined template. For each object \circ , ARTISAN (i) declares \circ as a mock object using Mockito [25], (ii) it specifies which object (if any) the mock should return using a stub, and (iii) it invokes the stubbed method sm on \circ based on the carved execution.

The GUI elements of an app require special treatment. To inject the pre-programmed mock objects inside GUI elements (which are handled by Robolectric), ARTISAN (i) retrieves the GUI elements from the activity using their known unique ID, (ii) it extracts the shadow objects simulating them, and (iii) passes the mock objects to the shadows.

IV. EMPIRICAL EVALUATION

We performed a preliminary evaluation of ARTISAN's effectiveness and efficiency by targeting the following research questions (RQs):

RQ1: Can ARTISAN carve unit tests from GUI tests?

RQ2: What is the cost of running ARTISAN?

RQ3: What are the characteristics of carved tests?

A. Experimental Benchmarks

In the evaluation, we used five open-source Android apps (Table I) with existing tests. We used open-source apps because through their public repositories they make available both their source code and the existing GUI tests, required by the evaluation of our technique. To identify relevant apps, we used a dataset of 1,002 apps with tests from related work [2]. The apps in the dataset are publicly available on GitHub and, to the best of our knowledge, the dataset was the largest set of apps with tests at the time we started evaluating ARTISAN.

We selected the five apps from the dataset as follows: First, we identified apps that contain GUI tests written in Espresso [32]; this step identified 245 relevant apps. Second, we filtered out apps that use programming languages other

TABLE I: Benchmark apps used in the empirical evaluation.

ID	Name	Category	Version	LOC (K)	GUI Tests
A1	BLABBERTABBER	Tools	1.0.10	2.4	9
A2	FIFTHELEMENT	Music	2.2.5	69.8	17
A3	OWL FLASH CARDS	Education	1.1	6.4	12
A4	PRISMACALLBLOCKER	Tools	1.2.3	12.0	67
A5	UK-GM	Education	1.2.1	5.3	47

than Java (e.g., Kotlin), as ARTISAN does not currently support them; this step left us with 180 apps. Third, we sorted the 180 apps in descending order based on the number of GUI tests associated with the apps and processed the list of apps starting from the one having the highest number of tests. We discarded any app that we could not build and for which the available GUI tests did not pass or showed flakiness. We determined whether all tests were passing and were not flaky by checking whether all tests passed in ten runs of the tests. We stopped as soon as we identified five apps. We could not build or run some of the top-ranked apps as those had outdated dependencies, required an API key for a third-party service of the app, or interacted with servers not reachable anymore.

Table I summarizes the main elements of the apps we considered in the evaluation. For each app, the table reports an identifier for the app (ID), its name (Name), the category of the app (Category), its version (Version), the number of source and test code lines (in thousands) (LOC(K)), and the number of existing GUI tests (GUI).

B. Experimental Settings

To answer the RQs, we ran ARTISAN on the five apps considered on a dedicated workstation with 128GB of memory, an Intel i9-9900K 3.60GHz processor, and running Ubuntu 18.04. To execute the GUI tests, we used an Android Nexus 5X emulator running API 28. We used API 28 as it was compatible with all the selected apps according to their supported Android API versions. In the RQs, we evaluated ARTISAN using the default carving strategy (i.e., the strategy that selects as carving targets only the first occurrence of method invocations with the same fully qualified method signature within a trace).

C. Results

RQ1: Can ARTISAN carve tests from GUI tests?: Table II reports the results of running ARTISAN on the benchmark apps. For each app, the table reports the identifier of the app (ID), the number of traces collected from running the GUI tests associated with the app (Traces), the number of method invocations in the traces (Method Invocations), the number of traces that ARTISAN could parse while carving tests (Parsed Traces), the number of method invocations selected by the carving strategy (Targets), the number of tests carved (Carved Tests), and the statement and branch coverage achieved by GUI and carved tests (GUI Tests and Carved Tests columns under the Statement Coverage and Branch Coverage headers). Additionally, the table relates the coverage of carved tests with the one of GUI tests (columns labeled with Included under

the *Statement Coverage* and *Branch Coverage* headers) by reporting the percentage of coverage from carved tests that also appears in GUI tests.

