

Mining Past-Time Temporal Rules From Execution Traces

David Lo^{1,2} Siau-Cheng Khoo² Chao Liu³

¹Singapore Management University

²National University of Singapore

³Microsoft Research, Redmond

Issue on Software Specifications

- o Documented specifications are often lacking, poor, outdated and incomplete

Hard deadlines & `short-time-to-market`

Productivity == LOC or completed project

High turn-over rate of IT professionals

Difficulties & programmer's reluctance in writing formal specs

(Ammons et al., POPL'02, Yang et al., ICSE'06)

The Specification Problem

- o **Contributes to high software costs**

 - Program comprehension = 50% of maintenance cost

 - High maintenance cost = 90% total cost

 - (Erlikh, 2000; Cimitile & Canfora, 2001)

 - US GDP software component = 214.4 billion USD.

- o **Causes challenges in ensuring correctness of systems**

 - Difficulty in verifying correctness of systems

 - US National Institute of Standards and Technology

 - 59.5 Billions annual lost due to bugs

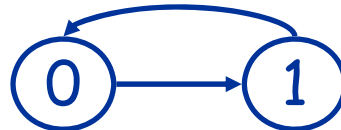
Specification Mining (SM)

A process to discover protocols that a code exhibit, often through an analysis of its execution traces (ABL02 [POPL])

Benefits:

Aid Program Comprehension and Maintenance
Aid Program Verification

Automaton-based SM



RR01 [ICSE], CW98 [TOSEM]
ABL02 [POPL], AMBL03 [PLDI],
WML02 [ISSTA], AXPX07 [FSE]
MP05 [ICEECS], LK06 [FSE]

Rule-based SM

<Lock> -> <Unlock>

YEBBD06 [ICSE]
LKLO8 [DASFAA, JSME]

Only **future-time** temporal rules are mined

Past-Time Temporal Rules

Whenever a series of events *pre* occurs, previously, another series of events *post* happened before, denoted as: $pre \rightarrow_p post$

Among **most-widely used** temporal logic expressions (Dwyer, ICSE'99)

Why Important ?

- Past-time temporal exp. \rightarrow complex future-time temporal exp. (Laroussinie et al., TCS'95, LICS'02)
- Not minable by existing algorithms mining future time rules (Yang et al. ICSE'06, Lo et al. JSME'08]
- Many interesting properties are more intuitively expressed in past-time
- Many interesting properties are non-symmetric

- o Whenever a **file is used** (read or written), it needs to be **opened** before.

file_used \rightarrow_p file_open

- o Whenever **SSL_read** is performed, **SSL_init** needs to be invoked before.

ssl_read \rightarrow_p ssl_init

- o Whenever a valid client **request** a non-sharable resource and the resource is **not granted**, previously the resource had been **allocated** to another client that **requested** it.

request, not_granted \rightarrow_p request, grant

- o Whenever money is **dispensed** from an ATM, previously, **card** was inserted, **pin** was entered, user was **authenticated** and account **balance suffices**.

dispense \rightarrow_p card, pin, authenticate, balance_suffice

Outline

- o **Motivation and Introduction**
- o **Concepts**
 - Past-Time LTL, Statistical Significance
 - Soundness and Completeness
- o **Mining Past-Time Rules**
 - Mining Strategy, Pruning Properties
 - Removal of Redundant Rules
 - Mining Framework
- o **Preliminary Experiments**
- o **Discussion**
- o **Related Work**
- o **Conclusion & Future Work**

Concepts

Past-Time Linear Temporal Logics (PLTL)

- o Linear Temporal Logic (LTL)
 - Logic that works on program paths
 - A path corresponds to an execution trace
- o Past-Time Linear Temporal Logic (PLTL)
 - Add LTL with past time operators
 - More succinct than LTL
- o Temporal operators under consideration
 - ' G ' - Globally
 - ' F ' - Once in the future
 - ' X ' - Next (immediate)
 - ' F^{-1} ' - Once in the past
 - ' X^{-1} ' - Previous (immediate)

PLTL- Examples

o $X^{-1}F^{-1}$ (file_open)

Meaning: At a time in the past file is opened

o $G(\text{file_read} \rightarrow X^{-1} F^{-1} (\text{file_open}))$

Meaning: Globally whenever file is read, at a time in the past file is opened

o $G((\text{account_deducted} \wedge XF (\text{money_dispensed})) \rightarrow (X^{-1}F^{-1} (\text{balance_suffice} \wedge (X^{-1}F^{-1} (\text{cash_requested} \wedge (X^{-1}F^{-1} (\text{correct_pin} \wedge (X^{-1}F^{-1} (\text{insert_debit_card}))))))))))$

Meaning: Globally whenever one's bank account is deducted and money is dispensed (from an ATM), previously user inserted debit card, entered correct pin, requested for cash and account balance suffices.

