

# Mining Past-Time Temporal Rules from Execution Traces

David Lo<sup>\*</sup>  
School of Information Systems  
Singapore Management  
University  
davidlo@smu.edu.sg

Siau-Cheng Khoo  
School of Computing  
National University of  
Singapore  
khoosc@comp.nus.edu.sg

Chao Liu  
Microsoft Research, Redmond  
chaoliu@microsoft.com

## ABSTRACT

Specification mining is a process of extracting specifications, often from program execution traces. These specifications can in turn be used to aid program understanding, monitoring and verification. There are a number of dynamic-analysis-based specification mining tools in the literature, however none so far extract past time temporal expressions in the form of rules stating: “whenever a series of events occurs, previously another series of events has happened”. Rules of this format are commonly found in practice and useful for various purposes. Most rule-based specification mining tools only mine future-time temporal expression. Many past-time temporal rules like “whenever a resource is used, it was allocated before” are asymmetric as the other direction does not holds. Hence, there is a need to mine past-time temporal rules.

In this paper, we describe an approach to mine significant rules of the above format occurring above a certain statistical thresholds from program execution traces. The approach start from a set of traces, each being a sequence of events (*i.e.*, method invocations) and resulting in a set of significant rules obeying minimum thresholds of support and confidence. A rule compaction mechanism is employed to reduce the number of reported rules significantly. Experiments on traces of JBoss Application Server shows the utility of our approach in inferring interesting past-time temporal rules.

**Categories and Subject Descriptors:** D.2.1 [Software Engineering]: Requirements/Specifications – Tools; D.2.7 [Software Engineering]: Distribution, Maintenance and Enhancement – Restructuring, reverse engineering and reengineering

**General Terms:** Algorithms, Design, Experimentation

**Keywords:** Specification Mining, Dynamic Analysis, Past-Time Temporal Rules, Data Mining

## 1. INTRODUCTION

Different from many engineering products that rarely change, software changes often throughout its lifespan. This phenomenon has been well studied under the umbrella notion of software evolution.

<sup>\*</sup>The work was done while the author was with School of Computing, National University of Singapore.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

WODA – Workshop on Dynamic Analysis, July 21, 2008  
Copyright 2008 ACM 978-1-60558-054-8/08/07 ...\$5.00.

Software maintenance effort deals with the management of such changes, ensuring that the software remains correct while additional features are incorporated. Maintenance cost can contribute up to 90% of software development cost [10]. *Reducing maintenance cost* and ensuring a program *remains correct during evolution* are certainly two worthwhile goals to pursue.

A substantial portion of maintenance cost is due to the difficulty in understanding an existing code base. Studies show that program comprehension can contribute up to 50% of the maintenance cost [13, 32]. A challenge to software comprehension is the maintenance of an accurate and updated specification as program changes. As a study shows, documented specifications often remain unchanged during program evolution [7]. One contributing factor is the short time-to-market requirement of software products. Multiple cycles of software evolution can potentially render the outdated specification invalid or even misleading.

To ensure correctness of a software system, model checking [6] has been proposed. It accepts a model and a set of formal properties to check. Unfortunately, difficulty in formulating a set of formal properties has been a barrier to its wide-spread adoption [3]. Adding software evolution to the equation, the verification process is further strained. First, ensuring correctness of software as changes are made is not a trivial task: a change in one part of a code, might induce unwanted effects resulting in bugs in other parts of the code. Furthermore, as a system changes and features are added, there is a constant need to add new properties or modify outdated properties to render automated verification techniques effective in detecting bugs and ensuring the correctness of the system.

Addressing the above problems, there is a need for techniques to automatically reverse engineer or mine formal specifications from program. Recently, there has been a surge in software engineering research to adopt machine learning and statistical approaches to address these problems. One active area is specification discovery [3, 29, 1, 23], where software specification is reverse-engineered from program traces. Employing these techniques ensures specifications remain updated; also it provides a set of properties to verify via formal verification tools like model checking.

Most specification miners extract specifications in the form of automata [23, 3, 29, 1] or temporal rules [34, 26]. Usually a mined automata express the whole behavior of a system under analysis. Mined rules express strongly-observed constraints each expressing a property which holds with certain statistical significance.

Rules mined in [34, 25, 26] express future-time temporal expressions. Yang *et al.* mine two event rules of the form: “Whenever an event occur, eventually another event occur in the future” [34]. Lo *et al.* mine temporal rules of arbitrary length of the form: “Whenever a series of event occur, eventually another series of event occur in the future” [25, 26]. In this work, we extend the above work by

mining past-time temporal expressions of this format:

“Whenever a series of events *pre* occurs, previously,  
another series of events *post* has happened”

The above rule is denoted by  $pre \hookrightarrow_P post$ , where *pre* and *post* corresponds respectively to the premise (pre-condition) and consequent (post-condition) of the rule. These set of rules can be expressed in past-time temporal logic (Linear Temporal Logic (LTL) + past time operators [19]) and belong to two of the most used families of temporal logic expressions for verification (*i.e.*, precedence and chain-precedence) according to a survey by Dwyer *et al.* [8]. Some example specifications of the above format are as follows:

1. Whenever a file is used, it was opened before.
2. Whenever a socket is written, it was initialized before.
3. Whenever `SSL_read` is performed, `SSL_init` was invoked before.
4. Whenever a client request a resource and the resource is not granted, previously, the resource had been allocated to another client that requested it.
5. Whenever money is dispensed in an Automated Teller Machine (ATM), previously, card was inserted, pin was entered, user was authenticated and account balance was checked.

