



Model inference of Mobile Applications with dynamic state abstraction

Sébastien Salva, Patrice Laurencot, Stassia Zafimiharisoa

► To cite this version:

Sébastien Salva, Patrice Laurencot, Stassia Zafimiharisoa. Model inference of Mobile Applications with dynamic state abstraction. Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing 2015, 612, 2016, Studies in Computational Intelligence (SCI). hal-02019290

HAL Id: hal-02019290

<https://hal.uca.fr/hal-02019290>

Submitted on 14 Feb 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Chapter 1

Model inference of Mobile Applications with dynamic state abstraction

Sébastien Salva and Patrice Laurençot and Stassia R. Zafimiharisoa

Abstract We propose an automatic testing method of mobile applications, which also learns formal models expressing navigational paths and application states. We focus on the quality of the models to later perform analysis (verification or test case generation). In this context, our algorithm infers formal and exact models that capture the events applied while testing, the content of the observed screens and the application environment changes. A key feature of the algorithm is that it avoids the state space explosion problem by dynamically constructing state equivalence classes to slice the state space domain of an application in a finite manner and to explore these equivalence classes. We implemented this algorithm on the tool *MCrawlT* that was used for experimentations. The results show that *MCrawlT* achieves significantly better code coverage than several available tools in a given time budget.

Key words: Model inference, Automatic testing, Android applications, State abstraction.

1.1 Introduction

Desktop, Web and more recently mobile applications are becoming increasingly prevalent nowadays and a plethora are now developed for several heterogeneous platforms. All these pieces of software need to be tested to assess the quality of

Sébastien Salva
LIMOS CNRS UMR 6158, University of Auvergne, France, e-mail: sebastien.salva@udamail.fr

Patrice Laurençot
LIMOS CNRS UMR 6158, Blaise Pascal University, France e-mail: laurencot@isima.fr

Stassia R. Zafimiharisoa
LIMOS CNRS UMR 6158, Blaise Pascal University, France e-mail: s.zafimiharisoa@openium.fr

their features in terms of functionalities e.g., conformance, security, performance, etc. Manual testing is the most employed approach for testing them, but manual testing is often error-prone and insufficient to achieve high code coverage. These applications share a common feature that can be used for automatic testing: they expose GUIs (Graphical User Interface) for user interaction which can be automatically experimented and explored. Several works already deal with GUI applications testing e.g., desktop applications [9], Web applications [3] or mobile ones [2]. These approaches interact with applications in an attempt to detect bugs and eventually to record models, but all with the same purpose: to obtain good code coverage quickly.

The work, proposed in this paper, falls under this automatic testing category and tackles the testing of mobile applications but also, and above all the learning of models. Our study of model inference techniques has revealed that they often leave aside the notion of correctness of the learned models. This feature is not required for just detecting bug, but is mandatory if models are later used for analysis. Indeed, false models may easily lead to false positives. The quality of the model with regards to its level of abstraction and the amount of information it captures is important as well. Indeed, the more data we collect, the more precise an analysis can be done thereafter. Nevertheless, large amounts of data often lead to large models, up to a state space explosion problem. Based on these observations, we propose an algorithm that aims at learning exact models of mobile applications. We consider the PLTS model (Parameterised Labelled Transition System) to capture the different events made on GUIs. PLTS states also capture all the observed screen contents and notifications about the modifications of the application environment. These notifications signal system events e.g., local database modifications or remote server calls. All this amount of data provide a rich expressiveness that is used while learning the model and that may be later considered for precise model analysis. To avoid a state space explosion, our algorithm dynamically builds state equivalence classes while testing. Each time a new state is discovered, it dynamically re-adjusts the state equivalence relation and classes to limit the state set. These equivalence classes also help recognise similar states that do not require to be explored. Like some available tools [7, 8], our algorithm can also detect application crashes and create test cases for replaying bugs.

We proceed as follows: Section 1.2 briefly presents some related work before introducing an overview of our algorithm that we apply on a straightforward Android application example in Section 1.3. We define the model, the state equivalence relation, and we provide the model inference algorithm in Section 1.4. We give an empirical evaluation on Android applications in Section 1.5 and conclude in Section 1.6.

1.2 Related Work

Several papers dealing with automatic testing and model generation approaches of black-box systems were issued in the last decade. Due to lack of room, we only

present some of them relative to our work. Memon et al. [9] initially presented *GUI Ripper*, a tool for scanning desktop applications. This tool produces event flow graphs and trees showing the GUI execution behaviours. Only the click event can be applied and *GUI Ripper* produces many false event sequences which may need to be weeded out later. Furthermore, the actions provided in the generated models are quite simple (no parameters). Mesbah et al. [10] proposed the tool *Crawljax* specialised in Ajax applications. It produces state machine models to capture the changes of DOM structures of the HTML documents by means of events (click, mouseover, etc.). To avoid the state explosion problem, state abstractions must be given manually to extract a model with a manageable size. Furthermore, the concatenation of identical states proposed in [10] is done in our work by minimisation.