Notations and Scope of Mined Rules

- o Denote a past-time rule as $pre \rightarrow_p post$
- o Sample mappings btw. rule representations and PLTL expressions

Notation	LTL Notation
$a \hookrightarrow_p b$	$G(a \rightarrow X^{-1}F^{-1}b)$
$\langle a, b \rangle \hookrightarrow_p c$	$G((a \wedge XFb) \rightarrow (X^{-1}F^{-1}c))$
$a \hookrightarrow_p \langle b, c \rangle$	$G(a \rightarrow X^{-1}F^{-1}(c \wedge X^{-1}F^{-1}b))$
$\langle a, b \rangle \hookrightarrow_p \langle c, d \rangle$	$G((a \wedge XFb) \rightarrow (X^{-1}F^{-1}(d \wedge X^{-1}F^{-1}c)))$

- o Scope of minable temporal expressions

$rules := G(pre \rightarrow post)$
$pre := (event) (event \wedge XF(pre))$
$post := (event) (event \wedge X^{-1}F^{-1}(post))$

Statistical Significance Metrics

o Distinguish Significant Rules via Statistical Notions

- Support: The number of traces supporting the premise *pre*
- Confidence: The likelihood of the premise *pre* being preceded by the consequent *post*

Sample Traces

Identifier	Trace/Sequence
S1	(c, b, a, e, b, a)
S2	(c, b, e, a, e, b, c, a)
S3	(d, a)

Rule: $\langle b, a \rangle \rightarrow_p \langle c \rangle$

Support: 2

Corres. to S1 and S2

Confidence: 100%

All occurrences of $\langle b, a \rangle$ is preceded by $\langle c \rangle$

Rule: $\langle b, a \rangle \rightarrow_p \langle e \rangle$

Support: 2

Confidence: 50%

Soundness and Completeness

- o **Ensure Soundness and Completeness**
 - With respect to input traces and specified thresholds
- o **Sound**
 - All mined rules are statistically significant
- o **Complete**
 - All statistically significant rules are mined or represented
- o **Commonly used in data/pattern mining**

Mining Past Time Temporal Rules

High-Level Mining Strategy

- o Mining Option 1: Check for all 2-event rules ($n \times n$ of them) for statistical significance.
 - Not scalable for rules of arbitrary lengths.
- o Our Mining Strategy: Consider mining as a search-space traversal for significant rules
 - Explore the search space depth-first
 - Identify significant rules
- o Employ pruning strategies to throw away search space containing insignificant rules
- o Detect search spaces containing redundant rules early during the mining process

Anti-Monotone Pruning Strategies

[*Apriori – Support*]

$$Rx = p \rightarrow_p c; Ry = q \rightarrow_p c$$

$$p \sqsubseteq q$$

$$\frac{sup(Rx) < min_sup}{sup(Ry) < min_sup}$$

$$sup(Ry) < min_sup$$

Ry is not significant

[*Apriori – Confidence*]

$$Rx = p \rightarrow_p c; Ry = p \rightarrow_p d$$

$$c \sqsubseteq d$$

$$\frac{conf(Rx) < min_conf}{conf(Ry) < min_conf}$$

$$conf(Ry) < min_conf$$

Ry is not significant

Rx: $a \rightarrow_p z$; $sup(Rx) < min_sup$

Ry_s $\left. \begin{array}{l} a, b \rightarrow_p z \\ a, b, c \rightarrow_p z \\ a, c \rightarrow_p z \\ a, b, d \rightarrow_p z \\ \dots \end{array} \right\}$ **Non-significant**

Rx: $a \rightarrow_p z$; $conf(Rx) < min_conf$

Ry_s $\left. \begin{array}{l} a \rightarrow_p z, b \\ a \rightarrow_p z, b, c \\ a \rightarrow_p z, c \\ a \rightarrow_p z, b, d \\ \dots \end{array} \right\}$ **Non-significant**

Detecting Redundant Rules

$$Rx = p \rightarrow_{\rho} c; Ry = q \rightarrow_{\rho} d$$

$$p \dashv\dashv c \sqsubseteq q \dashv\dashv d$$

$$\text{sup}(Rx) = \text{sup}(Ry)$$

$$\text{conf}(Rx) = \text{conf}(Ry)$$

Rx is redundant

Redundant rules are identified and removed early during mining process.