It has been shown that past-time LTL can express temporal properties more succinctly than (pure) future-time LTL [18, 22]. Simple past-time LTL can correspond to more complicated equivalent future-time LTL, many of which are not minable by existing techniques mining for future-time LTL rules from traces [34, 26]. The subset of the past-time LTL mined by the approach presented in this paper is not minable by previous approach in [34, 26]. Also, for some past-time temporal rules, the corresponding future-time rules do not hold. Consider the first and second example above; it is not the case that when a file is opened or a socket is initialized, the file need to be eventually used or the socket need to be eventually written, respectively. Thus, our work is not meant to replace, but rather to complement mining future-time temporal rules.

In static inference of specification, Ramanathan *et al.* [31] mine specifications from program source code of the form: “Whenever an event occurs, previously, another series of events has happened”. Different from Ramanathan *et al.* we analyze program traces and we need to address the issue of repeated behaviors within program traces (due to loops and recursions). Ramanathan *et al.* uses an off-the-shelf data mining algorithm [2] which *ignores repetitions within a sequence*. Static analysis has a different set of issues related to difficulty in analyzing pointers & references [5] and the number of infeasible paths [4]. Also, our target specification format is more general in that it captures pre-conditions of *multiple* event length, hence enabling user to mine for more complex temporal properties. Also, as will be described later, we present a method that compacts significant rules by ‘early’ pruning of redundant rules resulting in a potentially combinatorial speed-up and reduction in the set of mined rules. In [31], all significant rules are first generated before redundant ones are removed. The large number of intermediary rules (exponential to the length of the longest rule) might make the algorithm not scalable enough to mine for rules of *long* length.

Our mining algorithm models mining as a search space exploration process. The input is a set of sequences of events, where an

event corresponds to an interesting method invocation to be analyzed. The output is a set of significant rules that obeys the minimum thresholds of support and confidence – which are commonly used statistics in data mining [14]. We define the support of a rule to be the number of traces where the rule’s premise is observed. We define the confidence of a rule to be the likelihood that the rule’s premise is preceded by the consequent. The algorithm first constructs rule of short length, and utilizes several properties to throw away sub-search space not yet traversed if a short-length rule is not significant.

In addition, we observe that some rules are redundant. To address this, we employ additional pruning strategies to throw away these redundant rules. We kept the more comprehensive longer rules that capture more information and hence subsume the shorter ones.

We guarantee that all mined rules are significant and non-redundant. Also, all significant and non-redundant rules are mined. In data mining, an algorithm meeting the first and second criteria above is referred to as being correct (sound) and complete respectively (*c.f.*, [14, 20]). In this paper, we refer to the above criteria as *statistical soundness and completeness*.

To demonstrate our ideas, in this paper, we experimented with traces from components of JBoss Application Server [16]. The experiments show the utility of our approach in mining specifications of an industrial program.

The paper is organized as follows. In Section 2, we discuss the semantics of mined rules. Section 3 describes our mining algorithm. Section 4 describes the experiments performed. Section 5 describes related work. Section 6 discusses some future work and we finally conclude in Section 7.

## 2. CONCEPTS AND DEFINITIONS

This section introduces preliminaries on past-time LTL and formalize the scope of rules minable by our approach. Also, notations used in this paper are described.

**Past-Time Linear-time Temporal Logic** Our mined rules can be expressed in past-time Linear Temporal Logic (LTL) [18, 19, 22]. Past-time temporal logic is an extension of (future-time) LTL [30, 15]. LTL is a logic that works on possible program paths. A possible program path corresponds to a program trace. A path can be considered as a series of events, where an event is a method invocation. For example, (file\_open, file\_read, file\_write, file\_close), is a 4-event path.

There are a number of LTL operators, among which we are only interested in the operators ‘G’, ‘F’ and ‘X’ and some of their past-time counterparts ‘F<sup>-1</sup>’ and ‘X<sup>-1</sup>’. The operator ‘G’ specifies that *globally* at every point in time a certain property holds. The operator ‘F’ specifies that a property holds at that point in time or *at a time in the future*. The operator ‘X’ specifies that a property holds at the *next* event.

The operator ‘F<sup>-1</sup>’ specifies that a property holds at that point in time or *at a time in the past*. The operator ‘X<sup>-1</sup>’ specifies that a property holds at the *previous* event.

Let us consider three examples listed in Table 1.

Our mined rules state whenever a series of premise events occurs it was preceded by another series of consequent events. A mined rule denoted as  $pre \hookrightarrow_P post$ , can be mapped to its corresponding LTL expression. Examples of such correspondences are shown in Table 2.

Mapping to common English language expressions and for uniformity purpose, in both the premise and consequent of the rule the time goes forward to the future (*e.g.*, *a* is followed by *b*, is preceded by, *c* is followed by *d*). In the corresponding past-time LTL

$X^{-1}F^{-1}(file\_open)$
Meaning: <i>At a time in the past</i> file is opened
$G(file\_read \rightarrow X^{-1}F^{-1}(file\_open))$
Meaning: <i>Globally</i> whenever file is read, <i>at a time in the past</i> file is opened
$G((account\_deducted \wedge XF(money\_dispensed)) \rightarrow (X^{-1}F^{-1}(balance\_suffice \wedge (X^{-1}F^{-1}(cash\_requested \wedge (X^{-1}F^{-1}(correct\_pin \wedge (X^{-1}F^{-1}(insert\_debit\_card))))))))))$
Meaning: <i>Globally</i> 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 to be dispensed and account balance was checked and it sufficed.

