

2006

# Automated caching of behavioral patterns for efficient run-time

Natalia Stakhanova

*Iowa State University*

Samik Basu

*Iowa State University*, sbasu@iastate.edu

Robyn R. Lutz

*Iowa State University*, rlutz@iastate.edu

Johnny S. Wong

*Iowa State University*, wong@iastate.edu

Follow this and additional works at: [http://lib.dr.iastate.edu/cs\\_techreports](http://lib.dr.iastate.edu/cs_techreports)



Part of the [Computer Sciences Commons](#)

---

## Recommended Citation

Stakhanova, Natalia; Basu, Samik; Lutz, Robyn R.; and Wong, Johnny S., "Automated caching of behavioral patterns for efficient run-time" (2006). *Computer Science Technical Reports*. 211.

[http://lib.dr.iastate.edu/cs\\_techreports/211](http://lib.dr.iastate.edu/cs_techreports/211)

This Article is brought to you for free and open access by the Computer Science at Iowa State University Digital Repository. It has been accepted for inclusion in Computer Science Technical Reports by an authorized administrator of Iowa State University Digital Repository. For more information, please contact [digirep@iastate.edu](mailto:digirep@iastate.edu).

---

# Automated caching of behavioral patterns for efficient run-time

## **Abstract**

Run-time monitoring is a powerful approach for dynamically detecting faults or malicious activity of software systems. However, there are often two obstacles to the implementation of this approach in practice: (1) that developing correct and/or faulty behavioral patterns can be a difficult, labor-intensive process, and (2) that use of such pattern-monitoring must provide rapid turn-around or response time. We present a novel data structure, called extended action graph, and associated algorithms to overcome these drawbacks. At its core, our technique relies on effectively identifying and caching specifications from (correct/faulty) patterns learnt via machine-learning algorithm. We describe the design and implementation of our technique and show its practical applicability in the domain of security monitoring of sendmail software.

## **Disciplines**

Computer Sciences

# Automated caching of behavioral patterns for efficient run-time monitoring

Natalia Stakhanova<sup>1</sup> Samik Basu<sup>1</sup> Robyn R. Lutz<sup>1,2\*</sup> Johnny Wong<sup>1</sup>

<sup>1</sup>*Department of Computer Science  
Iowa State University  
Ames, IA 50011 USA*

<sup>2</sup>*Jet Propulsion Laboratory  
California Institute of Technology  
Pasadena, CA 97709 USA*

*{ndubrov, sbasu, rlutz, wong}@iastate.edu*

## Abstract

*Run-time monitoring is a powerful approach for dynamically detecting faults or malicious activity of software systems. However, there are often two obstacles to the implementation of this approach in practice: (1) that developing correct and/or faulty behavioral patterns can be a difficult, labor-intensive process, and (2) that use of such pattern-monitoring must provide rapid turn-around or response time. We present a novel data structure, called extended action graph, and associated algorithms to overcome these drawbacks. At its core, our technique relies on effectively identifying and caching specifications from (correct/faulty) patterns learnt via machine-learning algorithm. We describe the design and implementation of our technique and show its practical applicability in the domain of security monitoring of sendmail software.*

## 1 Introduction

Run-time monitoring is a proven technique for enhancing the reliability of a system. It observes the behavior of the system during execution and detects anomalous deviations from normal or expected behavior. Early indications of these deviations from expected behavior are frequently useful from a reliability perspective. They may indicate possible movement of the system from a safe to an unsafe state (e.g., from an aerodynamically stable to unstable state [21]), from a secure to an insecure state (e.g., a sequence of system calls characterizing behavior of malicious intruder [18]), or from a low-risk to a high-risk state (e.g., an unexpected load of users [8]). By providing early warning of possible, imminent risk in the dynamic execution environment, run-time moni-

toring is able to complement efforts to increase reliability of software via traditional testing and sound software development practices. In essence, run-time monitoring utilizes the knowledge of normal and/or abnormal system behavior and identifies problems if the system execution deviates from known normal behavior or follows a pre-specified abnormal scenario.

In this aspect, run-time monitoring resembles intrusion detection which aims at discovering malicious deviations from the expected program behavior. In intrusion detection field existing approaches can be classified into (a) misuse-based (b) anomaly-based and (c) specification-based [24]. Misuse-based technique relies on pre-specified attack signatures, and any execution sequence matching with a signature is flagged as abnormal. An anomaly-based approach, on the other hand, typically depends on normal patterns, and any deviation from normal is classified as malicious or faulty. Unlike misuse-based detection, anomaly-based techniques can detect previously unknown abnormalities. However, anomaly-based approaches rely on machine learning techniques which can only classify pre-specified, fixed-length behavioral patterns, and suffer from the disadvantage of a high rate of false positives [17]. Specification-based techniques operates in a similar fashion to anomaly-based method and detect deviations from the specified legitimate system behavior. However, as opposed to anomaly detection, specification-based approach requires user guidance in developing model of valid program behavior in a form of specifications. This process, though tedious and reliant on user-expertise, can handle variable-length sequences and is, therefore, more accurate than anomaly-based techniques.

