

AUTOMATED MODEL-BASED ANDROID GUI TESTING USING MULTI-LEVEL GUI COMPARISON CRITERIA

happyBiker, haxxorman

Flatulenza, University of Prone

homemadebike.9999999@stud.unipronax.fritt, lexusGTsportage.99999999@stud.unipronax.fritt

January 25, 2017

- 1 INTRODUZIONE
 - GUI Android



GUI TESTING AUTOMATICO(1)

PERCHÈ DEDICARE PARTICOLARE ATTENZIONE ALLE APPLICAZIONI ANDROID?

Attualmente le applicazioni android comprendono quasi il 90% del mercato di software mobile.

COS'È IL GUI TESTING?

È una metodologia di testing che esplora una parte dello spazio degli stati del programma testato dando gli input al programma tramite l'interfaccia grafica.

PERCHÈ AUTOMATIZZARE IL TESTING?

Il motivo principale è il costo, l'esecuzione manuale dei test ha il costo eccessivo per molte applicazioni. Inoltre l'esecuzione manuale dei test richiede tanto tempo e a volte non è in grado di individuare tutti gli errori.

GUI TESTING AUTOMATICO(2)

Attualmente sul mercato si trovano dei tool per il GUI testing automatico che generano gli input in modo aleatorio, tentando di generare una sequenza di input che provoca il malfunzionamento del software. Nel testo viene citato un tool chiamato *Android Monkey* che permette questo tipo di test.

Dato che si tratta essenzialmente di una ricerca locale nello spazio degli stati, questi metodi sono soggetti alle problematiche degli algoritmi di ricerca locale. Per esempio possono valutare lo stesso stato più volte.

L'approccio proposto dagli autori è di costruire un grafo degli stati della GUI, utilizzando dei criteri di equivalenza (GUICC) per distinguere gli stati.

MODELLO GUI

Un altro approccio prevede l'uso del modello della GUI. Il modello GUI è un automa a stati finiti, dove gli stati dell'automa corrispondono agli stati dell'interfaccia grafica. L'automa transisce da uno stato all'altro al verificarsi di un evento. Il modello GUI viene rappresentato da un multigrafo dove i nodi corrispondono agli stati e gli archi corrispondono alle transizioni. Ogni arco è etichettato dall'evento che provoca la transizione. Il concetto fondamentale necessario per la costruzione di questo grafo è il criterio di equivalenza degli stati.

Le applicazioni Android spesso generano parte dell'interfaccia grafica in modo dinamico. Per modellare correttamente la GUI bisogna tener conto delle componenti dinamiche.



Un criterio d'equivalenza debole come "Nome Activity" in questo caso non sarebbe in grado di accorgersi della differenza tra gli stati.

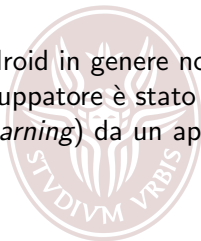
MULTI-LEVEL GUI COMPARISON CRITERIA (MULTI-LEVEL GUICC)

La struttura a livelli definita dagli autori del paper è composta da:

- **Nome pacchetto.** Quando l'activity corrente non appartiene allo stesso pacchetto dell'applicazione testata significa che l'applicazione testata è stata terminata oppure ha invocato un'altra applicazione, in ogni caso lo stato corrente è considerato lo stato terminale.
- **Nome Activity.** Quando l'activity corrente ha un nome diverso da quella precedente si tratta di uno stato diverso.
- **Composizione Widget non-eseguibili**
- **Composizione Widget eseguibili**
- **Contenuti**

GENERAZIONE DEL MODELLO

Dato che le applicazioni Android in genere non contengono un modello della GUI costruito dallo sviluppatore è stato sviluppato un metodo per dedurre il modello (*Model learning*) da un applicazione tramite *reverse engineering* (*GUI Ripping*).