Overall, ARTISAN carved 2,087 tests from 152 GUI tests. The carved tests cover 45.28% and 41.33% of the statements and branches that are covered by the GUI tests. The overall number of targets is 3,609, and the number of method invocations in the traces is 311,661. The difference between the number of method invocations in the traces and the number of targets is due to the fact that the traces contain a large number of method invocations whose definition is not inside the AUT (i.e., the methods are defined in the Java standard library, third-party libraries, or the Android framework) and due to the default carving strategy we adopted.

The difference between the number of targets and the number of carved tests is caused by some limitations in the implementation of the technique (see Section IV-F) and our design choice to reject unit tests in which the method under tests are not *directly* called. Additionally, some targets cannot be carved by ARTISAN as those are methods in anonymous classes (e.g., callback handler definitions for GUI elements) that cannot be directly invoked in Java. To carve those targets, ARTISAN could use a preprocessing step that refactors the code of the app such that those methods are not part of anonymous classes. However, we did not implement such a solution as it could lead to unwanted changes by the developers. There is a need for studies and interviews with developers to investigate this aspect, and we leave those studies as a possible direction for future work.

The design choices, limitations in the implementation of the technique (Section IV-F), the focus on ARTISAN on executions originating in the main thread [33], and the focus on Android activities are the reasons why the coverage of carved tests is not 100%. Nevertheless, ARTISAN still provides unit tests that offer a considerable coverage of the GUI tests.

Table II also compares (columns *Included* under the *State*ment Coverage and Branch Coverage headers) the statement and branch coverage achieved by GUI and carved tests. Specifically, we look at how much of the coverage in carved tests also appears in the GUI tests used for generating them. In other words, these columns reveal whether the carved tests cover portions of the apps that are not covered by the GUI tests (i.e., lead to spurious coverage). We computed this information by extending JaCoCo [34], the coverage tool we used in the experiments. All the branches covered in the carved tests are also covered by the GUI tests. In terms of statement coverage, instead, there are a few statements (in two of the five apps) that are covered by the carved tests but not by the GUI tests. We analyzed the tests leading to the discrepancies and identified that the cause behind the discrepancy is a different behavior of some of the Android API methods when they execute on an Android device and the JVM (via Robolectric). This situation can appear because Robolectric is a partial model of the Android framework.

TABLE II: Results of running ARTISAN on the benchmark apps.

ID	Traces Method	Mathod Imposations	Darsad Traces	s Targets	Carved Tests	Statement Coverage (%)		Branch Coverage (%)			
	Traces	Traces Memoa Invocations	Tursea Traces			GUI Tests	Carved Tests	Included	GUI Tests	Carved Tests	Included
A1	9	1,654	9	55	48	45	12	100	35	11	100
A2	17	35,495	17	1,559	1,004	49	18	100	35	15	100
A3	12	13,306	12	221	126	64	30	98	36	14	100
A4	67	31,151	67	1,479	821	64	29	99	52	23	100
A5	47	230,055	3	295	88	90	50	100	65	16	100

RQ1 answer: Yes, ARTISAN can carve tests from GUI tests. Additionally, carved tests cover 45.28% and 41.33% of the statements and branches that are covered by the GUI tests and rarely have spurious coverage.

RQ2: What is the cost of running ARTISAN?: Table III provides details on the execution time for running ARTISAN on the benchmark apps. For each app, the table provides the identifier for the app (ID), the time to execute the GUI tests before instrumenting the app (GUI Tests Execution Time Before Instrumentation), the time needed to instrument the app (Instrumentation Time), the time to execute the GUI tests after instrumenting the app (GUI Tests Execution Time After Instrumentation), and the time to carve the tests (Carving Time). The time values reported in Table III are averages across 10 runs of ARTISAN.

ARTISAN was able to generate carved tests in less than one hour for each app considered. The technique was the fastest when analyzing A1 and, in this case, ARTISAN took only 53 seconds. The technique took the longest when analyzing A2, roughly 42 minutes. The time to instrument the apps is negligible for the selected apps, especially considering that instrumentation is a one-time activity. The overhead introduced by the instrumentation when running the GUI tests is 16.34%. We believe that the overhead is reasonable as it is low and did not affect the test execution behavior (i.e., the tests passed before and after the instrumentation). The time to carve tests, instead, varies significantly between apps. We analyzed the causes behind the variance and identified that the main contributing factors are the presence of static method invocations in the traces and the size of the carved executions. Static method invocations affect the carving time as all the invocations to the same static classes need to be (conservatively) taken into account when analyzing the targets those invocations precede. Larger sets of required method invocations also affect carving time as it takes longer to extrapolate connected components from the graphs to generate the corresponding carved executions.

