Reverse Engineering Models of Concurrent Communicating Systems From Event Logs

Sébastien Salva (sebastien.salva@uca.fr) LIMOS, Clermont Auvergne University, France



Sixteenth International Conference on Software Engineering Advances, oct. 2021

Who am I?

Public void setUp(){
Identity id=new Identity("salva");}

```
Public void testid (){
assertEquals(id.surname, ''sébastien'');
assertEquals(id.name, ''salva'');
assertEquals(id.labo, ''LIMOS'');
assertEquals(id.univ ''University Clermont Auvergne, France'');
```

Introduction



Introduction

Model Learning :

• Generation of behavioral models from a black box application (by retro-engineering).

Use of Models? : for documentation, analysis, auto generation of some test cases, etc.

Several limitations:

- Need of Working app. or execution traces
- Difficult to extract accurate conversations (a.k.a. sessions) when applications are made up of concurrent components
- May build spaghetti models (unredable and large models)

Introduction

Assumptions :

- System under learning = Concurrent Communicating Systems made up of components
- No knowledge of the components
- correlation mechanisms employed to propagate context IDs and keep track of the process contexts among components
- But, we don't know them

Proposal:

- Passive model learning approach and tool to recover Input Output Labelled Transition Systems (IOLTSs) from event logs.
- Algorithm to automatically retrieve conversations along with correlations from event logs

Paper presentation

- 1. Overview
- 2. Conversation Extraction from event logs
- 3. Model Generation
- 4. (preliminary) Evaluation
- 5. Conclusion





How to find Correlation keys ?

- Using a brute- force search: time consumming, not effiscient -> No
- Our proposal: algorithm is based upon a formalisation of the notion of correlation patterns and guided towards the most relevant conversation sets by evaluating conversation quality.

Correlation patterns :

- Key based correlation:
 - /login(from:="cl", to:="ShopS", id:="tocken", account:="l") ok(from:="ShopS", to:="cl", id:="tocken" trans:="t1")

Chained correlation:

/order(from:="cl", to:="ShopS", trans:="t1", item:="a")

Function based correlation ans Time-based correlation

/login(t:= 1) ok(t:=1) /order(t:=1) with t= floor(time(event)/T, T=5s

From Correlation patterns & Conversation quality

- From patterns, extraction of conversation invariants = properties on conversations and correlation key sets that must hold
- Definition of 4 Conversation Quality metrics = metrics between 0 and 1 to express properties of the conversation structures
 - Ex: request followed by response ?

Conversation quality metric example :

$$0 < m_1(\sigma) = \frac{|ReqwResp(\sigma)| + 1}{|Req(\sigma)| + 1} \le 1$$
$$0 < m_2(\sigma) = \frac{|RespwReq(\sigma)| + 1}{|Resp(\sigma)| + 1} \le 1$$

m1 evaluates the ratio of requests associated to some responses

m2 measures the ratio of responses following a prior request

(1)

(2)

- Algorithm overview:
- 1. Coverage of the successive events of an event log and search for potential correlation keys
- -> several candidates
- 2. Computation of invariants and quality metrics
- 3. If invariants do not hold or low quality remove candidate
- 4. Go back to 1 for every living candidate

LOG EXAMPLE

/login(from:="cl", to:="ShopS", id:="tocken", account:="l") ok(from:="ShopS", to:="cl", id:="tockén" trans:="t1") /order(from:="cl", to:="ShopS", trans:="t1",item:="a") /stock(from:="ShopS", to:="StockS", trans:="t1", item:="a")



requests followed by responses Good quality

/login alone ? Response followed by request ? Lower quality ^{13/22}

ICSEA 2021 sebastien.salva@uca.fr

LOG EXAMPLE

/login(from:="cl", to:="ShopS", id:="tocken", account:="1") ok(from:="ShopS", to:="cl", id:="tocken" trans:="t1") /order(from:="cl", to:="ShopS", trans:="t1",item:="a") /stock(from:="ShopS", to:="StockS", trans:="t1", item:="a") ok(from:="StockS", to:="ShopS", trans:="t1", item:="a") ok(from:="ShopS", to:="cl", trans:="t1", content:="stock") /supply(from:="ShopS", to:="WS", trans:="t1", key:="k1",item:="a") ok(from:="WS", to:="ShopS", trans:="t1", kev:="k1")/supplyWS(from:="WS", to:="WS1", key:="k1", key2:="k2",item:="a") /supplyWS(from:="WS", to:="WS2", key:="k1", key2:="k3",item:="a") /supplyWS(from:="WS", to:="WS3", key:="k1", key2:="k4",item:="a") ok(from:=WS1,to:="WS", key:="k1", key2:="k2") ok(from:=WS2,to:="WS", key:="k1", key2:="k3") Unavailable(from:="WS3", to:="WS", key:="k1", kev2:="k4")/login(from:="cl", to:="ShopS", id:="tocken2", account:="12") ok(from:="ShopS", to:=cl, id:="tocken2", trans:="t2") /order(from:="cl", to:="ShopS", trans:="t2",item:="b") ok(from:="ShopS", to:="cl", trans:="t2" content:="no stock")