In this paper, we present a monitoring technique which combines the advantages of two intrusion detection approaches: anomaly-based and specification-based detection. Instead of manually developing possible variable-length legal behavioral patterns of a system, the approach relies on machine-learning tech-

---

\*The author is supported in part by NSF grants 0204139 & 0205588

nique to automatically classify system behavior at runtime, as *correct* or *incorrect* and infer classification of variable-length sequences.

To efficiently maintain the results of classification, we propose a novel structure *EXtended ACTION graph (Exact)* that appropriately combines multiple sequences classified by machine learning technique into variable-length patterns and memorizes them for future reference. In our framework, we have two **Exact**: one for storing normal patterns and the other for abnormal patterns. Sequences are classified using **Exact**, and the machine learning algorithm is only invoked if necessary. The following summarizes the contributions of this work:

1. *Exact structure.* **Exact** allows compact and exact representation of variable-length sequences.
2. *Automatic development of specifications.* While machine-learning technique automatically classifies fixed-length patterns, **Exact** caches the results of classification in such a way that variable-length sequences can be classified in future.
3. *Efficiency.* We describe efficient algorithms for insertion of new patterns into **Exact** graph and identification of existing patterns using **Exact** graph.

The remainder of the paper is organized as follows: We present a brief overview of related work in Section 2. Section 3 presents an overview of the integrated framework and its components. Section 4 describes the **Exact** structure and the novel algorithms used in the framework. Section 5 gives a brief description of machine learning-based classification. Analysis of **Exact** and SVM followed by the experimental results are presented in Sections 6 and 7 respectively. We conclude with the discussion of the significance of the results in Section 8.

## 2 Related Work

The practical benefits of the automatic detection of software errors and vulnerabilities have been noted by many researchers and a number of techniques ranging from static program analysis [27, 12], model checking [11, 15, 1, 14], theorem proving [25], to run-time monitoring of software executions [6, 13, 3] have been proposed over the last two decades. As our work is based on dynamic analysis of software and machine learning based classification, we will primarily focus on related work that proposes run-time monitoring and/or classification technique for software debugging.

In the field of dynamic analysis for detecting source-code errors, Ernst et al.[6] have developed a

dynamic invariant detection technique to determine fault-revealing properties of programs. Subsequent work by Hangal and Lam [13] used this technique to detect code errors by dynamically extracting invariants and checking for their violations through program execution. The approach associates a set of expressions at various program points to derive invariants that satisfy all expression-valuations. The program behavior is then checked against invariants for violations. More rigorous work in this direction was done later by Brun et al.[3]. These approaches are close in spirit to dynamic program analysis and are specifically designed to detect errors in source-code, while our approach relies on observable system behavior (control and data sequences).

The approach used in this paper was inspired by the technique proposed by Sekar et al. [24]. Their work aimed at augmenting machine learning techniques with high-level specifications to achieve a high degree of precision in detecting anomalies in software behavior. The authors showed that the sliding window technique [10] using a machine learning algorithm may be excessively error-prone due to its inability to classify sequences of varying length. They thus manually developed high-level specifications (as finite state machines) of software systems and annotated them using information learnt via machine-learning techniques. However, such manual development of specification is tedious and requires expert-knowledge. In contrast, we propose to generate specifications in the form of variable-length patterns automatically classified via machine learning. Specifically, we use one-class SVM, capable of classifying fixed-length patterns, to identify pattern-classification and infer the classification of variable-length patterns from aggregation of the results.

While sliding window technique is a common way of modeling system data, there have been several approaches to dealing with variable-length patterns. Debar[5] proposed generation of variable-length sequences based on suffix trees augmented with a number of occurrences of each subsequence. Similarly, Marceau[20] employed suffix tree as underlying structure for constructing finite state machine with states representing predictive sequences of variable length. Kosoresow and Hofmeyr[16] manually constructed finite automaton based on variable-length patterns and applied it for detection process.

Eskin et al[7] proposed an alternative algorithm for determining optimal sliding window size depending on the data context as different window sizes might be optimal at different points in the process. The approach based on Sparse Markov Transducers (SMTs), extension of probabilistic suffix trees, allows to consider a mixture of possible trees and estimate the best

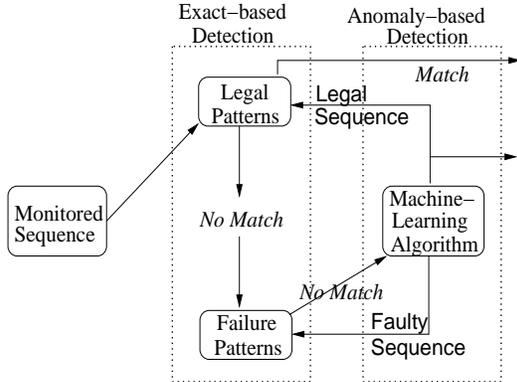


Figure 1: Framework architecture.

tree for the given data. While proposed algorithm provides a good prediction for variable-length patterns in a particular data set, it is static in nature as it does not allow an update of the prediction tree as more system data becomes available.