Google's *Monkey* [7] is a random testing tool that is considered as a reference in many papers dealing with Android application automatic testing. However, it cannot simulate complex workloads such as authentication, hence it offers light code coverage in such situations. *Dynodroid* [8] is an extension of Monkey supporting system events. No model is provided. Amalfitano et al. [1] proposed *AndroidRipper*, a crawler for crash testing and for regression test case generation. A simple model, called GUI tree, depicts the observed screens. Then, paths of the tree not terminated by a crash detection, are used to re-generate regression test cases. Yang et al. proposed the tool *Orbit* [12] whose novelty lies in the static analysis of Android application source code to infer the events that can be applied on screens. Then, a classical crawling technique is employed to derive a tree labelled by events. to generate approximate models. The algorithm is composed of a testing engine which executes applications to check if event sequences meet the model under generation until a counterexample is found. An active learning algorithm repeatedly asks the testing engine observation sequences to infer and eventually regenerate the model w.r.t. all the event and observation sequences.

To prevent from a state space explosion, the approaches [9, 10, 12] require state-abstractions given by users and specified in a high level of abstraction. Choi et al. [4] prefer using the approximate learning algorithm L^* . These choices are particularly suitable for inferring models for comprehension aid, but these models often are over approximations and given in a high level of abstraction, which may lead to many false positives with test case generation. In this paper, we focus on the inference of exact models. As in [10, 1], we consider the notion of state abstraction that we formally define to limit the state space domain to be explored. But, our algorithm also dynamically re-adjusts state equivalence classes to restrain the exploration and constructs a state abstraction according to the content of the application.

1.3 Overview

In the following, we present an overview on our model inference algorithm. Beforehand, we give some assumptions on mobile applications considered to design our approach:

Mobile application testing: we consider black-box applications which can be exercised through screens. It is possible to dynamically inspect application states to collect Widget properties. The set of UI events enabled on a screen should be collected as well. If not, Widgets provide enough information (type, etc.) to determine the set of events that may be triggered. Furthermore, any new screen can be observed and inspected (including application crashes). The application environment modifications (databases, network traffic, etc.) can be observed with probes,

Application reset: we assume that mobile applications and their environments (database, remote servers or mocked servers, Operating Systems) can be reset,

Back mechanism availability: several operating systems or applications (Web navigators, etc.) also propose a specialised mechanism, called the *back mechanism* to let users going back to the previous state of an application by undoing its last action. We do not consider that this mechanism is necessarily available and, if available, we assume that it does not always allow to go back to the previous state of an application (wrong implementation, unreachable state, etc.). Most of the other methods assume that the back mechanism always works as expected [1, 8], but this is frequently not the case.

1.3.1 Terminology

Mobile applications depict screens which represent application states, the number of states being potentially infinite. Screens are built by application components; here we take back the notation used with Android applications, i.e. *Activities*. The later display screens by instantiating Widgets (buttons, text fields, etc.) which are organised into a tree structure. They also declare the available events that may be triggered by users (click, swipe, etc.). A Widget is characterised by a set of properties (colour, text values, etc.). Hence, one Activity can depict several screens, composed of different Widgets or composed of the same Widgets but having different properties.

Figure 1.1 depicts the screens of an Android application example used throughout the paper. This application converts colour formats from RGB to HSL (hue-saturation-lightness) and vice-versa by means of two radio buttons *r1* and *r2*. When the button *Convert* is pressed, the value entered in the blank text field *txt* is converted and the result appears in the red text field *result*. The chosen colour is also displayed in a colour-box which is depicted at the screen bottom. This application is composed of one Activity which can display an infinite number of screens composed of different text fields values and colour-boxes.

1.3.2 Algorithm overview

Figure 1.2 introduces an overview of our algorithm which is composed of two parts. The algorithm is framed on the task-pool paradigm (Figure 1.2(a)). Tasks are placed



Fig. 1.1 Colour Converter Android application

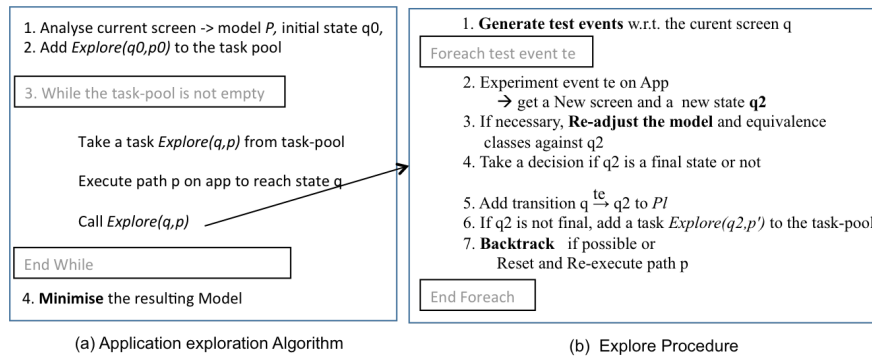


Fig. 1.2 Overview of the Model inference algorithm

into the task-pool, implemented as an ordered list, and each can be executed in parallel. A task $Explore(q, p)$ corresponds to one screen to explore. A screen is transcribed by the state q gathering all the Widget properties composing the screen and p is a path allowing to reach q from the initial state q_0 . When there is no more task to do, the exploration implicitly ends. The resulting model is then minimised to be more readable.