RQ2 answer: Based on the results of our evaluation, we believe that the time cost of running ARTISAN is low. For the apps considered, the technique always terminated within an hour, and, in some cases, within a few minutes.

RQ3: What are the characteristics of carved tests?: To characterize carved tests, we consider the size of the unit tests

TABLE III: Time cost of running ARTISAN.

ID	GUI Tests Execution Time Before Instrumentation	ARTISAN				
		Instrumentation Time	GUI Tests Execution Time After Instrumentation	Carving Time		
A1	34s	10s	36s	7s		
A2	2m56s	28s	3m26s	38m07s		
A3	1m49s	18s	1m41s	31s		
A4	6m10s	15s	6m24s	3m12s		
A5	4m51s	19s	6m53s	6m16s		

and the number of mocks contained in them. We focus on test size as larger tests can be harder to maintain and can lead to test smells [35]. We also focus on the number of mocks as they are not always straightforward to set up [3]; hence, having tests with them can potentially help app developers.

For each app considered, Figure 3 reports the size of the carved tests. The chart reports the size of the tests on the y-axis and uses the log scale. We computed the size of each test by counting the number of statements in the tests. The results reported in Figure 3 are promising. For four out of the five apps considered, the median number of statements in the tests is less than 10. A2 is the app with the largest number of statements per test. We observed that this app requires setting some values in its database for a larger number of tests which, in turn, led to an increase in the size of the tests. Considering that most tests have a reasonable size, we believe that ARTISAN can provide developers with tests that might be useful for debugging.

The total number of mocks in the carved tests is 241. The ratio between the number of tests and the number of mocks follows the ratio of developer-written tests in some of the apps analyzed by related work on test doubles in Android [3].

Based on our experimental results, we argue that the tests generated by ARTISAN, both in terms of test size and the mocks they provided, could be actionable for developers. However, to confirm this hypothesis, studies and interviews with developers are necessary; hence, we suggest and envision them as possible future work.

RQ3 answer: The tests carved by ARTISAN tend to be concise and provide mocks that are required for testing certain parts of the apps. Considering the tests' characteristics, we believe that the tests could be actionable for developers.

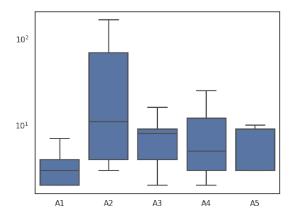


Fig. 3: Test size (y-axis) for carved tests generated by ARTI-SAN. The graph uses a logarithmic scale.

D. Discussion

ARTISAN is an end-to-end technique to perform tests carving for Android apps. Implementing such a technique required a significant technical effort, the mastery of several technologies, and knowledge from different domains besides Android. For instance, we employed byte code modification for instrumenting and tracing the execution, automated build systems for setting up the evaluation benchmarks, graph theory and program analysis to carve the execution traces, and Java code generation to synthesize the unit tests.

There are parts of ARTISAN that provide a solid base for future work and others that can be further refined. For instance, instrumentation and tracing of apps, as well as experiment automation, worked smoothly. Similarly, the carving algorithm and the generation of mocks produced excellent results. However, the carving algorithm also has some cost when carving apps that heavily employ loops or implement overly complex methods. In these cases, carving might take too long and produce excessively long unit tests that can quickly become hard to inspect and manage. We postulate that in these cases, it might be beneficial to combine action-based and state-based carving so that unit tests can directly load large objects created with long method sequences, or loops, from memory. Another way to reduce the size of the carved unit tests might be using other dynamic techniques, such as delta debugging [36], or employing purity analysis to identify and filter out method invocations that do not introduce any relevant dependency. Additionally, different carving strategies could be explored to make action-based carving more efficient and effective.