### 3 Multi-Level Classifier

Our model for monitoring execution sequences of a software system consists of a two-level classification mechanism (Figure 1). Sequences in this context are defined over the observable actions performed by the system, e.g., commands issued by a controller or system calls invoked by a device driver. Observable actions represent the alphabets of the monitored sequences.

Specification of correct/legal and faulty behavioral patterns are provided in the first level in the form of **Exact**. In the event that the sequence to be monitored matches the specifications, the second level classification is not invoked. A sequence that matches legal specifications is allowed to execute unaltered while a faulty sequence is blocked and/or appropriate evasive actions (such as intrusion response) are fired.

If the sequence is *not* found in the specification module, the second-level classifier is used. We then rely on machine-learning techniques to determine whether the sequence is normal or anomalous. In either case, the sequence is recorded in the corresponding **Exact** specifications for future reference. In the domain of software reliability monitoring, a faulty behavior may result in re-visiting the requirements and correcting the implementation. In that case, the revised implementation can be monitored against prior faulty behaviors, recorded in the specification, to rule out the presence of the same errors.

One of the important features of our model is that the technique can be deployed with empty or partially

filled **Exact** in the first level. As more sequences are classified by the second level, the **Exact**-level is populated automatically. There are three advantages to this technique. *First*, in addition to acting as a cache for pre-specified classification results, **Exact** also allows future classification of patterns of any size. *Second*, **Exact** is similar to low-level specifications of system behavior. In other words, the framework is generating specifications of system behavior automatically. *Finally*, these **Exact** specifications can be used to retrain the second-level classifier as more new patterns become available.

## 4 Extended Action Graph: Exact

As noted in Section 3, **Exact** is used to record previously classified behavioral patterns. In **Exact**, which is a graphical representation of multiple sequences of varying length, states are annotated by observable actions of the system to be monitored, and a sequence of states represents a behavioral pattern.

**Definition 1 (Exact)** *An Extended Action Graph is a tuple  $E = (S, S_0, \rightarrow, \Sigma, L)$  where  $S$  is the set of states,  $S_0 \subseteq S$  is the set of start states,  $\Sigma$  is a set of binary numbers used to represent transition,  $\rightarrow \subseteq S \times \Sigma \times S$  is the set of transition relations, and  $L : S_0 \rightarrow \Sigma$  is a mapping of start states to a binary vector.*

A sequence in **Exact** is represented by  $s_1, s_2, \dots, s_n$  where each  $s_i$  has a transition to  $s_{i+1}$ . Consider the example in Figure 2(a). The action graph, that was generated by three sequences  $s_1, s_2, s_3, s_2$ , and  $s_2, s_4, s_5, s_1, s_3, s_6$ , and  $s_1, s_3, s_6$ , has six states and two start states  $s_1$  and  $s_2$ . Each transition and the start states are labeled by a binary vector; e.g.,  $L(s_1) = 101$ .

However, not all the sequences in **Exact** are classified as valid and valid sequences form a superset of the known sequences. In the above example,  $s_1, s_2, s_3$  and  $s_1, s_3, s_6$  are valid patterns, and the graph also contains the sequence  $s_1, s_2, s_3, s_6$  which is not valid.

To rule out invalidity, we use the transition label  $\sigma \in \Sigma$ , a binary vector, whose  $k$ -th element is denoted by  $\sigma[k]$ . If there exists a transition  $s_i \xrightarrow{\sigma} s_j$  where  $\sigma_i[k] = 1$ , then  $s_i, s_j$  are said to be consecutive alphabets in the  $k$ -th valid sequence. Note that the first sequence is identified by setting the rightmost bit to 1, i.e., 001 is the identifier for the first sequence, 010 is the identifier for the second sequence and so on. In Figure 2(a),  $s_1, s_2$  are consecutive states in the first pattern while  $s_1, s_3$  are consecutive states in the second and third valid sequences. Every valid sequence is assigned a start state:  $s_2$  is the start state of