Corr keys :={id,key,key2}

2 conversations :

Conv:={

/login() ok() /order() /stock() ok() ok()/supply() ok()
/supplyWS()/supplyWS() /supplyWS() ok() ok() unavailable(),

/login() ok() /order() ok()

Trace partionning

Detection of components with analysis of emitters and receivers in events

Generation of as many trace sets as components found

Algorithm to segment the conversations into sub-sequences, each capturing the behaviours of one component only

Events Example

```
/login(from:="cl", to:="ShopS )
/order(from:="cl", to:="ShopS")
/stock(from:="ShopS",
to:="StockS") ok(from:="StockS",
to:="ShopS") ok(from:="ShopS",
to:="cl").....
```

T(ShopS)={ ?/login() !ok() ?/order() !/stock() ?ok() !ok() !/supply() ?ok(), ?/login() !ok() ?/order() !ok() }

T(<mark>StockS</mark>)={ ?/stock() !ok() }

T(WS)={ ?/supply() !ok() !/supplyWS() !/supplyWS() !/supplyWS() ?ok() ?ok() ?unavailable() }

T(WS1)=T(WS2)={ ?/supplyWS() !ok() }

T(WS3)={ ?/supplyWS() !unavailable() }

IOLTS Generation

Every trace set lifted to the level of IOLTS

1 trace set -> 1 IOLTS obtained by transforming the traces into IOLTS paths having the same initial state only

T(ShopS)={ ?/login() !ok() ?/order() !/stock() ?ok() !ok() !/supply() ?ok(), ?/login() !ok() ?/order() !ok() }



IOLTS Generalisation

IOLTS generalisation by merging their equivalent states. State merging performed by means of the k-Tail algorithm assembles the states sharing the same k-future, i.e., the same event sequences having the maximum length k



Exemple: Final Results

LOG EXAMPLE

/login(from:="cl", to:="ShopS", id:="tocken", account:="1") ok(from:="ShopS", to:="cl", id:="tocken" trans:="t1") /order(from:="cl", to:="ShopS", trans:="t1",item:="a") /stock(from:="ShopS", to:="StockS", trans:="t1", item:="a") ok(from:="StockS", to:="ShopS", trans:="t1", item:="a") ok(from:="ShopS", to:="cl", trans:="t1", content:="stock") /supply(from:="ShopS", to:="WS", trans:="t1", key:="kl",item:="a") ok(from:="WS", to:="ShopS", trans:="t1", key:="k1")/supplyWS(from:="WS", to:="WS1", key:="k1", key2:="k2",item:="a") /supplyWS(from:="WS", to:="WS2", key:="k1", key2:="k3",item:="a") /supplyWS(from:="WS", to:="WS3", key:="k1", key2:="k4",item:="a") ok(from:=WS1,to:="WS", key:="k1", key2:="k2") ok(from:=WS2,to:="WS", key:="k1", key2:="k3") Unavailable(from:="WS3", to:="WS", key:="k1", key2:="k4")/login(from:="cl", to:="ShopS", id:="tocken2", account:="12") ok(from:="ShopS", to:=cl, id:="tocken2", trans:="t2") /order(from:="cl", to:="ShopS", trans:="t2",item:="b") ok(from:="ShopS", to:="cl", trans:="t2" content:="no stock")







Conducted on 6 IoT systems integrating varied devices and gateways communicating over HTTP and UDP



- RQ1: Can the approach extract relevant conversation sets?
- RQ2: What is the performance of our algorithm?

RQ1: Can The Approach Extract Relevant Conversation Sets?

Event logs ~ 2200 events

Manually analysed the event logs S1 to S6-> correct correlation key sets

Computed Precision and Recall on the generated conversations

8	Correlation Key	Correlation Key
	Set Recall	Set Precision
<i>S</i> 1	100%	81%
<i>S</i> 2	100%	76%
<i>S</i> 3	100%	80%
<i>S</i> 4	100%	100%
<i>S</i> 5	100%	100%
<i>S</i> 6	100%	90%

Results: provides good recall and precision, but sometimes returns several results, choice can be conducted with quality metrics ordering

RQ2: what is the performance of our algorithm?

1. We took the 20 first conversations of S1 and augmented them using 40 to 10000 events;



The curve follows a quadratic curve and reveals that our approach performs well in practice.

- RQ2: what is the performance of our algorithm?
- 2. We measured execution times with regard to the number of conversations in the event logs from 10 to 200 conversations of 2 events.



cubic polynomial curve, shows that execution times quicker increase with regard to the number of conversations.

Conclusion

- Design of a model learning approach specialised into the recovery of formal models from event logs generated by communicating systems made up of concurrent components
- The generated IOLTS can be later used as documentation or for automatics analyses (security testing, etc.)

Limitations:

• Need of a good balance between model size, readability and precision. For instance, the generated IOLTSs may be very large on account of similar event sequences having different parameter values.

Thanks

• Questions ?