GENERAZIONE DEL MODELLO(2)

Per costruire il grafo della GUI vengono eseguiti iterativamente questi passaggi:

- Generare un evento ammissibile nello stato corrente.
- Verificare tramite i criteri d'equivalenza se lo stato della GUI è cambiato. In caso affermativo ci sono tre casi:
 - Il nuovo stato è già presente come nodo nel grafo. Allora si aggiunge un arco dallo stato precedente al nuovo stato etichettato dall'evento che porta in questo stato.
 - Non esiste un nodo corrispondente, allora abbiamo scoperto uno stato nuovo. Viene generato un nuovo nodo e generato l'arco come nel primo caso.
 - Lo stato risultato è uno stato terminale, allora l'applicazione è terminata in seguito all'evento. (Terminazione normale o anormale dovuta ad un errore/eccezione).

In caso contrario si aggiunge un *loop* (un arco che origina e termina nello stesso nodo) etichettato dall'evento.

ARCHITETTURA DEL SISTEMA

MOTORE DI TESTING

Viene eseguito lato controllore.

- **Esecutore test** Sceglie l'input da mandare al software testato secondo un algoritmo(BFS).
- **Generatore del grafo** descritto sopra.
- **Error checker** Controlla la presenza di errori.

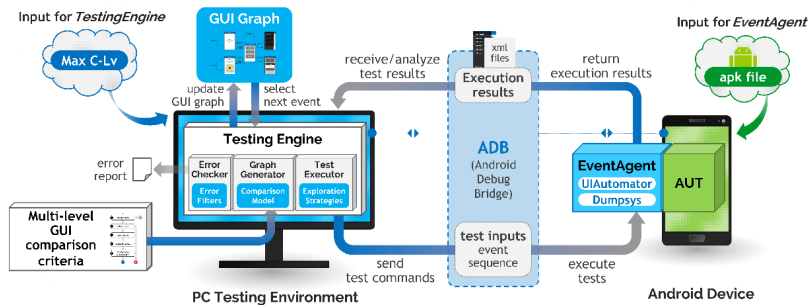
STRATO DI COMUNICAZIONE

Serve per trasferire informazioni fra la macchina controllore e il dispositivo Android che esegue il software testato.

EVENT AGENT

Software eseguito sul dispositivo Android. Riceve messaggi dal controllore e genera input al software testato corrispondenti.

GUI TESTING 1 : FUNZIONAMENTO DEL FRAMEWORK



GUI TESTING 2 : INPUT E OUTPUT

figura in sfondo Questo framework riceve in input:

- file **Android APK** (eseguibile dell'App);
- **Max C-Lv**

e restituisce in output:

- un insieme di **sequenze di eventi**, eseguite sull'applicazione ;
- un **grafo della GUI**, rappresentante gli stati esplorati e testati dell'app;
- un **error report**, se il framework ha riscontrato qualche errore nell'app ,

che è di fatto il fine ultimo del testing.



OBIETTIVO: SCOVARE L'ERRORE

Automatizzare il processo di error-detection nella GUI di una app porta vantaggi :

- ❶ di tempo, perchè non bisogna testare manualmente ogni possibilità;
- ❷ economici, perchè si evita che l'utente finale incappi in *uncovered errors*, con ripercussioni sulla *fame* dell'app.

Lo scopo del **model-based GUI testing** è di rilevare i *runtime errors* in modo automaticizzato.

Ovviamente, migliore è il **modello** su cui si basa il GUI testing, maggiori sono le probabilità di individuare tutti i possibili *runtime errors*.

SONO UTILI I GUICC ?

DOMANDA 1

In che modo i *GUICC* influenzano la generazione del *behavior model* delle app Android ?

DOMANDA 2

I GUICC influenzano la **copertura di codice** durante la fase di GUI testing ?

DOMANDA 3

I GUICC influenzano la capacità di scovare errori durante la fase di GUI testing ?

DOMANDA 1

DOMANDA 1

In che modo i *GUICC* influenzano la generazione del *behavior model* delle app Android ?

Sono state considerate varie apps, sia open source che proprietarie, e ne sono stati generati vari GUI graphs, ognuno basato su un diverso livello di "specificità" dei GUICC .

Nella raccolta dei dati , i grafi generati con C-Lv 1 non sono stati considerati, poichè è un criterio troppo poco discriminante per ottenere risultati validi.

RECALL

C-Lv 1 : due ScreenNodes differiscono per il package di appartenenza.