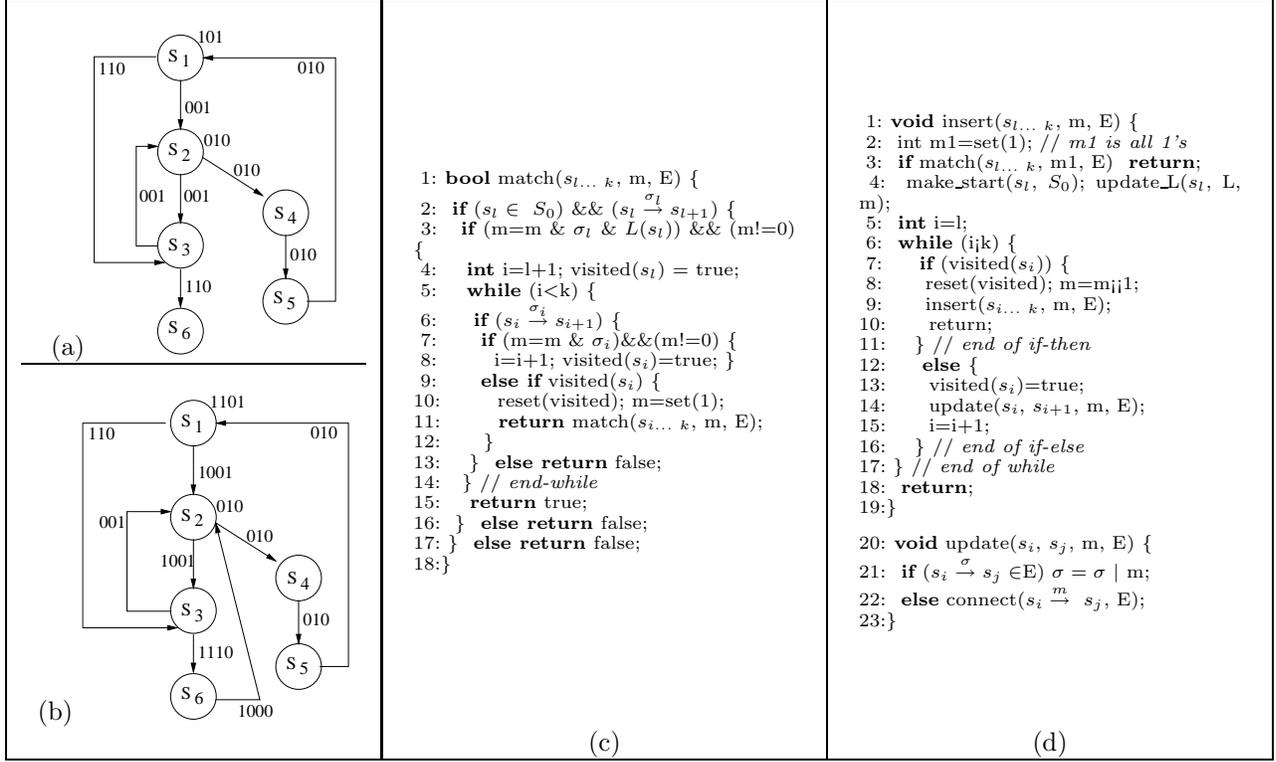


Figure 2: (a) Example of an **Exact** graph that was generated by three sequences  $s_1, s_2, s_3, s_2$ , and  $s_2, s_4, s_5, s_1, s_3, s_6$ , and  $s_1, s_3, s_6$ . The graph has six states and two start states  $s_1$  and  $s_2$ . (b) **Exact** graph in Figure 2(a) shown after insertion of  $s_1, s_2, s_3, s_6, s_2, s_4, s_5$ . (c) Pseudo-code for **Exact** search. (d) Pseudo-code for **Exact** insert.

the second valid sequence. Formally, a valid sequence is defined as follows:

**Definition 2 (Validity)** A sequence  $s_1, s_2, \dots, s_n$  is said to be valid if  $s_1 \in S_0$ ,  $L(s_1) = \sigma_s$  and

$$\exists k \forall 1 \leq i < n : s_i \xrightarrow{\sigma_i} s_{i+1} \wedge (\sigma_s[k] = 1 \wedge \sigma_i[k] = 1)$$

In other words, there exists a specific element in the vector-label of each transition in this sequence and the vector-identifier of the start state which is equal to 1.

Furthermore, via transitivity, if  $S_1 = s_i, s_{i+1}, \dots, s_{i+n}$  is a valid sequence and  $S_2 = s_j, s_{j+1}, \dots, s_{j+m}$  is another valid sequence such that  $s_{i+n} = s_j$  then  $s_i, s_{i+1}, \dots, s_{(i+n)-1}, s_j, s_{j+1}, \dots, s_{j+m}$  is also a valid sequence.

Validity takes care of unbounded (one or more) repetition of the alphabets in a sequence, e.g., in Figure 2(a)  $s_1, s_2, s_3, s_2, s_3, \dots$  is a valid sequence. In the above,  $s_2$  is said to be the root of the loop. A sequence representing *exit* from a loop is classified as a *new* sequence starting from the root of the loop. In Figure 2(a),  $s_2, s_4, s_5, \dots$  is the valid sequence exiting from the loop rooted at  $s_2$ . Note that the transitivity relation in Definition 2 can be used to identify valid

sequences with bounded repetition (from a valid sequence with finite looping and a valid exit sequence). For example in Figure 2(a),  $s_1, s_2, s_3, s_2$  and  $s_2, s_4, s_5$  are valid sequences and they form, via transitivity, a new valid sequence  $s_1, s_2, s_3, s_2, s_4, s_5$ .