The exploration of one state (Figure 1.2(b)) is done by the Explore procedure. A set of test events (parameter values combined with an event set), which match the current application state, is firstly generated. The current screen is experimented with every test event to produce new screens. However, this step may lead to an infinite set of states to explore. To avoid this well-known issue, the algorithm slices the state space domain into a finite state equivalence class set by means of an equivalence relation (defined in Section 1.4.1). A state which belongs to a previously discovered equivalence class is marked as final otherwise it has to be explored. Intuitively, for every new built state q_2 (step 2), the algorithm eventually readjusts the

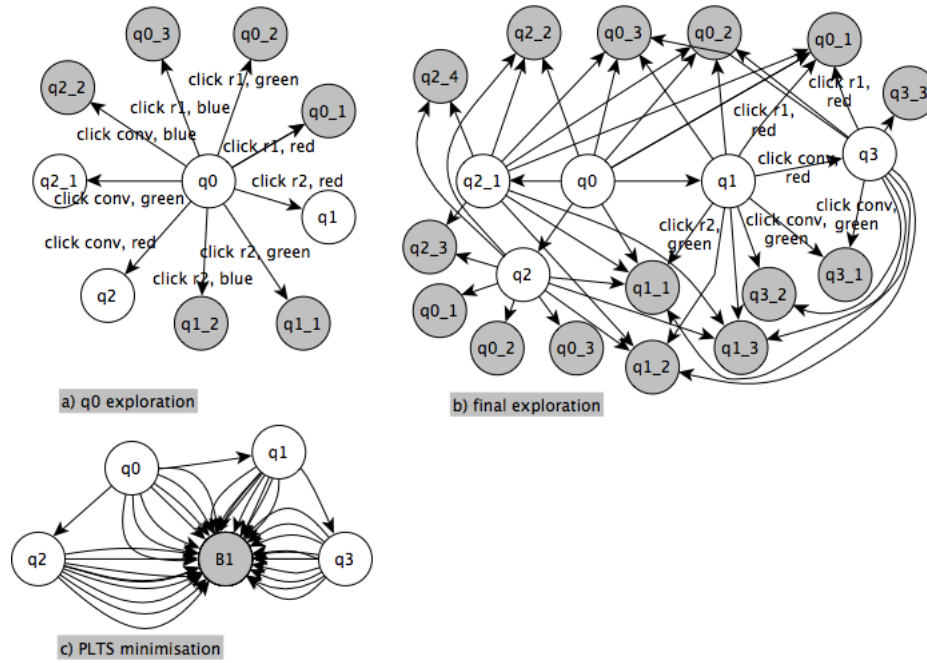


Fig. 1.3 PLTS generation example

state abstraction to limit the state set size (step 3). It scans the detected equivalence classes and checks if some of them (three or more in the algorithm) are different only on account of one Widget property. If so, it has detected a Widget property which may lead to the construction of several equivalence classes and states to explore (up to an infinite set of states). Consequently, it readjusts the equivalence relation, classes and the model by masking this Widget property. This means that this property is no more taken into account for the equivalence class computation. Therefore, the new state q_2 belongs automatically to an already discovered equivalence class and so, it will not be explored. No new equivalence class is built either. Then, the algorithm checks if new states have to be explored (step 4). Finally, the algorithm tries to backtrack the application to go back to its previous state by undoing the previous action. If it doesn't work, the application and its environment (OS, databases, etc.) are reset and the previous path p is used to reach the state which is currently under exploration.

Figure 1.3 illustrates with simplified graphs (no PLTSs) how the algorithm works on the example of Figure 1.1. For simplicity, only three values are considered for testing: colour red ($rgb=255,0,0$ or $hsv=0,100,50$), green and blue. We also assume that these values are always used for testing in the same order. The equivalence relation is: *two states are equivalent if they have the same Widget properties, except those related to text field values*. These last properties are usually not considered for conceiving state abstractions since these often lead to a large potentially infinite

set of states. Furthermore, if a Widget property takes more than two values in the different equivalence classes then the relation has to be re-adjusted.

1. Initially, we have a state q_0 which corresponds to the beginning of the application (Figure 1.1(a)) and the corresponding equivalence class $[q_0]$. A list of test events is generated from q_0 : the events click on the radio-buttons r_1 , r_2 or on the button $conv$ are combined with the colour values that are injected into the blank text field txt . The application is firstly experimented with the click on r_1 and with the red colour value. This produces a new screen and a state $q_{0.1}$ which belongs to the same equivalence class $[q_0]$ because only the text field txt is modified. This new state is marked as final (in grey in Figure 1.3) and is not explored. The application is backtracked to return to q_0 . With the other colours, we also reach final states $q_{0.2}$ and $q_{0.3}$ which are marked as final because they belong to the same equivalence class $[q_0]$.

Then, the radio-button r_2 is clicked, with the red colour. We obtain a new state q_1 (Figure 1.1(b)) and a new equivalence class $[q_1]$ since r_2 is now enabled. Therefore, we get a new task $Explore(q_1, p')$. Once more, the application is backtracked. When using the other colour values, we obtain the states $q_{1.1}$ and $q_{1.2}$ that are marked as final since these belong to $[q_1]$. No task is created.