**Table 1: Past-time LTL Expressions and their Meanings**

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

**Table 2: Rules and their Past-time LTL Equivalences**

expression we need to reverse the order of c and d. Also note that although the operator ‘X’ might seem redundant, it is needed to specify rules such as  $\langle a \rangle \hookrightarrow_P \langle b, b \rangle$  where the ‘b’s refer to *different occurrences of ‘b’*. The set of LTL expressions minable by our mining framework is represented in the Backus-Naur Form (BNF) as follows<sup>1</sup>:

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

**Basic Notations** Let  $I$  be a set of distinct events considered in which an event corresponds to a behavior of interest, e.g. method call. Input to our mining framework is a set of traces. A trace corresponds to a sequence or an ordered list of events from  $I$ . For formality, we refer to this set of traces as a sequence database denoted by  $SeqDB$ . Each trace or sequence is denoted by  $\langle e_1, e_2, \dots, e_{end} \rangle$  where  $e_i \in I$ .

We define a pattern  $P$  to be a series of events. We use  $first(P)$  to denote the first event of  $P$ . A pattern  $P_1 \uparrow P_2$  denotes the concatenation of patterns  $P_1$  and  $P_2$ . A pattern  $P_1 \langle \langle e_1, e_2, \dots, e_n \rangle \rangle$  is considered a *subsequence* of another pattern  $P_2 \langle \langle f_1, f_2, \dots, f_m \rangle \rangle$  denoted as  $P_1 \sqsubseteq P_2$  if there exist integers  $1 \leq i_1 < i_2 < \dots < i_n \leq m$  such that  $e_1 = f_{i_1}, e_2 = f_{i_2}, \dots, e_n = f_{i_n}$ .

### 3. MINING PAST TIME TEMPORAL RULES

Each temporal rule of interest has the form  $P_1 \hookrightarrow_P P_2$ , where  $P_1$  and  $P_2$  are two series of events.  $P_1$  is referred to as the *premise* or *pre-condition* of the rule, while  $P_2$  is referred to as the *consequent* or *post-condition* of the rule. The rules correspond to temporal constraints expressible in past-time LTL notations. Some examples are shown in Table 2.

In this paper, since a trace is a series of events, where an event corresponds to a software behavior of interest, e.g., method call, we formalize a trace as a sequence and a set of input traces as a sequence database. We use the sample trace or sequence database in Table 3 as our running example to illustrate the concepts behind generation of temporal rules.

#### 3.1 Concepts & Definitions

Mined rules are formalized as past-time Linear Temporal Logic expressions with the format:  $G(\dots \rightarrow X^{-1}F^{-1}\dots)$ . The semantics of past-time LTL described in Section 2 will dictate the semantics

<sup>1</sup>post is in reversed order

Identifier	Trace/Sequence
S1	$\langle c, b, a, e, b, a \rangle$
S2	$\langle c, b, e, a, e, b, c, a \rangle$
S3	$\langle d, a \rangle$

**Table 3: Example Database – DBX**

of temporal rules described here. Noting the meaning of the temporal operators illustrated in Table 1, to be precise, a mined past-time temporal rule expresses:

“Whenever a series of events *occurs starting at a point in time (i.e. a temporal point)*, previously, another series of events has happened.”

From the above definition, to generate temporal rules, we need to “peek” at interesting temporal points and “see” what series of events are likely to occur *before*. We first formalize the notion of temporal points and the related notion of occurrences.

**DEFINITION 3.1 (Temporal Points).** Consider a sequence  $S$  of the form  $\langle a_1, a_2, \dots, a_{end} \rangle$ . All events in  $S$  are indexed by their positions in  $S$ , starting at 1 (e.g.,  $a_j$  is indexed by  $j$ ). These positions are called temporal points in  $S$ . For a temporal point  $j$  in  $S = \langle a_1, \dots, a_n \rangle$ , the suffix  $\langle a_j, \dots, a_n \rangle$  is called the  $j$ -suffix of  $S$ .

**DEFINITION 3.2 (Occurrences & Instances).** Given a pattern  $P$  and a sequence  $S$ , the occurrences of  $P$  in  $S$  are defined by a set of temporal points  $T$  in  $S$  such that for each  $j \in T$ , the  $j$ -suffix of  $S$  is a super-sequence of  $P$  and  $first(P)$  is indexed by  $j$ . The set of instances of pattern  $P$  in  $S$  is defined as the set of  $j$ -suffixes of  $S$ , for each  $j \in T$ .

**Example.** Consider a pattern  $P \langle b, a \rangle$  and the sequence  $S1$  in Table 3 (i.e.,  $\langle c, b, a, e, b, a \rangle$ ). The occurrences of  $P$  in  $S1$  form the set of temporal points  $\{2, 5\}$ , and the corresponding set of instances are  $\{\langle b, a \rangle, \langle b, a, e, b, a \rangle\}$ .

We define database projection operations to capture events occurring before specified temporal points. The following are two different types of projections and their associated support notions.