#### 4.1 Searching and constructing **Exact**

Figure 2(c) presents the algorithm to find out whether a given sequence is a valid sequence in **Exact**. Procedure **match** takes as input the given sequence  $s_{l...k}$ , a bit-vector  $m$  and the **Exact** and returns true if the sequence is present as a valid sequence in **Exact**. **Exact** is deterministic, i.e., for every pair of states there exists at most one transition. Absence of non-determinism makes the complexity of searching for a valid sequence linear in the size of the given sequence.

Figure 2(d) presents the algorithm for insertion of a new sequence in **Exact** graph. Procedure **insert** takes as arguments the sequence to be inserted, a bit vector  $m$  identifying the new sequence and the graph **Exact**. The procedure **match** is always invoked before inserting any new sequence to avoid duplicate insertions. As such, the worst case complexity of the insertion algorithm is  $O(r \times n)$  where  $r$  is the number of repeated occurrences of alphabets in the sequence

of length  $n$ .

## 4.2 Illustrative Example

Let  $s_1, s_2, s_3, s_6, s_2, s_4, s_5$  be a sequence to be inserted in the example **Exact** in Figure 2(a). First we break the sequence up into substrings  $s_1, s_2, s_3, s_6, s_2$  (up to the first repeated alphabet  $s_2$ ) and  $s_2, s_4, s_5$  (following the transitivity rule). Then we insert these two subsequences following these steps:

1. This sequence  $s_1, s_2, s_3, s_6, s_2$  is not present in the **Exact** graph shown in Figure 2(a) as there is no transition from  $s_6$  to  $s_2$ . The **match** algorithm (Figure 2(c)), in this case, makes an *early* detection of its absence as the prefix  $s_1, s_2, s_3, s_6$  of the given sequence is not valid in the exact. Observe that, for this prefix, bit-wise “and”-ing of the transition labels and the start state label results in 0 ( $101 \text{ AND } 001 \text{ AND } 001 \text{ AND } 110 = 000$ , see Figure 2(a)) and our **match** algorithm (Figure 2(c)) returns false. As such,  $s_1, s_2, s_3, s_6, s_2$  is inserted as a new sequence. Further, as  $s_1$  is already present in the set of start states, its start-state label is updated using the new sequence identifier 1000. Recall that the identifier for the first sequence is 001, identifier for the second sequence is 010 and for the third sequence is 100. As such the identifier for the new (fourth) sequence is 1000.
2. Each transition of the new sequence is added to **Exact** graph with identifier 1000. We start with transition  $s_1 \rightarrow s_2$ . It already exists in the graph and its identifier is 001. Applying *bitwise-OR* of the existing and new transition label we obtain 1001 and update the transition with this new label (Figure 2(b)). We continue in this fashion until we reach transition substring  $s_6, s_2$ . There is no transition between  $s_6$  and  $s_2$ . A new transition ( $s_6 \rightarrow s_2$ ), therefore, is added with the transition label 1000.
3. Due to the repeated appearance of  $s_2$ , the second sequence  $s_2, s_4, s_5$  is set up to be inserted as a new sequence with a new sequence number, 10000. However, its insertion is avoided as the sequence  $s_2, s_4, s_5$  is already a valid sequence in **Exact**. The updated **Exact** graph is shown in Figure 2(b).

## 5 Second-level Classifier

Any machine learning technique can be applied as a second-level classifier in our framework. For our case study, we used one-class support vector machine

(SVM) which allows usage of *unlabeled data*, i.e., unsupervised learning. As opposed to its more classical version, two-class SVM [26], one-class SVM relies on maximally separating all data from origin using a hyperplane. See [23] for details.

The unlabeled data, in our case, are sequences of observable actions representing the system behavior we are interested in monitoring. Observability may be defined in different ways in different application domain; for example, in-flight stability or collision avoidance controllers, we use the SVM to classify *pairs of input and output control signals* [19], while in host-based software intrusion detection, we are interested in classifying *sequences of system calls*. However, in all cases (in adaptive and/or secure systems), we rely on machine-learning technique to classify behavioral patterns as normal and abnormal depending on how well they fit in the learnt data domain.