When the $conv$ button is clicked, a value appears in the text field $result$ and a colour is depicted in the colour-box (Figure 1.1(c)). We obtain a state q_2 , which has to be explored, and a new equivalence class $[q_2]$. Next, $conv$ is clicked with the green colour. The state $q_{2.1}$ is built with a new equivalence class $[q_{2.1}]$. This state is not marked as final since the colour-box displays a new colour. This process should continue for every colour and in particular with the blue one, which produces the state $q_{2.2}$. A state space explosion may happens here. But the algorithm detects that three equivalence classes are different only on account of the same property $colourbox.colour$. The algorithm readjusts the equivalence relation to limit the state set size. Intuitively, the equivalence relation becomes *two states are equivalent if they have the same Widget properties, except those related to text field values and to the colourbox.colour property*. Then, it updates states and equivalence classes to match this new relation. As a consequence, $[q_2]$, $[q_{2.1}]$ and $[q_{2.2}]$ are now merged into $[q_2]$. The new state $q_{2.2}$ now belongs to an existing equivalence class and is hence marked as final. The first task $Explore(q_0, p)$ is finished and we obtain the graph depicted in Figure 1.3a,

2. we assume that the task $Explore(q_1, q_0 \xrightarrow{\text{click } r_2, txt=red} q_1)$ is picked out to explore q_1 . A list of test event, which is the same as previously is constructed. From the state q_1 , when the button $conv$ is clicked with the red colour value, a new state q_3 is added because the colour box appears. When $conv$ is clicked with other colour values and events we only obtain final states, since they belong to previously discovered equivalence classes,
3. the same reasoning is followed on states q_2 , $q_{2.1}$ and q_3 , but only final states are added (no task). We obtain the PLTS of Figure 1.3b,
4. the task-pool is empty. The PLTS is finally minimised [6]. Here, the final states are merged to one unique state as illustrated in Figure 1.3c.

In this short example, we have shown that our algorithm avoids the state explosion problem and ends once at least one state of all the detected equivalence classes is explored. In the following, we describe formally the model, the equivalence relation, and the algorithm.

1.4 Model inference Algorithm

1.4.1 Mobile application modelling with PLTS

We use PLTSs as models for representing mobile applications. A PLTS is a kind of state machine extended with variables and guards on transitions. Beforehand, we assume that there exist a domain of values denoted D and a variable set X taking values in D . The assignment of variables in $Y \subseteq X$ to elements of D is denoted with a mapping $\alpha : Y \rightarrow D$. We denote D_Y the assignment set over Y .

Definition 1 (PLTS). A PLTS (Parameterised Labelled Transition System) is a tuple $\langle V, I, Q, q_0, \Sigma, \rightarrow \rangle$ where:

- $V \subseteq X$ is the finite set of variables, $I \subseteq X$ is the finite set of parameters used with actions,
- Q is the finite set of states, such that a state $q \in Q$ is an assignment over D_V , q_0 is the initial state composed of the initial condition D_{V_0} ,
- Σ is the finite set of valued actions $a(\alpha)$ with $\alpha \subseteq D_I$,
- $\rightarrow \subseteq Q \times \Sigma \times Q$ is the transition relation. A transition $(q, a(\alpha), q')$ is also denoted $q \xrightarrow{a(\alpha)} q'$.

The behaviour of a PLTS P is characterised by its sequences of valued actions starting from its initial state q_0 . These sequences are also called the traces of P :

Definition 2 (PLTS Traces). Let $P = \langle V, I, Q, q_0, \Sigma, \rightarrow \rangle$ be a PLTS. $Traces(P) = Traces(q_0) = \{a_1(\alpha_1) \dots a_n(\alpha_n) \mid \exists q_1, \dots, q_n, q_0 \xrightarrow{a_1(\alpha_1)} q_1 \dots q_{n-1} \xrightarrow{a_n(\alpha_n)} q_n \in (\rightarrow)^*\}$.

We model mobile application behaviours with PLTSs by encoding (UI) events with actions. We also store the properties collected from screens (Widget properties) and notifications about the application environment changes in states with variable assignments:

UI events representation

We interact with mobile applications by means of events, e.g. a click on a button, and by entering values into editable Widgets. We capture such events with PLTS actions of the form $event(\alpha)$ with $\alpha = \{widget := w, w_1 := val_1, \dots, w_n := val_n\}$ an

assignment over D_I ; the parameter *widget* denotes the Widget name on which is applied the event and the remaining variables are assignments of Widget properties. We denote the triggering of the back mechanism with the action $back(\alpha)$ with α an empty assignment.

Application state representation

We specialise PLTS states to store the content of screens (Widget properties) in such a way as to facilitate the construction of equivalence classes of states. We split the set of Widget properties into two categories: we gather in the set W the Widget properties that indicate a strong application behaviour modification and that take only few values e.g., Widget visibility, size, etc. The others that usually take a lot of different values such as the properties about text field values, are placed into W^c . This separation affects the state representation: we denote wp the assignment composed of properties in W , while the assignment wo is composed of the other Widget properties. A PLTS state q is then a specific assignment of the form $act \cup wp \cup wo \cup env \cup end$ where:

- act is an assignment returning an Activity name,
- (wp, wo) are two Widget property assignments. The union of wp and wo gives all the property values of an application screen displayed by act .
- env is a boolean assignment indicating whether the application environment has been modified,
- end is a boolean assignment marking a state as final or not.