$$Rx: a \rightarrow_{\rho} b, c, d$$

Ry_s

$$a \rightarrow_{\rho} b$$

$$a \rightarrow_{\rho} c$$

$$a \rightarrow_{\rho} b, c$$

$$a \rightarrow_{\rho} b, d$$

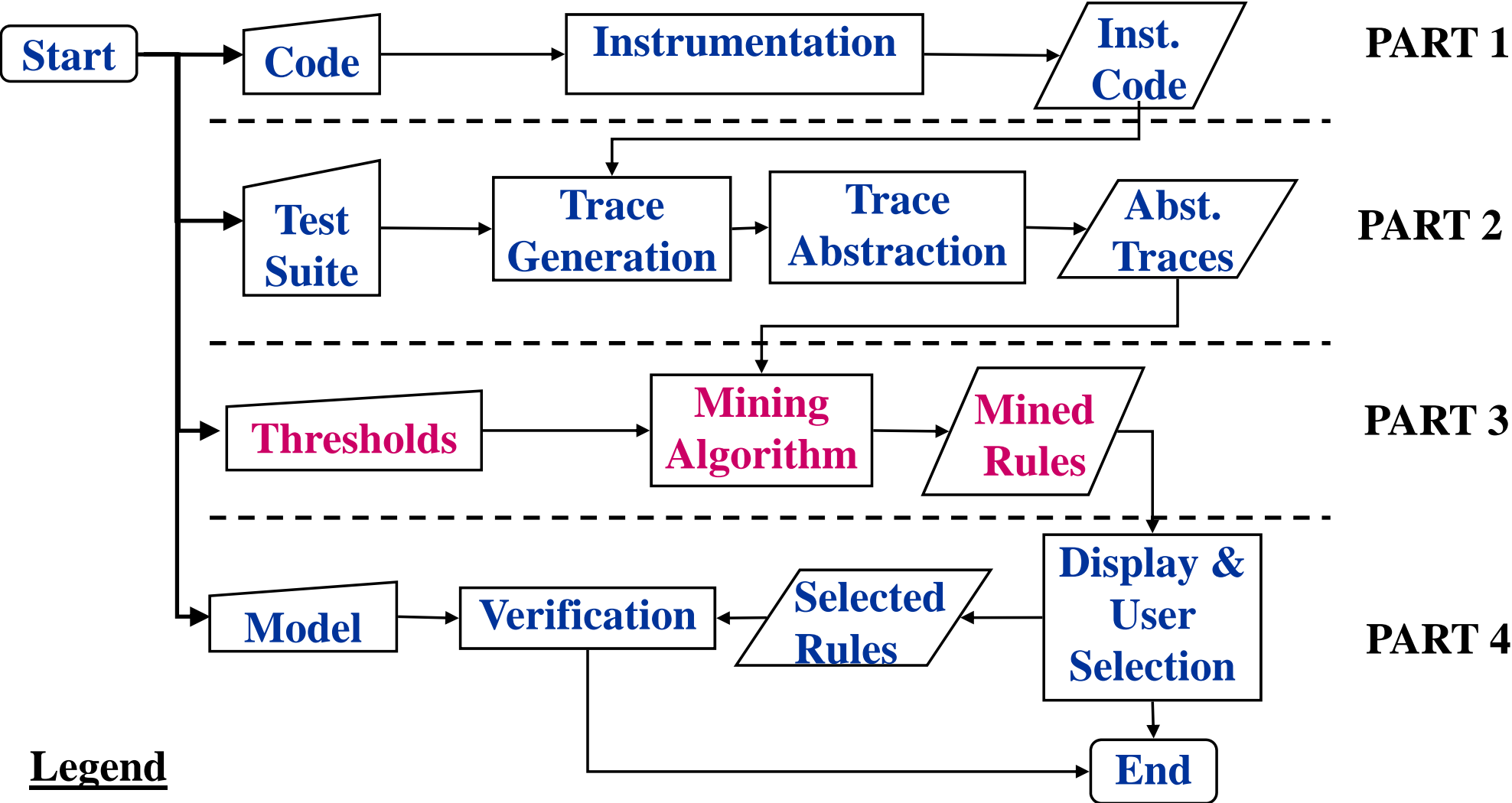
....

Redundant
iff
sup and conf are
the same

Algorithm Steps

- o **Step 1:** Generate a **pruned** set of **significant pre-conditions** satisfying the minimum **support** threshold.
- o **Step 2:** For each pre-condition, find occurrences of pre in the trace database.
- o **Step 3:** For each pre-condition, generate a **pruned** set of **significant post-conditions** satisfying the minimum **confidence** threshold.
- o **Step 4:** Remove remaining rules that are redundant. Note that many/most redundant rules have been removed at step 1 and 3.