For the purpose of discussion, we illustrate the application of SVM classifier via an example. Let the observed input stream be **Istream**  $\equiv s_1, s_2, s_3, s_2, s_3, s_4, s_2$  and **Exact** in Figure 2(b) failed to recognize **Istream** as a valid sequence. First, we break-up **Istream** following the transitivity relationship in Definition 2 of Section 4, i.e., **Seq**<sub>1</sub>  $\equiv s_1, s_2, s_3, s_2$  and **Seq**<sub>2</sub>  $\equiv s_2, s_3, s_4, s_2$ . Note that the break-up point is at  $s_2$  which appears in **Seq**<sub>1</sub> and **Seq**<sub>2</sub>, and is the first alphabet to be repeated in **Istream**. SVM can only take fixed length sequences as input and as such we apply classic sliding window technique to provide inputs to the SVM. Let the sliding window size be 3, then SVM is fed with subsequences: (i)  $s_1, s_2, s_3$ , (ii)  $s_2, s_3, s_2$  (from **Seq**<sub>1</sub>), (iii)  $s_2, s_3, s_4$  and (iv)  $s_3, s_4, s_2$  (from **Seq**<sub>2</sub>). Finally, **Seq**<sub>1</sub> and **Seq**<sub>2</sub> are termed as normal if and only if all their subsequences are classified by SVM as normal. Note that, break-up of **Seq**<sub>1</sub> and **Seq**<sub>2</sub> using sliding window does not adversely effect end result, i.e., if any subsequence of **Seq**<sub>1</sub>/**Seq**<sub>2</sub> is classified as anomalous, then the corresponding sequence is conservatively classified as anomalous. Furthermore, the sequences **Seq**<sub>1</sub> and **Seq**<sub>2</sub> provide an easy way of inserting **Istream** in the corresponding **Exact**.

## 6 Analysis of Exact

As **Exact** represents variable-length sequences, the comparison with models based on the fixed-length patterns is challenging; the main challenge being the difference in the number of fixed-length and variable-length sequences generated from the same data set. The comparison is also aggravated by the fact that classification of variable-length patterns in **Exact** depends entirely on the underlying fixed-length classifier (SVM in this case).

In this section, we present a comparative study of number of sequences being examined by **Exact** and the SVM classifier. We consider two possible cases: one where the SVM, used in conjunction with **Exact**, acts as the backend for our framework (backend SVM) and the other where SVM acts alone (stand-alone SVM). The comparison will form the basis for results presented in the Section 7.

For the purpose of analysis, we will consider average length of **Exact** sequences; the average length computed using the weighted mean where the weight of a length denotes the number (frequency) of sequences of the corresponding length. We represent this length as  $L$ . Let the fixed sliding window size of SVM be  $W$ .

**Exact Vs Backend-SVM** The two possible scenarios of interest are  $W > L$  and  $W < L$ . For  $W = L$ , the number of **Exact** sequence and SVM sequence is identical.

1.  $W > L$ : In this case, several **Exact** sequences are combined to form one SVM sequence. Consider first the case where  $x$  **Exact** sequences fit *exactly* in one SVM sequence of size  $W$ . In other words,  $xL - (x - 1) = W$  (the subtraction of  $x - 1$  from  $xL$  is required to account for overlap between two consecutive **Exact** sequences). Therefore,

$$x = \frac{W - 1}{L - 1} \quad (1)$$

In other words, the number of **Exact** sequences is greater than the number of SVM sequences and classification of one SVM sequence influences the classification of  $x$  **Exact** sequences.

Secondly, consider the case where  $(W - 1)/(L - 1)$  is not a whole number, i.e. the **Exact** sequences is not a integer-multiple of SVM sequences. Let  $x$  be the smallest number of **Exact** sequences such that  $xL - (x - 1) > W$  and  $\forall y < x : yL - (y - 1) < W$ . Therefore, the number of SVM sequences corresponding to  $x$  **Exact** sequences is,  $xL - (x - 1) - W + 1 = x(L - 1) - W + 2$ . Then the number of **Exact** sequences is greater than the number of SVM sequences if  $x < (W - 2)/(L - 2)$ ; otherwise the number of SVM sequences is greater than the number of **Exact** sequences. In case of former, one SVM sequence classification influences one **Exact** sequence classification while in latter, one SVM sequence can potentially effect  $x$  **Exact** sequences.

2.  $W < L$ : In this case, the number of **Exact** sequences is less than the number of SVM se-

quences. Specifically, the number of SVM sequences corresponding to one **Exact** sequence is  $(L - W + 1)$  and therefore, one SVM sequence classification can effect the classification of one **Exact** sequence.

**Exact Vs. Stand-alone SVM** Next, we consider the number of sequences examined by SVM if it is deployed alone. Given that the total length of the input stream is  $IS$ , the total number of SVM sequences is  $N = IS - W + 1$ . If the same input stream is input to our framework – **Exact** with backend-SVM – the total number of **Exact** sequences is  $(IS - 1)/(L - 1)$ , i.e.  $(N + W - 2)/(L - 1)$ . The number of sequences examined by SVM alone is greater than the number of **Exact** sequences in our framework if  $N > (W - 2)/(L - 2)$ .

Also, note that if  $W < L$ , the number of sequences examined by SVM, when deployed alone, can be potentially greater than the number of sequences examined by SVM, when deployed in conjunction to **Exact**. Specifically, the situation requires  $N > L - W + 1$  and can be explained from the fact that number of sequences classified by backend-SVM depends on the number of **Exact** sequences when  $W < L$ .