For readability, a state $q = act \cup wp \cup wo \cup env \cup end$ is denoted (act, wp, wo, env, end) . This state structure eases the definition of the state equivalence relation given below:

Definition 3 (State equivalence relation). Let $P = \langle V, I, Q, q_0, \Sigma, \rightarrow \rangle$ be a PLTS and for $i = (1, 2)$ let $q_i = (act_i, wp_i, wo_i, env_i, end_i)$, be two states in Q . We say that q_1 is equivalent to q_2 , denoted $q_1 \sim q_2$ iff $act_1 = act_2$, $wp_1 = wp_2$ and $env_1 = env_2$. $[q]$ denotes the equivalence class of equivalent states of q . Q/\sim stands for the set of all equivalence classes in Q .

This definition gives a very adaptable state equivalence relation whose meaning can be modified by altering the assignments wp . If we take back our example, one can consider two states $q_1 \approx q_2$ which are different only because they do not include the same assignments of the Widget property *colourbox.colour* ($act_1 = act_2$, $env_1 = env_2$, but $wp_1 \neq wp_2$). We have two equivalence classes $[q_1], [q_2]$. The equivalence relation is adaptable in the sense that wp_i can be changed as follows: if we consider that *colourbox.colour* takes too much values and implies to much different equivalence classes, *colourbox.colour* can be shifted from wp_i to wo_i ($i = 1, 2$) in states. We obtain $wp_1 = wp_2$, $q_1 \sim q_2$ and only one equivalence class $[q_1]$. Intuitively, our algorithm uses this adjustment process to dynamically reduce the equivalence class domain and the state exploration according to the screen content.

1.4.2 Model inference Algorithm

In this section, we describe more precisely the Explore procedure (second part of the overview in Figure 1.2) whose pseudo code is given in Algorithm 1. Due to lack of room, the task-pool management algorithm can be found in [11].

As stated above, this procedure aims at visiting one state to augment the PLTS under construction, denoted P , with new transitions and states and to eventually produce new tasks $Explore(q, p)$ added to the task-pool. Its steps are explained below:

- Test data generation and execution (lines 4-11): the current screen is analysed to generate a set of events expressing how to complete Widgets with values and to trigger an event. In short, our algorithm generates a set of events of the form $\{event(\alpha) \mid event \text{ is an event, } \alpha \text{ is an assignment}\}$. It starts collecting the events that may be applied on the different Widgets of the current screen. Then, it constructs assignments of the form $w_1 = v_1 \wedge \dots \wedge w_n = v_n$, with (w_1, \dots, w_n) the list of editable Widget properties found on the screen and (v_1, \dots, v_n) , a list of test values. Instead of only using random values, we propose to use several data sets: a set *User* gathering values manually chosen such as logins and passwords, a set *RV* composed of values well known for detecting bugs e.g., String values like "&", "", or null, and of random values. A last set, denoted *Fakedata*, is composed of fake user identities. Furthermore, we adopted a Pairwise technique [5] to derive a set of assignment tuples over these data sets. Assuming that errors can be revealed by modifying pairs of variables, this technique strongly reduces the coverage of variable domains by constructing discrete combinations for pair of parameters only. Then, each $event(\alpha)$ is applied on the current screen to produce new ones (application crash included). Each screen is analysed to retrieve Widget properties and the activity which produces this screen. Probes are requested to detect if the application environment were modified. These data are formalised by the state q_2 ,
- Model readjustment: the Explore procedure now checks whether the re-adjustment of P and of the state equivalence classes is required (lines 9-12). We denote $C^{Wprop}(Q/\sim)$ the number of assignments of the same Widget property $Wprop$ found in the set of equivalence classes Q/\sim . $C^{Wprop}(Q/\sim) = card(\{\alpha = (Wprop := val) \mid [q] \in Q/\sim, q = (act, wp, wo, env, end), \alpha \in wp\})$. For each assignment $\alpha = (Wprop := val)$ in wp_2 , we check how much values the Widget property $Wprop$ takes in the equivalence classes: if $Wprop$ takes more than 2 values in Q/\sim (if $card(C^{Wprop})(Q/\sim) > 2$), then we re-adjust the state representation. In every state $q = (act, wp, wo, env, end)$ of $Q \cup \{q_2\}$, the assignments of the form $(Wprop := val)$ are shifted from wp to wo (procedure Readjust in Algorithm 1 line 11). The equivalence classes are also transformed in accordance (procedure Readjust line 12),
- PLTS completion: a new transition $q \xrightarrow{event(\alpha)} q_2$ is added to the PLTS P (lines 13-20). q_2 is marked as final if q_2 belongs to an existing equivalence class. Otherwise (line 17), q_2 has the assignment $(end := false)$ and a new task $Explore(q_2, p')$