**DEFINITION 3.3 (Projected-past & Sup-past).** A database projected-past on a pattern  $p$  is defined as:

$SeqDB_P^{past} = \{(j, px) \mid \text{the } j^{th} \text{ sequence in } SeqDB \text{ is } s, \text{ where } s = px \uparrow sx, \text{ and } sx \text{ is the minimum suffix of } s \text{ containing } p\}$

Given a pattern  $P_X$ , we define  $\sup^{past}(P_X, SeqDB)$  to be the size of  $SeqDB_{P_X}^{past}$  (i.e., the number of sequences in  $SeqDB$  containing  $P_X$ ). Reference to the database is omitted, i.e., we write it

	Identifier	Trace/Sequence
(a)	S1	(1, (c, b, a, e))
	S2	(2, (c, b, e, a, e))
	Identifier	Trace/Sequence
(b)	S1 <sub>1</sub>	(1, (c, b, a, e))
	S1 <sub>2</sub>	(1, (c))
	S2 <sub>1</sub>	(2, (c, b, e, a, e))
	S2 <sub>2</sub>	(2, (c))

**Table 4:** (a);  $DBX_{\langle b, a \rangle}^{past}$  & (b);  $DBX_{\langle b, a \rangle}^{past-all}$

as  $\text{sup}(P_X)$ , if the database is clear from the context, e.g., it refers to input sequence database  $\text{SeqDB}$ .

**DEFINITION 3.4 (Projected-past-all & Sup-past-all).** A database projected-past-all on a pattern  $p$  is defined as:  $\text{SeqDB}_p^{past-all} = \{(j, px) \mid \text{the } j^{\text{th}} \text{ sequence in SeqDB is } s, \text{ where } s = px ++ sx, \text{ and } sx \text{ is an instance of } p \text{ in } s \text{ and } \text{first}(sx) = \text{first}(p)\}$

Given a pattern  $P_X$ , we define  $\text{sup}^{past-all}(P_X, \text{SeqDB})$  to be the size of  $\text{SeqDB}_{P_X}^{past-all}$ . Reference to the database is omitted if it is clear from the context.

Definition 3.3 captures events occurring before the last temporal point. Definition 3.4 captures events occurring before each temporal point.

**Example.** To illustrate the above concepts, we project and project-all the example database  $DBX$  with respect to the pattern  $\langle b, a \rangle$ . The results are shown in Table 4 (a) & (b) respectively.

The two projection methods' associated notions of  $\text{sup}^{past}$  and  $\text{sup}^{past-all}$  are different. Specifically,  $\text{sup}^{past-all}$  reflects the number of occurrences of  $P_X$  in  $\text{SeqDB}$  rather than the number of sequences in  $\text{SeqDB}$  supporting  $P_X$ .

**Example.** Consider the example database,  $\text{sup}^{past}(\langle b, a \rangle, DBX) = |DBX_{\langle b, a \rangle}^{past}| = 2$ . On the other hand,  $\text{sup}^{past-all}(\langle b, a \rangle, DBX) = |DBX_{\langle b, a \rangle}^{past-all}| = 4$ .

From the above notions of temporal points, projected databases and pattern supports, we can define the support and confidence of temporal rules.

**DEFINITION 3.5 (Support & Confidence).** Consider a temporal rule  $R_X$  ( $\text{pre}_X \hookrightarrow_P \text{post}_X$ ). The support of  $R_X$  is defined as the number of sequences in  $\text{SeqDB}$  where  $\text{pre}_X$  occurs, which is equivalent to  $\text{sup}^{past}(\text{pre}_X, \text{SeqDB})$ . The confidence of  $R_X$  is defined as the likelihood of  $\text{post}_X$  happening before  $\text{pre}_X$ . This is equivalent to the ratio of  $\text{sup}^{past}(\text{post}_X, \text{SeqDB}_{\text{pre}_X}^{past-all})$  to the size of  $\text{SeqDB}_{\text{pre}_X}^{past-all}$ .

**Example.** Consider  $DBX$  and a temporal rule  $R_X, \langle b, a \rangle \hookrightarrow_P \langle c \rangle$ . From the database, the support of  $R_X$  is the number of sequences in  $DBX$  supporting (or is a super-sequence of) the rule's pre-condition –  $\langle b, a \rangle$ . There are 2 of them – see Table 4 (a). Hence support of  $R_X$  is 2. The confidence of the rule  $R_X (\langle b, a \rangle \hookrightarrow_P \langle c \rangle)$  is the likelihood of  $\langle c \rangle$  occurring before each temporal point of  $\langle b, a \rangle$ . Referring to Table 4(b), we see that there is a  $\langle c \rangle$  occurring before each temporal point of  $\langle b, a \rangle$ . Hence, the confidence of  $R_X$  is 1.

Significant rules to be mined must have their supports greater than the  $\text{min\_sup}$  threshold, and their confidences greater than the  $\text{min\_conf}$  threshold.

In mining program properties, the confidence of a rule (or property), which is a measure of its certainty, matters the most (c.f., [34]). Support values are considered to differentiate high confidence rules from one another according to the frequency of their occurrences

in the traces. Rules with confidences less than 100% are also of interest due to the imperfect trace collection and the presence of bugs and anomalies [34]. Similar to the assumption made by work in statistical debugging (e.g., [9]), simply put, if a program behaves in one way 99% of the time, and the opposite 1% of the time, the latter likely corresponds to a possible bug. Hence, a high confidence and highly supported rule is a good candidate for bug detection using program verifiers or runtime monitors.

We add the notions of support and confidence to past-time temporal rules. The formal notation of past-time temporal rules is defined below.