Mining Framework



Legend

User Input
 Process
 Intermediate Result

Preliminary Experiments

Experiment Setups - JBoss Application Server

- o **JBoss Application Server (JBoss AS)**
 - One of the most widely used J2EE application server
 - Analyze the transaction and security component
- o **Program Instrumentation & Trace Generation**
 - Instrument the application using JBoss-AOP
 - Run regression tests from JBoss AS distribution
- o **Transaction component**
 - 2551 events, 64 unique events
 - min_sup: 25, min_conf: 90%
 - Mining time: 30 seconds , Mined non-redundant rules: 36
- o **Security component**
 - 4115 events, 60 unique events
 - min_sup: 15, min_conf: 90%
 - Mining time: 2.5 seconds, Mined non-redundant rules: 4

A Rule from JBoss Transaction

Premise	→ P	Consequent
<code>TransactionImpl.isDone()</code>		<code>TxManagerLocator.getInstance()</code> <code>TxManagerLocator.locate()</code> <code>TxManagerLocator.tryJNDI()</code> <code>TxManagerLocator.usePrivateAPI()</code> <code>TxManager.getInstance()</code> <code>TxManager.begin()</code> <code>XidFactory.newXid()</code> <code>XidFactory.getNextId()</code> <code>XidImpl.getTrulyGlobalId()</code> <code>TransImpl.assocCurrentThread()</code> ... 5 events ... <code>TxManager.getTransaction()</code>

Whenever a **transaction** is checked for **completion** (premise), previously **transaction manager** is located (ev 1-4 consequent), **transaction manager & impl** are initialized (ev 5-6,10-12), **ids** are acquired (ev 7-9,13-15) and **transaction object** is obtained from the **manager** (ev 16).

A Rule from JBoss Security

Premise	→ P	Consequent
<code>SimplePrincipal.toString()</code>		<code>XLoginConfImpl.getConfEntry()</code>
<code>SecAssoc.getPrincipal()</code>		<code>PolicyConfig.get()</code>
<code>SecAssoc.getCredential()</code>		<code>XLoginConfImpl\$1.run()</code>
<code>SecAssoc.getPrincipal()</code>		<code>AuthenInfo.copyAppConfEntry()</code>
<code>SecAssoc.getCredential()</code>		<code>AuthenInfo.getName()</code>
		<code>ClientLoginModule.initialize()</code>
		<code>ClientLoginModule.login()</code>
		<code>ClientLoginModule.commit()</code>
		<code>SecAssocActs.setPrincipalInfo()</code>
		<code>SetPrincipalInfoAction.run()</code>
		<code>SecAssocActs.pushSubjectContext()</code>
		<code>SubjectThreadLocalStack.push()</code>

Whenever **principal and credential info is required** (the premise), previously **config. info is checked** to determine the auth. service availability (ev 1-5), **actual authentication events are invoked** (ev 6-8) and **principal info is bound to the subject** (ev 9-12)

Discussions

o Setting min-sup/conf threshold

- Appropriate values depend on application
- Mining as an iterative process

o Sound and Complete

- With respect to trace and specified thresholds
- If trace is not complete or buggy so does the results
- Confidence provide a measure of tolerance to buggy traces

o Scalability

- Algorithm works better with many shorter traces than one very long trace
- It's better to split a trace to sub-traces
 - Focus on immediate inter-component interaction (Mariani et al., ICSE'08)
 - Trace abstraction (Ammons et al., POPL'02)

Related Work

o Daikon

- Complement Daikon by mining temporal constraints

o Mining Automata

- Many work: ABL02, RR01, MP05, AXPX07, LK06, ...
- Diff: Focus on statistically significant property rather than overall behavior

o Mining Future-Time Temporal Rules

- Many work: YEBBD06, LKL08, ...
- Diff: Mining past-time temporal rules

o Mining Sequence Diagram: BLL06, LMK07, LM08, ...

o Mining from Code: RGJ07, WN05, ...

o Data Mining: S99, AS94, YHA03, WH04, LKL07, ...

Conclusion

- o Propose a new approach to mine **past-time temporal rules** using **dynamic analysis**, not minable by existing tools:

Whenever a series of events *pre* occurs, **previously**, another series of events *post* happened before, denoted as: $pre \rightarrow_p post$

- o Address the problem of runtime costs by employing **smart pruning strategies**.
 - Throw away **insignificant rules en-masse**
 - Throw away **redundant rules en-masse**
- o **Preliminary experiments** on traces of JBoss AS show utility of the technique to discover **program behavioral rules/specifications**

Future Work

o User Guided Mining

- Let user provide more information to the mining process aside from the significance thresholds
- Mining Scenario-Based Triggers and Effects
 - (- ASE'08 - to-appear) - Mining Sequence Diagram

o Mining more complex LTL expressions

- Incorporating both future and past-time temporal rules

o Improving the scalability of the technique

- Abstraction technique
- Pruning strategies

o Experimenting with more case studies

Thank you

Comments ? Questions ? Advices ?