Algorithm 1: Explore Procedure

```

1 Procedure Explore( $q, p$ );
2  $Events = GenEvents$ , analyse the current screen to generate the set of events
  { $event(\alpha) \mid event$  is an event,  $\alpha$  is an assignment};
3 foreach  $event(\alpha) \in Events$  do
4   Experiment  $event(\alpha)$  on  $App \rightarrow$  new screen  $Inew$ ;
5   Analyse  $Inew \rightarrow$  assignments  $act_2, wp_2, wo_2$ ;
6   Analyse the application environment  $\rightarrow env_2$ ;
7    $q_2 = (act_2, wp_2, wo_2, env_2, end := null)$ ;
8   foreach  $\alpha = \{Wprop := val\} \in wp_2$  do
9     if  $card(C^{Wprop}(Q/\sim) \cup \{\alpha\}) > 2$  then
10      Read just( $Q \cup \{q_2\}, Wprop$ );
11      Read just( $Q/\sim, Wprop$ );
12   if  $Inew$  reflects a crash or there exists  $[q'] \in Q/\sim$  such that  $q_2 \in [q']$  then
13     {Add a transition  $q \xrightarrow{event(\alpha)} q_2 = (act_2, wp_2, wo_2, env_2, end := true)$  to  $\rightarrow_P$ ;
14     } (in critical section)
15   else
16     {Add a transition  $t = q \xrightarrow{event(\alpha)} q_2 = (act_2, wp_2, wo_2, env_2, end := false)$  to  $\rightarrow_P$ ;
17      $Q/\sim = Q/\sim \cup \{q_2\}$ ;
18     Add the task (Explore( $q_2, p.t$ )) to the task-pool;
19     } (in critical section)
20   Backtrack( $q_2, q, p$ );

```

is added to the task pool. Since the algorithm is highly parallelisable, we use critical sections to modify the PLTS P (which is shared among threads),

- **Application backtracking:** to apply the next event, the Explore procedure calls the Backtrack one (line 21) to reach the previous screen and state q . Its algorithm is given in Algorithm 2. Here the notion of application environment really makes a difference to achieve an exact model: if the current state q_2 has an assignment ($env := false$), it reflects the fact that the application environment has not been modified, therefore the Backtrack procedure calls the back mechanism to undo the most recent action (if available). We observe a new screen and check whether it is equivalent to the previous screen stored in q (we compare their Widget properties). Otherwise, the application and its environment are reset and we re-execute the path p to reach the state q (Algorithm 2, line 7) (here, we assume that the application is deterministic though).

1.4.3 PLTS minimisation

Our algorithm performs a minimisation on the first generated PLTS to achieve a more readable model. We have chosen a bisimulation minimisation technique since this one still preserves the functional behaviours represented in the original model while reducing the state space domain. A detailed algorithm can be found in [6]. In short, this algorithm constructs sets (blocks) of states that are bisimilar equivalent (any action from one of them can be matched by the same action from the other and

Algorithm 2: Backtrack procedure Algorithm

```

1 Procedure Backtrack( $q_2 = (act2, wp2, wo2, env2, end2), q = (act, wp, wo, env, end), p$ );
2 if  $env2 = (env := false)$  and the back mechanism is available then
3   Call the back mechanism  $\rightarrow$  screen  $I_{New}$ ;
4   Analyse  $I_{New} \rightarrow$  assignments  $rc', wp', wo'$ ;
5   Analyse the application environment  $\rightarrow env'$ ;
6   if  $act \neq act'$  or  $wp \neq wp'$  or  $wo \neq wo'$  or  $env \neq env'$  then
7      $\lfloor$  Reset and Execute  $App$  by covering the actions of  $p$ ;
8   else
9      $\lfloor$  Add a transition  $t = q_2 \xrightarrow{back(\alpha)} q$  to  $\rightarrow_{Tree}$ ;
10 else
11  $\lfloor$  Reset and Execute  $App$  by covering the actions of  $p$ ;

```

the arrival states are again bisimilar). Figure 1.3c depicts the (simplified) minimised PLTS of the application example. Here, final states are aggregated into one block of states.

1.4.4 Algorithm correctness, complexity and termination

We express the correctness of our model inference method in term of trace equivalence between the inferred PLTS and the traces of the application under test:

Proposition 1. *Let P be a PLTS constructed with our model inference algorithm from a deterministic mobile application App . We have $Traces(P) \subseteq Traces(App)$.*

The proof is given in [11]. Intuitively, our algorithm constructs a PLTS P with these steps:

1. *Generation of PLTS:* from a given state q , every new event applied on the application is modelled with a unique transition whose arrival state q_2 is new or final. We do not merge states and hence we construct a PLTS P ,
2. *Correct use of the back mechanism:* we call this mechanism with care: it is called only if the environment of the application (databases, remote servers, etc.) were not modified with the execution of the last action. Indeed, if we apply the back mechanism even so, we necessarily reach a new state since the application environment is modified. Secondly, we check if the state of the application obtained after the call of the back mechanism is really the previous state of the application. If one of these conditions is not met, we reset the application and its environment and we re-execute the path p to reach the state q ,
3. *Minimisation with trace equivalence:* we apply a bisimulation minimisation technique to produce a PLTS MP from P such that the two PLTS are bisimilar and consequently trace equivalent as well.