## 7 Case Study

We evaluated our model in the security domain using synthetic sendmail data provided by the UNM [9]. Sendmail data is an unlabeled collection of system calls. It consists of a normal data sets which contain only legal (normal) patterns and trace data sets containing normal patterns as well as anomalies. We considered three intrusion trace data sets: snsnd-mailcp, decode and fwdloops. The one-class SVM classifier was trained on the normal data set(training set), tested on the trace sets (test sets) and implemented using libsvm tool [4] and the window size of 8.

**Data Sets.** Table 1 presents the pattern of data being used for evaluation pupose in terms of number of sequences. The training data set contains 30792 normal fixed-length sequences. On the other hand, using **Exact**, the number of variable-length sequences is 3314. The decrease in the number of sequences is due to the fact that **Exact** partitions sequences using repetitions and as such can handle variable-length sequences (see Section 4). We, then, processed the normal and abnormal patterns of the test data set to generate two test sets: one for stand-alone SVM, containing fixed-length sequences obtained through sliding window technique, and one for **Exact**, con-

	stand-alone SVM			Exact (variable-length sequences)		
	snsndmailcp	decode	fwdloop	snsndmailcp	decode	fwdloop
Number of normal sequences in train data set	30792	30792	30792	3314	3314	3314
Number of sequences in test data set	1098	2983	2499	78	405	204
Number of anomalous sequences in test data set	264	741	387	24	92	43

Table 1: Information on sendmail normal and intrusive trace data sets

taining variable-length sequences generated in **Exact** fashion. (row 2 in Table 1).

Finally, the last row shows the number of sequences that are in the test data set but are not present in the training data set. For example, out of 1098 fixed-length sequences for **snsndmailcp**, there are 264 sequences which are not present in fixed-length sequences of training data. For the purpose of evaluation, we can conservatively assume that sequences not present in the training data set are anomalous; the goal is to identify all such anomalous sequences.

**Efficiency.** In these experiments we focused primarily on the rate of populating the **Exact** with normal(legal) and abnormal(anomalous) patterns. To evaluate our technique we monitored the stage at which each sequence was classified. We examined two scenarios:

1. Both **Exact** graphs representing normal and abnormal specifications are initially empty
2. Partial specification is available initially, i.e., the **Exact** graph corresponding to a normal specification is populated with 10% of the patterns from the normal data set.

The results for both scenarios are presented in Figures 3, 4 and 5.

Figure 3 shows the frequency at which both levels of classifiers (**Exact** and SVM) were invoked for classifying the incoming sequences. Since simulation started with empty **Exact** graph, almost every incoming sequence is classified at the second-level classifier. However, the access rate of second-level classifier rapidly decreases as more patterns were stored in the **Exact**. Consequently, the number of sequences classified at the **Exact** graph level increases. Figure 4 shows the number of new patterns added to empty **Exact** over the same run of decode trace set. The majority of patterns were recorded within about 200 sequences. After that almost all patterns were found at the **Exact** level.

The result corresponding to the second scenario where normal **Exact** graph is partially populated is

shown in Figure 5. As opposed to the previous figure, the access rate of the second-level classifier in the beginning of the run is low while the **Exact** graph access rate is high. This is explained by the partial presence of the sequences in the normal **Exact** specifications. However, since only partial normal patterns were added to the specifications, the second-level classifier was still accessed whenever new normal or anomalous sequence is found.

In this scenario we benefited from the available specifications having populated the **Exact** in advance. This shortened the start-up time necessary to store a sufficient number of patterns (Table 6)<sup>1</sup>. In fact, the processing time for 405 sequences was 2 times faster with the populated specifications (7 sec) than with the empty specifications (16 sec). Note it is the SVM classifier access that requires most of this time.

**Accuracy.** As the **Exact** graph provides a succinct representation of learned through machine-learning technique variable-length patterns, we focused in these experiments on the comparison of the accuracy of our structure to the accuracy of SVM tested on model built using sliding window technique.

As evaluation criteria we considered *detection rate* (ratio of detected anomalies to the total number of anomalies presented in the set) and *false positive rate* (FP) (number of normal instances incorrectly identified as anomalous). To determine the accuracy of the classification, recall that (Table 1), the test set was labeled in the following way [10]: instances present in the normal data set were labeled as normal, other instances were marked as anomalies.

As Table 3 show, classification results of fixed-length patterns for stand-alone SVM and **Exact** integrated with SVM are similar. For example, for **snsndmailcp**, the detection rate is 98% for both stand-alone SVM and back-end SVM used in **Exact**. The results confirm the fact that **Exact** structure, while recording variable-length patterns, truly represents information given by the backend machine-learning based classifier in compact fashion. Existing 1–2% variation in the results is explained the poten-

<sup>1</sup>Average over 10 runs.