**DEFINITION 3.6 (Past-Time Temporal Rules).** A temporal rule  $R_X$  is denoted by  $\text{pre} \hookrightarrow_P \text{post} (\text{sup}, \text{conf})$ . The series of events  $\text{pre}$  and  $\text{post}$  represent the rule's pre- and post-condition and are denoted by  $R_X.\text{Pre}$  and  $R_X.\text{Post}$  respectively. The notions  $\text{sup}$ , and  $\text{conf}$  represent the support, and the confidence of  $R_X$  respectively. They are denoted by  $\text{sup}(R_X)$  and  $\text{conf}(R_X)$  respectively.

**Example.** Consider  $DBX$  and the rule  $R_X, \langle b, a \rangle \hookrightarrow_P \langle c \rangle$  shown in the previous example. It has support of 2 and confidence of 1. It is denoted by  $\langle b, a \rangle \hookrightarrow_P \langle c \rangle (2, 1)$ .

## 3.2 Monotonicity and Non-Redundancy

Our algorithm is a member of the family of pattern mining algorithms, e.g. [2, 33]. Monotonicity (a.k.a. apriori) properties have been widely used to ensure efficiency of many pattern mining techniques (e.g., [2, 33]). Different mining algorithm often require new or additional apriori property. Fortunately, past-time temporal rules obey the following apriori properties:

**THEOREM 1 (Monotonicity Property – Support).** If a rule  $\text{ev}_S \hookrightarrow_P \text{ev}_C$  does not satisfy the  $\text{min\_sup}$  threshold, neither will all rules  $\text{ev}_Q \hookrightarrow_P \text{ev}_C$  where  $\text{ev}_Q$  is a super-sequence of  $\text{ev}_S$ .

**THEOREM 2 (Monotonicity Property – Confidence).** If a rule  $\text{ev}_S \hookrightarrow_P \text{ev}_C$  does not satisfy the  $\text{min\_conf}$  threshold, neither will all rules  $\text{ev}_S \hookrightarrow_P \text{ev}_D$  where  $\text{ev}_D$  is a super-sequence of  $\text{ev}_C$ .

To reduce the number of rules and improve efficiency, we define a notion of rule redundancy defined based on super-sequence relationship among rules having the same support and confidence values. This is similar to the notion of closed patterns applied to sequential patterns [33].

**DEFINITION 3.7 (Rule Redundancy).** A rule  $R_X$  ( $\text{pre}_X \hookrightarrow_P \text{post}_X$ ) is redundant if there is another rule  $R_Y$  ( $\text{pre}_Y \hookrightarrow_P \text{post}_Y$ ) where:

- (1)  $R_X$  is a sub-sequence of  $R_Y$  (i.e.,  $\text{post}_X ++ \text{pre}_X \sqsubseteq \text{post}_Y ++ \text{pre}_Y$ )
- (2) Both rules' support and confidence are the same

Also, in the case that the concatenations are the same (i.e.,  $\text{post}_X ++ \text{pre}_X = \text{post}_Y ++ \text{pre}_Y$ ), to break the tie, we call the one with the longer premise as being redundant (i.e., we wish to retain the rule with a shorter premise and longer consequent).

To illustrate redundant rules, consider the following set of rules describing an Automated Teller Machine (ATM):

R1	$\text{money\_dispensed} \hookrightarrow_P \text{card\_inserted, enter\_pin, pin\_correct cash\_request}$
R2	$\text{money\_dispensed} \hookrightarrow_P \text{card\_inserted}$
R3	$\text{money\_dispensed} \hookrightarrow_P \text{enter\_pin}$
R4	$\text{money\_dispensed} \hookrightarrow_P \text{card\_inserted, enter\_pin}$
R5	$\text{money\_dispensed} \hookrightarrow_P \text{enter\_pin, cash\_request}$

If all of the above rules have the same support and confidence values, rules R2-R5 are redundant since they are represented by rule R1. To keep the number of mined rules manageable, we remove redundant rules. Noting the combinatorial nature of redundant rules, removing redundant rules can drastically reduce the number of reported rules.

A simple approach to reduce the number of rules is to first mine a full-set of rules and then remove redundant ones. However, this “late” removal of redundant rules is inefficient due to the exponential explosion of the number of intermediary rules that need to be checked for redundancy. To improve efficiency, it is therefore necessary to identify and prune a search space containing redundant rules “early” during the mining process. The following two theorems are used for ‘early’ pruning of redundant rules. The proofs are available in our technical report [27].

**THEOREM 3 (Pruning Redundant Pre-Conds).** *Given two pre-conditions  $P_X$  and  $P_Y$  where  $P_X \sqsubset P_Y$ , if  $SeqDB_{P_X}^{past} = SeqDB_{P_Y}^{past}$  then for all sequences of events  $post$ , rules  $P_X \hookrightarrow_P post$  is rendered redundant by  $P_Y \hookrightarrow_P post$  and can be pruned.*

**THEOREM 4 (Pruning Redundant Post-Conds).** *Given two rules  $R_X$  ( $pre \hookrightarrow_P P_X$ ) and  $R_Y$  ( $pre \hookrightarrow_P P_Y$ ) if  $P_X \sqsubset P_Y$  and  $(SeqDB_{pre}^{past-all})_{P_X}^{past} = (SeqDB_{pre}^{past-all})_{P_Y}^{past}$  then  $R_X$  is rendered redundant by  $R_Y$  and can be pruned.*

Utilizing Theorems 3 & 4, many redundant rules can be pruned ‘early’. However, the theorems only provide sufficient conditions for the identification of redundant rules – there are redundant rules which are not identified by them. To remove remaining redundant rules, we perform a post-mining filtering step based on Definition 3.7.