GRAFI GENERATI CON VARI CRITERIA LEVELS

No	App name	C-Lv2		C-Lv4		C-Lv5	
		#SN	#EE	#SN	#EE	#SN	#EE
1	<i>Alogcat</i>	5	45	15	247	76	269
2	<i>Anycut</i>	8	33	8	33	9	42
3	<i>Mileage</i>	16	117	69	532	81	618
4	<i>Sanity</i>	1	4	2	7	49	552
5	<i>DalvikExplorer</i>	16	178	30	301	S/E	S/E
6	<i>MyLock</i>	2	24	5	51	10	101
7	<i>Bequick</i>	2	7	60	250	71	351
8	<i>Syncmypix</i>	4	11	20	96	20	115
9	<i>TippyTipper</i>	4	29	13	102	19	175
10	<i>WhoHasMyStuff</i>	7	37	24	143	26	180

No	App name	C-Lv2		C-Lv3 ~ C-Lv5		
		#SN	#EE	C-Lv	#SN	#EE
1	<i>Google Translate</i>	5	41	4	55	871
2	<i>Advanced Task Killer</i>	4	27	3	12	178
3	<i>Alarm Clock Xtreme Free</i>	4	81	4	23	363
4	<i>GPS Status & Toolbox</i>	3	12	5	49	592
5	<i>Music Folder Player Free</i>	7	37	5	30	295
6	<i>Wifi Matic</i>	2	9	5	20	160
7	<i>VNC Viewer</i>	6	43	5	7	60
8	<i>Unified Remote</i>	6	94	5	20	160
9	<i>Clipper</i>	2	17	5	36	195
10	<i>Lifetime Alarm Clock</i>	2	5	5	60	529

OSSERVAZIONI

- C-Lvs superiori generano più nodi e più archi ;
- C-Lv 2 è un criterio meno "espressivo" di C-Lv5 ;
- ATTENZIONE: app come Sanity o come DalvikExplorrer (gestione della *Dalvik virtual machine*) generano stati infiniti e il grafo esplode .

DOMANDA 1

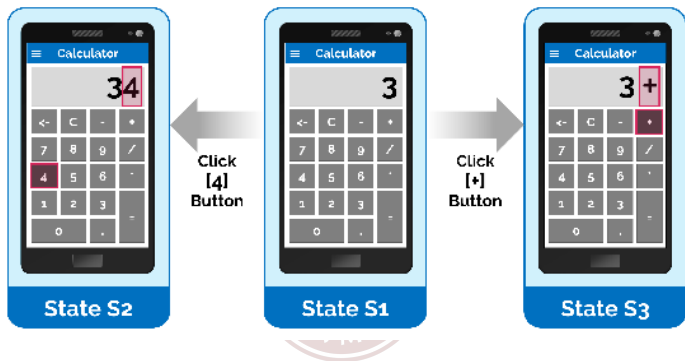
DOMANDA 1

In che modo i *GUICC* influenzano la generazione del *behavior model* delle app Android ?

RISPOSTA

Criteri più "specializzati" possono rendere il GUI states graph più completo ed esaustivo , ma livelli eccessivamente fine-grained di selezione non migliorano la qualità del grafo .

UNA SEMPLICE CALCOLATRICE ...



Una calcolatrice è una app particolare : è indistinguibile per ogni **C-Lv** precedente al quinto

Quando si inizia a comparare **in base al testo inserito** però il numero di stati esplode : ogni nuovo carattere (numero o operatore) inserito corrisponde ad un nuovo ScreenNode.

WIDER IS NOT BETTER ...

Testare applicazioni simili a Calculator è complicato , poichè più specifico è **Max C-Lv** , più grande è il grafo risultante.

Grafi molto grandi talvolta contengono ScreenNodes diversi da cui possono essere avviati eventi uguali , quindi i test perdono molto tempo su input simili.

Inoltre non è detto che con grafi più grandi si riescano a coprire tutte le funzionalità messe a disposizione da un'app.

COME MISURARE OGGETTIVAMENTE LA QUANTITÀ DI FUNZIONALITÀ TESTATE ?