Complexity and termination of the Algorithm: our algorithm builds at most $2 * 2^n$ equivalence classes, with n the number of Widget properties in W . In short,

we can have two different ($env := true$, $env := false$) and m^n different assignments over W if m is the maximum number of values that any Widget property can take. Nonetheless, when a property of W takes more than two values, our algorithm shifts it from the assignment wp to wo in states. Furthermore, since we explore one state per equivalence class, the algorithm ends and we have $2 * 2^n$ equivalence classes and not final states. We also have at most nm transitions (Pairwise testing [5]) for each. If N and M stand for the number of not final states and transitions, the whole algorithm has a complexity proportional to $\mathcal{O}(M + N + MN + M \log(N))$. Indeed, the Explore procedure covers every transition twice (one time to execute the event and one time to go back to the previous state) and every not final state is processed once. But, sometimes the back mechanism is not available. In this situation, the application is reset to go back to a state q by executing the events of a path p at worst composed of M transitions. In the worst case, this step is done for every state with a complexity proportional to NM . Furthermore, the minimisation procedure has a complexity proportional to $\mathcal{O}(M \log(N))$ [6].

1.5 Empirical Evaluation

We present here some experimentations on Android applications to answer on the following questions: does the algorithm offer good code coverage in a reasonable time delay? How are the models in terms of size and quality for analysis?

We have implemented our algorithm in a tool called *MCrawlIT* (Mobile Crawler Tool ¹). It takes packaged applications or source projects and user data e.g., logins and passwords required for the application execution. *MCrawlIT* is based on the testing framework *Robotium* ² which retrieves the Widget properties of a screen and simulates events.

To avoid any bias, we compare the effectiveness of *MCrawlIT* with the following available tools, *Monkey* [7] and *Dynodroid* [8], on applications taken as reference in the papers [12, 9, 4, 8] and whose source code is available (30 applications). The results of some other tools *Orbit* [12], *Guitar* [9] and *Swifthand* [4] are taken from the papers. It is important to note that *Monkey* is taken as a reference in most of the papers dealing with Android testing. Thereby, our results can be compared with other studies related to Android testing.

Code coverage: Table 1.1 reports the percentages of code coverage obtained with the different tools on 30 applications with a time budget of three hours. If we do a side by side comparison of *MCrawlIT* with the other tools, we observe that *Monkey* provides better code coverage for 8 applications, *SwiftHand* for 2 and *Dynodroid* for 5. In comparison to all the tools together, *MCrawlIT* provides better code coverage for 20 applications, the coverage difference being higher than 5% with 13 applications. These results show that *MCrawlIT* gives better code coverage

¹ available here <https://github.com/statops/mcrawlert.git>

² <https://code.google.com/p/robotium/>

Applications	Monkey	Orbit	Guitar	MCrawlT	SwiftHand	Dynodroid
NotePad	60	82		88		crash
Tippy_TipperV1	41	78		79		48
ToDoManager	71	75	71	81		34
OpenManager	29	63		65		crash
HelloAUT	71	86	51	96		76
TomDroid	46	70		76		42
ContactManager	53	91	71	68		28
Aardict	52	65		67		51
Musicnote	69			81	72.2	47
Explorer	58			74	74	crash
Myexpense	25			61	41.8	40
Anynemo	61			54	52.9	crash
Whohas	58			95	59.3	65
Mininote	42			26	34	39
Weight	51			34	62	56
Tippy_TipperV2	49			74	68	12
Sanity	8			26	19.6	1
Nectdroid	70.7			54		68.6
Alogcat	66.6			66		67.2
ACal	14			46		23
Anycut	67			71		69.7
Mirrored	63			76		60
Jamendo	64			46		3.9
Netcounter	47			56		70
Multisms	65			73		77
Alarm	77			72		55
Bomber	79			75		70
Adsroid	72			83		80
Aagtl	18			25		17
PasswordFor Android	58			61		58

Table 1.1 Code coverage (in %)

than the other tools and even offers good results against all the tools together on half the applications with comparable execution times. Table 1.1 also reveals that the obtained code coverage percentage is between 25% and 96%. We manually analysed the 8 applications which yield the less good results with *MCrawlT* to identify the underlying causes behind low coverage. This can be explained at least by these ways:

- Specific functionalities and unreachable code: several applications are incompletely covered either on account of unused code parts (libraries, packages, etc.) that are not called, or on account of functionalities difficult to start automatically,
- Unsupported events: several applications e.g., Nectdroid, Multism, Acal or Alogcat chosen for experimentation with *Dynodroid* take UI events as inputs but also system events such as *Android broadcast messages*. Our tool does not support these events yet. Moreover, *MCrawlT* only supports the event list also supported by the testing tool *Robotium* (viz. click and scroll). The long click is thus not supported but is used in some applications (Mininote and Contactmanager). In contrast, *Orbit* supports this event and therefore offers a better code coverage with the application Contactmanager.