### 3.3 Mining Steps

Our approach to mining a set of non-redundant rules satisfying the support and confidence thresholds is as follows:

- Step 1** Leveraging Theorems 1 & 3, we generate a *pruned* set of pre-conditions satisfying *min\_sup*.
- Step 2** For each pre-condition *pre*, we create a *projected-past-all* database  $SeqDB_{pre}^{past-all}$ .
- Step 3** Leveraging Theorems 2 & 4, for each  $SeqDB_{pre}^{past-all}$ , we generate a *pruned* set containing post-condition *post*, such that the rule  $pre \hookrightarrow_P post$  satisfies *min\_conf*.
- Step 4** Using Definition 3.7, we filter any remaining redundant rules.

At Step 1, the pruned set of pre-conditions contains those frequent pre-conditions that obey *min\_sup* threshold (following Theorem 1) and are without any super-sequence having the same projected database (following Theorem 3). At Step 3, the pruned set of post-conditions associated with a pre-condition *pre* contains those post-conditions *post* where  $pre \hookrightarrow_P post$  obeys *min\_conf* threshold (following Theorem 2) and are without any super-sequence having the same projected database (following Theorem 4).

In this work, we guarantee that *all mined rules are significant and non-redundant* and *all significant and non-redundant rules are mined* (i.e., statistical soundness and completeness). The first claim is true, due to the fact that we verify each rule’s significance and non-redundancy before outputting it. The second claim is true, due to the fact that the mining process traverses the search space of all possible rules and will only prune those sub-search spaces which are flagged by the four pruning theorems. These sub-search spaces only contain insignificant and redundant rules.

## 4. EXPERIMENTS

In this section we discuss our experiments on mining past-time temporal rules from traces of JBoss Application Server. It shows the utility of our method in recovering specifications of an industrial system.

JBoss AS is the most widely used J2EE application server. It contains over 100,000 lines of code and comments. The purpose of this study is to show the usefulness of the mined rules to describe the behavior of a real software system.

**Case 1: JBoss AS Security Component.** We instrumented the security component of JBoss-AS using JBoss-AOP and generated traces by running the test suite that comes with the JBoss-AS distribution. In particular, we ran the regression tests on Enterprise Java Bean (EJB) security implementation of JBoss-AS. Twenty-three traces of a total size of 4115 events, with 60 unique events, were generated. Running the algorithm on the traces with the minimum support and confidence thresholds set at 15 and 90% respectively, 4 non-redundant rules were mined. The algorithm completed within 2.5 seconds.

A sample of the mined rules is shown in Figure 1 (left). It describes authentication using Java Authentication and Authorization Service (JAAS) for EJB within JBoss-AS. Roughly it describes a rule that states: “Whenever principal and credential information is required (the premise of the rule), previously configuration information is checked to determine authentication service availability (event 1-5 in the consequent), actual authentication events are invoked (event 6-8) and principal information is bound to the subject being authenticated (event 9-12)”.

**Case 2: JBoss AS Transaction Component.** We instrumented the transaction component of JBoss-AS using JBoss-AOP and generated traces by running the test suite that comes with the JBoss-AS distribution. In particular, we ran a set of transaction manager regression tests of JBoss-AS. Each trace is abstracted as a sequence of events, where an event corresponds to a method invocation. Twenty-eight traces with a total size of 2551 events containing 64 unique events, were generated. Running the algorithm on the traces with the minimum support and confidence thresholds set at 25 traces and 90% respectively, 36 non-redundant rules were mined. The algorithm completed within 30 seconds.

A sample of the mined rules is shown in Figure 1 (right). The rule describes that: “Whenever a check is performed on whether transaction is completed (the premise of the rule), previously connection to a server instance (event 1-4 in the consequent), initialization and utilization of transaction manager and implementation (event 5-6,10-12), acquiring of ids (event 7-9,13-15) and obtaining of transaction from the manager (event 16) are performed before.”

## 5. RELATED WORK

One of the most well-known specification mining tool is Daikon [12]. It returns value-based invariants (e.g.,  $x > 5$ , etc.) by monitoring a fixed set of templates as a program is run. Different from Daikon, in this work, we consider temporal invariants capturing ordering constraint among events.

Most specification mining tools mine *temporal* specifications in the form of automata [29, 1, 23]. An automata specify a global behavior of a system. Different from work mining automata, mined rules describe *strongly observed* sub-behaviors of a system or properties that occur with statistical significance (i.e., appear with enough support and confidence).

In [34], Yang et al. present an interesting work on mining two-event future-time temporal logic rules (i.e., of the form  $G(a \rightarrow XF(b))$ , where  $G$ ,  $X$ , and  $F$  are LTL operators [15]), which are

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

Premise	→ <sub>p</sub>	Consequent
TransactionImpl.isDone()		TransManLocator.getInstance() TransManLocator.locate() TransManLocator.tryJNDI() TransManLocator.usePrivateAPI() TxManager.getInstance() TxManager.begin() XidFactory.newXid() XidFactory.getNextId() XidImpl.getTrulyGlobalId() TransImpl.assocCurrentThread() TransImpl.lock() TransImpl.unlock() TransImpl.getLocalId() XidImpl.getLocalId() LocalId.hashCode() TxManager.getTransaction()

**Figure 1: A sample rule from JBoss-Security (left) and another from JBoss-Transaction (right). Each of the rules are read from top to bottom, left to right.**

statistically significant with respect to a user-defined ‘satisfaction rate’. These rules express: “whenever an event occurs, eventually in the future another event occurs”. The algorithm presented, however, does not scale to mine multi-event rules of arbitrary length. To handle longer rules, Yang et al. suggest a partial solution based on concatenation of mined two-event rules. Yet, the method proposed might miss some multi-event rules or introduce additional rules that are not statistically significant – it is neither statistically sound nor complete.