E. Threats to Validity

As it is the case for most empirical evaluations, there are both external and construct threats to validity associated with the results we presented. In terms of external validity, our results might not generalize to other apps. In particular, we only considered five apps. This limitation is an artifact of the complexity involved in setting up the infrastructure to run the apps, which might require customized build configurations and manually inspecting the results of our analysis. We selected apps of different sizes that belong to different app categories and are already considered in related work to mitigate this threat. In terms of construct validity, there might be errors in the implementation of our technique. To mitigate this threat, we extensively inspected the evaluation results manually.

F. Limitations

In its current implementation, ARTISAN only supports Android apps written in Java. Additionally, ARTISAN synthesizes executable unit tests that lack test oracles, do not involve any Android components besides activities (e.g., does not handle Android services) and GUI elements, and consider only traces generated by the main Android thread. We argue that automatically generating test oracles, which is currently an open research problem, is outside the scope of this first work on carving unit tests for Android apps. Additionally, we believe that extending ARTISAN's implementation to handle apps written in Kotlin, more Android components, and multiple threads is mostly an engineering effort.

Regarding our evaluation, we did not involve developers to evaluate the quality of carved tests. We plan to perform such an evaluation in future work. Specifically, we plan to perform studies with developers to understand which carved tests best help developers and explore alternative carving strategies based on the results of the studies.

V. RELATED WORK

The idea of test carving was originally proposed by Elbaum et al. [11] as a means to generate unit tests, dubbed Differential Unit Tests, to spot regression errors in Java programs. Elbaum and co-authors identified three main test carving paradigms: state-based carving, in which carving takes place on the system under test's state recorded during execution; actionbased carving, in which carving takes place on the sequence of method invocations recorded during execution; and hybrid carving, which combines the previous two. However, they implemented only state-based carving. In this area, remarkable results have been achieved by Krikava and Vitek [12], who proposed GENTHAT for generating unit tests from execution traces of R libraries, Kampmann and Zeller [13], who proposed BASILISK for state-based carving of parameterized unit tests targeting C programs, and, Juvekar et al. [15], who created a program executing all public methods on a given object the same way as in a given program trace. Compared to those works, ARTISAN implements a different form of test carving, works on Android apps, generates unit tests across different platforms, and augments carved tests with automatically synthesized mocks.

ARTISAN generates focused unit tests from execution traces. Pasternak et al. [19], Saff et al. [17], and Thummalapenta et al. [18] achieved the same goal but with different techniques that, respectively, selected the interactions

to recreate the state of objects until a certain point in time, automatically created focused unit tests by test factoring and environmental mocks generation, and mined an extensive collection of execution traces to generate generic parameterized tests using dynamic symbolic execution to cover paths not contained in the traces. Unlike these techniques, ARTISAN is not limited to inter-object interactions implemented as method calls, considers method invocations invoked by frameworks, and generates focused unit tests.

Action-based carving sets unit tests' preconditions, i.e., objects' state, by re-executing specific sequences of method calls that have been observed during end-to-end testing. Therefore, action-based carving is a sensible solution to the object creation problem defined by Bach et al. [37]. The main differences between ARTISAN and the work done by Bach et al. lies in the fact that they first identified feasible method-call sequences for object creation statically and then selected the most desirable sequence using a search algorithm. Another difference is that Bach et al.'s approach is for C++ programs.

Alternative approaches make use of selective capture and replay techniques. For instance, Orso et al. proposed SCARPE [16] and JINSI [14] to capture parts of program execution for replay and isolate the instructions that lead to the failure to produce a minimal example that reliably replays it. Compared to those works, ARTISAN has a broader scope as it does not consider interactions involving a single component, can generate tests from both normal and exceptional executions, and can carve unit tests for Android apps.

VI. CONCLUSIONS AND FUTURE WORK

This paper presented ARTISAN, a technique to perform test carving for Android apps. ARTISAN carves unit tests that run on the JVM from GUI tests that run on a device. We evaluated ARTISAN based on 152 GUI tests and five apps and identified that the technique carves tests that achieve 45% of the original GUI tests' coverage and does so in an amount of time compatible with standard development practices.

In future work, we plan to perform studies and interviews with developers to understand which carved tests best help developers and explore alternative carving strategies based on the gathered insights. We also plan to investigate test suite reduction techniques to identify duplicate tests among carved tests. Finally, we plan to investigate techniques to carve oracles from end-to-end tests into oracles suitable for unit tests.