Per avere una misura oggettiva del numero di funzionalità testate , si può misurare il numero di linee di codice effettivamente eseguite durante un test sul grafo .

In generale , maggiore è il numero di linee di codice eseguite , maggiore è il numero di funzionalità testate.

DOMANDA 2

DOMANDA 2

I GUICC influenzano la **copertura di codice** durante la fase di GUI testing ?

I test sono stati effettuati sempre sullo stesso insieme di app open source e proprietarie.

LINEE DI CODICE ESEGUITE NELLE APP TESTATE

No	Package name	Method coverage			Statement coverage		
		A	C-Lv	M	A	C-Lv	M
1	<i>Alogcat</i>	46%	4	65%	39%	5	56%
2	<i>Anycut</i>	23%	4	69%	19%	4	55%
3	<i>Mileage</i>	22%	5	43%	18%	5	33%
4	<i>Sanity</i>	n/a					
5	<i>DalvikExplorer</i>	65%	4	70%	57%	4	64%
6	<i>MyLock</i>	11%	4	12%	10%	4	11%
7	<i>Bequick</i>	24%	5	39%	21%	5	39%
8	<i>Syncmypix</i>	10%	4	24%	6%	4	17%
9	<i>TippyTipper</i>	42%	5	65%	36%	5	61%
10	<i>WhoHasMyStuff</i>	39%	5	62%	35%	4	51%
Average		31%	50%		27%	43%	

OSSERVAZIONI

- le colonne *A* contengono la **percentuale di copertura di codice** di un test con Max *C-Lv* pari a 2 , ossia test eseguiti con il criterio degli activity name;
- le colonne *C-Lv* e *M* vanno lette insieme : *C-Lv* indica il livello a cui viene raggiunta la massima percentuale di codice eseguito *M*;
- ci sono alcune eccezioni , dovute però a specifici casi particolari;
- si nota chiaramente che **a livelli più specializzati di GUICC corrisponde una maggior copertura di codice**

DOMANDA 2

DOMANDA 2

I GUICC influenzano la **copertura di codice** durante la fase di GUI testing ?

RISPOSTA

Il livello di GUICC influenza in modo significativo la percentuale di codice eseguito .

DOMANDA 3

DOMANDA 3

I GUICC influenzano la capacità di scovare errori durante la fase di GUI testing ?

RUNTIME ERRORS

No	C-Lv	App name	Error type
1	C-Lv5	<i>Mileage</i>	Fatal error (Fatal signal 11 <F/libc>)
2		<i>Sanity</i>	Fatal exception (RuntimeException <E/AndroidRuntime>)
3		<i>Sanity</i>	Fatal exception (RuntimeException <E/AndroidRuntime>)
4	C-Lv4	<i>Mileage</i>	Fatal error (Fatal signal 11 <F/libc>)

- 4 tipi di errori fondamentali ;
- tutti rilevati con grafi "fine-grained" , ossia con criteri di comparazione a basso livello (C-Lv 4/5) .

DOMANDA 3

DOMANDA 3

I GUICC influenzano la capacità di scovare errori durante la fase di GUI testing ?

RISPOSTA

Oltre a influenzare la quantità di codice eseguito , il livello di GUICC influenza anche molto la capacità di "error detection" del test.

CONCLUSIONE

La maggior parte delle tecniche di *model-based Android GUI testing* esistenti non ha possibilità di utilizzo pratico, poichè vengono trascurati alcuni comportamenti dinamici frequenti nelle applicazioni moderne.

Questo è il primo studio che si è concentrato sulla *valutazione della generazione del modello* più che sulla applicazione del modello stesso su un qualche framework.

Da questa valutazione è emerso che i GUI Comparison Criteria sono molto vantaggiosi : si può scegliere il grado di dettaglio di un test (modificando il **Max C-Lv**) ed è possibile coprire un'alta percentuale di linee di codice eseguite , rilevando una altrettanto alta percentuale di **runtime errors**.

Manca un metodo che selezioni , per una data app , quali Comparison Criteria siano più adeguati per ottenere una più vasta copertura di codice e di errori. Tale metodo sarà oggetto di ricerca futura.