In [25, 26], Lo et al. extended the work by Yang et al. to mine future-time temporal rules of *arbitrary lengths*. The algorithm is statistically sound and complete. Rules of arbitrary lengths is able to capture more complex temporal properties. Often, simple properties are already known by the programmers while complex properties might be missed or might be an emergent behavior.

In [28], Lo et al. mine Live Sequence Chart (LSC) from program execution traces. LSC can be viewed as a formal form of a sequence diagram. In [28], the LSCs mined are of the format: “whenever a chart pre is satisfied, eventually another chart main is satisfied”. Different from standard temporal rules, LSCs impose specific constraints on the satisfaction of a chart (pre or main). When translated to LTL, LSC corresponds to a rather complex temporal expressions [17]. Also, different from this work, the work in [28] only mine LSC that express future time temporal expressions.

In [24], we proposed iterative patterns to mine frequent patterns of program behavior. Different from rules, patterns do not express any constraints. A rule on the other hand expresses a constraint that state when its premise is satisfied, its consequent is satisfied as well. For monitoring and verification purposes, constraints are needed.

There are several studies on extracting specifications from code (e.g. [21, 31]). The above techniques belong to the *static analysis* family. In contrast, we adopt a *dynamic analysis* approach in extracting specifications from execution traces. Static and dynamic analyses complement each other (c.f., [11]). Their pros and cons have been discussed in the literature [11, 5]. With dynamic analysis, even with the best algorithms, the quality of specification mined is only as good as the quality of traces. With static analysis, one is faced with the problem of pointers and infeasible paths. Some specifications pertaining to the dynamic behavior of a system can only be mined (or are much easier to mine) via dynamic analysis.

## 6. DISCUSSION

In this section, we discuss some related issues and potential fu-

ture work.

Similar to pattern mining algorithms (e.g., [2]), one issue is in setting suitable minimum support threshold. It is less straightforward than setting appropriate minimum confidence threshold. In our experiments, we used two different minimum support thresholds: one at 15 (65% of the number of traces), another at 25 (90% of the number of traces). When setting minimum support threshold at 21 (90% of the number of traces) for analyzing traces from the first experiment, we find no significant rules. In general, the more diverse the trace set is, the lower the minimum support threshold needs to be set. Similar to many data mining processes, we view our mining strategy to be an iterative process; at each step, user provides a set of thresholds, runs the miner and evaluates the mined result to decide whether another refinement, by mining at lower or higher thresholds, is needed. As a future work, we plan to formalize and provide tool support to this iterative process to help users decide on a suitable minimum support threshold.

In this work, we guarantee statistical soundness and completeness. Hypothetically, if the input trace set is sound and complete, we will mine a sound and complete set of specifications. Admittedly, traces generated by running a program are generally not complete. As with other dynamic analysis techniques, the effectiveness of the proposed technique is dependent on the quality and sufficiency of the input traces. For example, the pre-conditions of the rule shown in Figure 1 (left) involve multiple occurrences of *SecAssoc.getPrinciple()* and *SecAssoc.getCredential()*. This is the case, as within each trace, in the input trace set, the two methods are called two or more times in tandem. As a future work, we plan to look into employing a synergy of static and dynamic analysis techniques to allow generation of more complete trace set as input to the mining process.

In this study, we use only traces of a few thousand events. In future, we will experiment with traces of longer lengths obtained from a wider variety of software systems. In general, our algorithm will work better with many traces of shorter length, than few very long traces. Employing scenario extraction technique (c.f., [3]) will convert trace of long length to sub-traces of shorter lengths. We plan to employ this and use the modified techniques to mine rules from longer traces. The rules will be fewer as a scenario extraction technique usually employs additional constraints to decide on which events are related to another aside from temporal ordering observed in the traces alone. Also, to handle longer traces, user can provide additional constraints to guide the mining process further. We plan to experiment with these options as a future work.

Also, as a future work, we plan to integrate the proposed mining

technique into automated processes, for example, software verification, runtime monitoring or anomaly detection.

## 7. CONCLUSION

In this paper, we propose a technique to mine past-time temporal rules from program execution traces. The rules state: “Whenever a series of events occurs, previously another series of events has happened”. These rules capture important properties useful for verification, monitoring and program understanding.

Existing work on mining temporal rules focuses on future-time temporal expressions. Past-time temporal logic is more intuitive and compact to express some class of important properties. We consider our work to complement existing techniques mining future-time temporal expressions. To the best of our knowledge, this is the first work on mining past-time temporal rules from program execution traces where repetitions due to loop and recursion need to be considered. Our rule format is also more general than the precedence rule mined by static-analysis-based approach in [31].

Also, the problems of a potentially exponential runtime cost and a huge number of reported rules have been effectively mitigated by employing search space pruning strategies and elimination of redundant rules. Experiments on JBoss Application Server show the utility of our technique in recovering specifications of an industrial program.

**Acknowledgement** We would like to thank Jianyong Wang and Jiawei Han for allowing us to use BIDE, a closed sequential pattern miner [33]. Our proposed algorithm involves making calls to a modified extended version of BIDE (see our previous work in [26] on mining future-time temporal rules).