REFERENCES

- [1] Google, "Fundamentals of testing android apps." [Online]. Available: https://developer.android.com/training/testing/fundamentals
- [2] J.-W. Lin, N. Salehnamadi, and S. Malek, "Test automation in open-source android apps: A large-scale empirical study," in 2020 35th IEEE/ACM International Conference on Automated Software Engineering. New York, NY, USA: Association for Computing Machinery, 2020, pp. 1078–1089.
- [3] M. Fazzini, C. Choi, J. M. Copia, G. Lee, Y. Kakehi, A. Gorla, and A. Orso, "Use of test doubles in android testing: An in-depth investigation," in *Proceedings of the 44th International Conference on Software Engineering*, 2022.

- [4] K. Mao, M. Harman, and Y. Jia, "Sapienz: Multi-objective automated testing for android applications," in *Proceedings of the 25th International Symposium on Software Testing and Analysis*. New York, NY, USA: Association for Computing Machinery, 2016, pp. 94–105.
- [5] T. Su, G. Meng, Y. Chen, K. Wu, W. Yang, Y. Yao, G. Pu, Y. Liu, and Z. Su, "Guided, stochastic model-based gui testing of android apps," in *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*. New York, NY, USA: Association for Computing Machinery, 2017, pp. 245–256.
- [6] Y. Li, Z. Yang, Y. Guo, and X. Chen, "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). Piscataway, NJ, USA: Institute of Electrical and Electronics Engineers, 2019, pp. 1070–1073.
- [7] J. Wang, Y. Jiang, C. Xu, C. Cao, X. Ma, and J. Lu, "Combodroid: generating high-quality test inputs for android apps via use case combinations," in *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. New York, NY, USA: Association for Computing Machinery, 2020, pp. 469–480.
- [8] M. Fazzini, E. N. D. A. Freitas, S. R. Choudhary, and A. Orso, "Barista: A technique for recording, encoding, and running platform independent android tests," in 2017 IEEE International Conference on Software Testing, Verification and Validation. Piscataway, NJ, USA: Institute of Electrical and Electronics Engineers, 2017, pp. 149–160.
- [9] Y. Liu, Y. Lu, and Y. Li, "An android-based approach for automatic unit test," in *International Conference on Cyberspace Technology (CCT* 2014), 2014, pp. 1–4.
- [10] J. Cao, H. Huang, and F. Liu, "Android unit test case generation based on the strategy of multi-dimensional coverage," in 7th International Conference on Cloud Computing and Intelligent Systems, 2021.
- [11] S. Elbaum, H. N. Chin, M. B. Dwyer, and J. Dokulil, "Carving differential unit test cases from system test cases," in *Proceedings* of the 14th ACM SIGSOFT International Symposium on Foundations of Software Engineering. New York, NY, USA: Association for Computing Machinery, 2006, p. 253–264.
- [12] F. Krikava and J. Vitek, "Tests from traces: automated unit test extraction for R," in *Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis*. New York, NY, USA: Association for Computing Machinery, 2018, pp. 232–241.
- [13] A. Kampmann and A. Zeller, "Carving parameterized unit tests," in Proceedings of the 41st International Conference on Software Engineering: Companion Proceedings. Piscataway, NJ, USA / New York, NY, USA: Institute of Electrical and Electronics Engineers / Association for Computing Machinery, 2019, pp. 248–249.
- [14] A. Orso, S. Joshi, M. Burger, and A. Zeller, "Isolating relevant component interactions with jinsi," in *Proceedings of the 2006 international workshop on Dynamic systems analysis*. New York, NY, USA: Association for Computing Machinery, 2006, pp. 3–10.
- [15] S. Juvekar, J. Burnim, and K. Sen, "Path slicing per object for better testing, debugging, and usage discovery," EECS Department, University of California, Berkeley, Tech. Rep., Sep 2009. [Online]. Available: http://www2.eecs.berkeley.edu/Pubs/TechRpts/2009/EECS-2009-132.html
- [16] A. Orso and B. Kennedy, "Selective capture and replay of program executions," ACM SIGSOFT Software Engineering Notes, vol. 30, no. 4, pp. 1–7, 2005.
- [17] D. Saff, S. Artzi, J. H. Perkins, and M. D. Ernst, "Automatic test factoring for java," in *Proceedings of the 20th IEEE/ACM international Conference on Automated software engineering*. New York, NY, USA: Association for Computing Machinery, 2005, pp. 114–123.
- [18] S. Thummalapenta, J. de Halleux, N. Tillmann, and S. Wadsworth, "Dygen: Automatic generation of high-coverage tests via mining gigabytes of dynamic traces," in *Tests and Proofs*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2010, pp. 77–93.
- [19] B. Pasternak, S. Tyszberowicz, and A. Yehudai, "Genutest: a unit test and mock aspect generation tool," *International Journal on Software Tools for Technology Transfer*, vol. 11, no. 4, pp. 273–290, 2009.
- [20] G. Meszaros, xUnit test patterns: Refactoring test code. Pearson Education, 2007.
- [21] M. Fowler, "Testdouble." [Online]. Available: https://martinfowler.com/bliki/TestDouble.html
- [22] A. Authors, "Artifact for action-based test carving for android apps." [Online]. Available: https://zenodo.org/record/7285409