Quality and size of the models: Table 1.2 finally shows the number of states obtained with *MCrawlT*, *Orbit* [12] and *SwiftHand* [4] since they produce models

Applications	#PLTS states (MCrawlT)	#states after minimisation (MCrawlT)	#states (Orbit)	#states (Swift-Hand)
NotePad	13	8	7	
Tippy_TipperV1	37	18	9	
ToDoManager	6	2		
OpenManager	31	12	20	
HelloAUT	8	5	8	
TomDroid	12	6	9	
ContactManager	5	4	5	
Sanity	31	24		78
Musicnote	41	23		46
Explorer	96	74		195
Myexpense	52	37		149
Anynemo	139	106		169
Whohas	36	11		97
Mininote	45	19		169
Tippy_TipperV2	54	26		71
Weight	69	23		109

Table 1.2 Inferred model size

as well. Before minimisation, our tool generates larger and tacitly less comprehensive models than those obtained with *Orbit*. In term of quality of the learned models, we do not produce extrapolated models and we believe that those generated by *MCrawlT* offer more testing capabilities. Indeed, these models include states which store all the observed Widget properties (colours, texts, etc.) and notifications about the application environment changes. We have precisely chosen this feature to later perform test case generation. For instance, with this amount of information, we can construct test cases to apply events and to check the content of the resulting screen but also if remote servers are called, etc. Both *Orbit* and *SwiftHand* only store UI events. After minimisation, we obtain more compact and readable models whose sizes are comparable to the sizes of the models obtained with *Orbit*. This tends to show that our approach of producing larger but more detailed models that are after minimised, only offer advantages for model inference. In addition, *MCrawlT* constructs storyboards from these minimized models by replacing states with screenshots of the application.

All these experimental results on real applications tend to show that our tool is effective and can be used in practice since it produces equivalent or higher code coverages than the other tools.

1.6 Conclusion

In this paper, we present an algorithm which infers PLTS models from mobile applications. It constructs PLTSs that capture events and all the Widget properties extracted from the observed screens. Despite the huge amount of collected data, we avoid the state space explosion problem by using an equivalence relation and classes that are dynamically re-adjusted all along the algorithm execution with re-

gards to the screen content. Our experimental results show that our algorithm offers good code coverage quickly and can be used in practice. Furthermore, the generated models can be reused for precise model analysis. An immediate line of future work would be to apply this kind of algorithm for security breach detection. The exploration could be specialised to target some specific application parts (login step, etc.). Then, test cases could be automatically generated from test patterns to further explore specific states with the purpose of improving detection.

References

1. Amalfitano, D., Fasolino, A., Tramontana, P.: A gui crawling-based technique for android mobile application testing. In: Software Testing, Verification and Validation Workshops (ICSTW), 2011 IEEE Fourth International Conference on, pp. 252–261 (2011). DOI 10.1109/ICSTW.2011.77
2. Anand, S., Naik, M., Harrold, M.J., Yang, H.: Automated concolic testing of smartphone apps. In: Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering, FSE '12, pp. 1–11. ACM, New York, NY, USA (2012). DOI 10.1145/2393596.2393666
3. Artzi, S., Kiezun, A., Dolby, J., Tip, F., Dig, D., Paradkar, A., Ernst, M.: Finding bugs in web applications using dynamic test generation and explicit-state model checking. Software Engineering, IEEE Transactions on **36**(4), 474–494 (2010). DOI 10.1109/TSE.2010.31
4. Choi, W., Necula, G., Sen, K.: Guided gui testing of android apps with minimal restart and approximate learning. SIGPLAN Not. **48**(10), 623–640 (2013). DOI 10.1145/2544173.2509552. URL <http://doi.acm.org/10.1145/2544173.2509552>
5. Cohen, M.B., Gibbons, P.B., Mugridge, W.B., Colbourn, C.J.: Constructing test suites for interaction testing. In: Proc. of the 25th International Conference on Software Engineering, pp. 38–48 (2003)
6. Fernandez, J.C.: An implementation of an efficient algorithm for bisimulation equivalence. Science of Computer Programming **13**, 13–219 (1989)
7. Google: Ui/application exerciser monkey. [Http://developer.android.com/tools/help/monkey.html](http://developer.android.com/tools/help/monkey.html), last accessed jan 2015
8. Machiry, A., Tahiliani, R., Naik, M.: Dynodroid: An input generation system for android apps. In: Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2013, pp. 224–234. ACM, New York, NY, USA (2013). DOI 10.1145/2491411.2491450
9. Memon, A., Banerjee, I., Nagarajan, A.: Gui ripping: Reverse engineering of graphical user interfaces for testing. In: Proceedings of the 10th Working Conference on Reverse Engineering, WCRE '03, pp. 260–. IEEE Computer Society, Washington, DC, USA (2003)
10. Mesbah, A., van Deursen, A., Lenselink, S.: Crawling Ajax-based web applications through dynamic analysis of user interface state changes. ACM Transactions on the Web (TWEB) **6**(1), 1–30 (2012)
11. Salva, S., Laurençot, P., Zafimiharisoa, S.R.: Model inference of mobile applications with dynamic state abstraction. Tech. rep., LIMOS, <http://sebastien.salva.free.fr/RR-15-01.pdf> (2015)
12. Yang, W., Prasad, M.R., Xie, T.: A grey-box approach for automated gui-model generation of mobile applications. In: Proceedings of the 16th international conference on Fundamental Approaches to Software Engineering, FASE'13, pp. 250–265. Springer-Verlag, Berlin, Heidelberg (2013). DOI 10.1007/978-3-642-37057-1_19