## 8. REFERENCES

- [1] M. Acharya, T. Xie, J. Pei, and J. Xu. Mining API Patterns as Partial Orders from Source Code: From Usage Scenarios to Specifications. In *SIGSOFT FSE*, 2007.
- [2] R. Agrawal and R. Srikant. Mining sequential patterns. In *Proc. of IEEE Int. Conf. on Data Engineering*, 1995.
- [3] G. Ammons, R. Bodik, and J. R. Larus. Mining Specification. In *POPL*, 2002.
- [4] N. Baskiotis, M. Sebag, M.-C. Gaudel, and S. Gouraud. A machine learning approach for statistical software testing. In *Proc. of Int. Joint Conf. on Artificial Intelligence*, 2007.
- [5] L. C. Briand, Y. Labiche, and J. Leduc. Toward the reverse engineering of UML sequence diagrams for distributed java software. *IEEE TSE*, 32(9):642–663, 2006.
- [6] E. Clarke, O. Grumberg, and D. Peled. *Model Checking*. MIT Press, 1999.
- [7] S. Deelstra, M. Sinnema, and J. Bosch. Experiences in software product families: Problems and issues during product derivation. In *SPLC*, 2004.
- [8] M. Dwyer, G. Avrunin, and J. Corbett. Patterns in property specifications for finite-state verification. In *ICSE*, 1999.
- [9] D. Engler, D. Y. Chen, S. Hallem, A. Chou, and B. Chelf. Bugs as deviant behavior: A general approach to inferring errors in systems code. In *Proc. of Symp. on Operating Systems Principles*, 2001.
- [10] L. Erlikh. Leveraging legacy system dollars for e-business. *IEEE IT Pro*, pages 17–23, 2000.
- [11] M. Ernst. Static and dynamic analysis: Synergy and duality. In *Work. on Dynamic Analysis*, 2003.
- [12] M. Ernst, J. Cockrell, W. Griswold, and D. Notkin. Dynamically discovering likely program invariants to support program evolution. *TSE*, 27(2):99–123, 2001.
- [13] R. Fjeldstad and W. Hamlen. Application program maintenance-report to our respondents. *Tutorial on Software Maintenance*, pages 13–27, 1983.
- [14] J. Han and M. Kamber. *Data Mining Concepts and Techniques*. Morgan Kaufmann, 2006.
- [15] M. Huth and M. Ryan. *Logic in Computer Science*. Cambridge, 2004.
- [16] JBoss. <http://www.jboss.org>.
- [17] H. Kugler, D. Harel, A. Pnueli, Y. Lu, and Y. Bontemps. Temporal Logic for Scenario-Based Specifications. In *TACAS*, 2005.
- [18] F. Laroussinie, N. Markey, and P. Schnoebelen. Temporal logic with forgettable past. In *Proc. of IEEE Symp. on Logic in Computer Science*, 2002.
- [19] F. Laroussinie and P. Schnoebelen. A hierarchy of temporal logics with past. *Theoretical Computer Science*, 2(148):303–324, 1995.
- [20] J. Li, H. Li, L. Wong, J. Pei, and G. Dong. Minimum description length principle: Generators are preferable to closed patterns. In *AAAI*, 2006.
- [21] Z. Li and Y. Zhou. PR-miner: Automatically extracting implicit programming rules and detecting violations in large software code. In *SIGSOFT FSE*, 2005.
- [22] O. Lichtenstein, A. Pnueli, and L. D. Zuck. The glory of the past. In *Proc. Logics of Programs Workshop*, pages 196–218, 1985.
- [23] D. Lo and S.-C. Khoo. SMaRTIC: Towards building an accurate, robust and scalable specification miner. In *SIGSOFT FSE*, 2006.
- [24] D. Lo, S.-C. Khoo, and C. Liu. Efficient mining of iterative patterns for software specification discovery. *Proc. of SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, 2007.
- [25] D. Lo, S.-C. Khoo, and C. Liu. Mining temporal rules from program execution traces. In *Proc. of Int. Work. on Program Comprehension through Dynamic Analysis*, 2007.
- [26] D. Lo, S.-C. Khoo, and C. Liu. Efficient mining of recurrent rules from a sequence database. In *Proc. of Int. Conf. on Database Systems for Advanced Apps.*, 2008.
- [27] D. Lo, S.-C. Khoo, and C. Liu. Mining past-time temporal rules. *NUS SoC Technical Report - TRA5/08*, 2008.
- [28] D. Lo, S. Maoz, and S.-C. Khoo. Mining Modal Scenario-Based Specification from Execution Traces of Reactive Systems. In *ASE*, 2007.
- [29] D. Lorenzoli, L. Mariani, and M. Pezzè. Automatic Generation of Software Behavioral Models. In *ICSE*, 2008.
- [30] Z. Manna and A. Pnueli. *The Temporal Logic of Reactive and Concurrent Systems: Specification*. Springer, 1992.
- [31] M. Ramanathan, A. Grama, and S. Jagannathan. Path-sensitive inference of function precedence protocols. In *ICSE*, 2007.
- [32] T. Standish. An essay on software reuse. *IEEE TSE*, 5(10):494–497, 1984.
- [33] J. Wang and J. Han. BIDE: Efficient mining of frequent closed sequences. In *Proc. of IEEE Int. Conf. on Data Engineering*, 2004.
- [34] J. Yang, D. Evans, D. Bhardwaj, T. Bhat, and M. Das. Perracotta: Mining temporal API rules from imperfect traces. In *ICSE*, 2006.