- [23] Google, "Android view." [Online]. Availhttps://developer.android.com/reference/android/view/ able: View#findViewById(int)
- [24] activity." Avail-"Android [Online]. https://developer.android.com/reference/android/app/ able: Activity#startActivity(android.content.Intent)
- [25] S. Faber, B. Dutheil, R. Winterhalter, and T. van der Lippe, "Mockito." [Online]. Available: https://site.mockito.org/
- [26] Robolectric, "Robolectric." [Online]. Available: http://robolectric.org [27] R. Cytron, J. Ferrante, B. K. Rosen, M. N. Wegman, and F. K. Zadeck, "Efficiently computing static single assignment form and the control dependence graph," ACM Trans. Program. Lang. Syst., vol. 13, no. 4, pp. 451-490, 1991. [Online]. Available: https://doi.org/10.1145/115372.115320
- [28] A. Leitner, M. Oriol, A. Zeller, I. Ciupa, and B. Meyer, "Efficient unit test case minimization," in 22nd IEEE/ACM International Conference on Automated Software Engineering (ASE 2007), November 5-9, 2007, Atlanta, Georgia, USA, R. E. K. Stirewalt, A. Egyed, and B. Fischer, Eds. ACM, 2007, pp. 417–420. [Online]. Available: https://doi.org/10.1145/1321631.1321698
- [29] R. Vallée-Rai, P. Co, E. Gagnon, L. Hendren, P. Lam, and V. Sundaresan, "Soot: A java bytecode optimization framework," in CASCON First Decade High Impact Papers. Armonk, NY, USA: IBM Corp., 2010, p. 214-224.

- [30] J. Clause, W. Li, and A. Orso, "Dytan: a generic dynamic taint analysis framework," in Proceedings of the 2007 international symposium on Software testing and analysis, 2007, pp. 196-206.
- [31] J. Newsome and D. X. Song, "Dynamic taint analysis for automatic detection, analysis, and signaturegeneration of exploits on commodity software." in NDSS, vol. 5, 2005, pp. 3-4.
- Google, "Espresso." [Online]. https://developer.android.com/training/testing/espresso
- "Processes and threads overview." [Online]. Available: https://developer.android.com/guide/components/processes-and-threads
- [34] JaCoCo, "Jacoco." [Online]. Available: https://www.jacoco.org
- [35] G. Grano, F. Palomba, D. Di Nucci, A. De Lucia, and H. C. Gall, "Scented since the beginning: On the diffuseness of test smells in automatically generated test code," Journal of Systems and Software, vol. 156, pp. 312-327, 2019.
- [36] A. Zeller, "Isolating cause-effect chains from computer programs," in Proceedings of the 10th ACM SIGSOFT Symposium on Foundations of Software Engineering. New York, NY, USA: Association for Computing Machinery, 2002, pp. 1–10.
- [37] T. Bach, R. Pannemans, and A. Andrzejak, "Determining methodcall sequences for object creation in C++," in Proceedings of the 13th IEEE International Conference on Software Testing, Validation and Verification. Piscataway, NJ, USA: Institute of Electrical and Electronics Engineers, 2020, pp. 108-119.