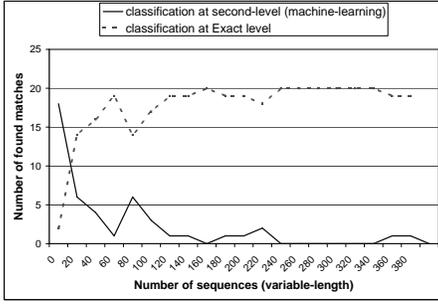


Figure 3: Initialization: empty **Exact**

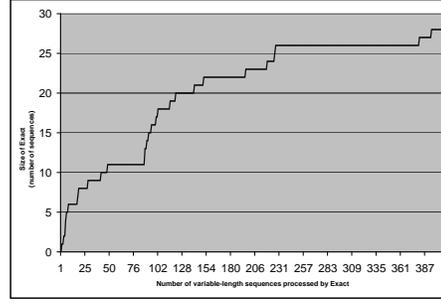


Figure 4: Populating **Exact** (*decode* intrusion).

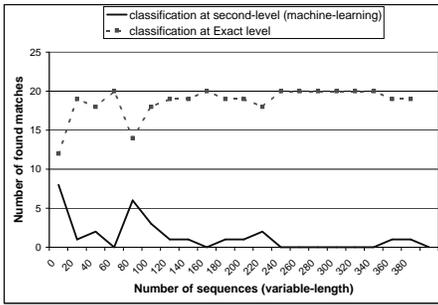


Figure 5: Initialization: **Exact** partially populated with 10% of normal sequences (*decode* intrusion).

	Total time(sec)	Backend SVM running time(sec)
<b>Exact</b> with empty specs	16	12
<b>Exact</b> with partially populated specs	7	5

Figure 6: Running time(*decode* intrusion).

tial difference in the number of sequences examined by stand-alone SVM and back-end SVM used in **Exact** as discussed in Section 6.

The classification results are also given in terms of variable-length patterns, stored by **Exact** (Table 4). Examining Table 4, we notice that prediction results are slightly different from the corresponding percentages given in Table 3. For example, detection rate of **Exact** for `snsndmailcp` intrusion given in fixed-length patterns is 98% while the corresponding number of detected variable-length sequences is 24 out of 24. Recall that that there are 24 anomalous **Exact** sequences in the test data set for `snsndmailcp` (Table 1). This happens when several SVM sequences, including those that are correctly classified as anomalous and those that represent missed intrusions, are effectively combined into one **Exact** sequence resulting in an anomalous **Exact** sequence and providing a higher detection rate.

An opposite scenario is represented by `decode` intrusion, where detection rate in fixed-length patterns is 100% which corresponds to 90 out 92 variable-length **Exact** sequences. Closer inspection reveals that the result is as expected and can be explained

by the fact that **Exact** records sequences depending on the classification result from backend SVM classifier. There are couple of occurrences of one particular **Exact** sequence in the test data set which is not present in the training data set. Hence, this sequence is classified as anomaly (counted as one of the anomalous patterns among 92 anomalies: see Table 1). It turns out that the length of the sequence is 2 due to two consecutive system call-invocation. As such the SVM using sliding window size 8 does not consider this sequence independently; instead it combines the sequence with another (next) **Exact** sequence and performs classification. As the combined sequences are classified as normal by SVM, the **Exact** also records the combined sequence as normal. This is acceptable as the main purpose of **Exact** is to memorize variable length sequence and closely follow SVM classifier. Note that if SVM classifier used window size of 2, then the above scenario will be removed.

The number of variable-length sequences falsely recognized as positive in **Exact** is also different from the corresponding percentages given for fixed-length sequences. This is due to the fact that several SVM sequences can represent one **Exact** sequence, thus, sig-

	snsndmailcp	decode	fwdloop
Exact graph of normal specifications	5	16	13
Exact graph of faulty specifications	14	31	39

Table 2: Maximum length of **Exact** binary vectors

	stand-alone SVM			Exact (results from the backend SVM based on fixed-length sequences)		
	snsndmailcp	decode	fwdloop	snsndmailcp	decode	fwdloop
Detection rate	98%	99%	99%	98%	100%	100%
FP rate	11%	7%	10.7%	13%	8%	10%

Table 3: Accuracy of classification with empty **Exact** shown in fixed-length sequences.

	empty <b>Exact</b>			<b>Exact</b> populated with 10% of normal sequences		
	snsndmailcp	decode	fwdloop	snsndmailcp	decode	fwdloop
Number of detected sequences	24 out of 24	90 out of 92	42 out of 43	24	90	42
FP sequences	21 out of 54	62 out of 313	75 out of 161	0	1	9

Table 4: Accuracy of **Exact** classification shown in variable-length sequences.

nificantly reducing a total number of variable-length sequences in comparison to those in fixed-length. The detailed analysis of such dichotomy was presented Section 6. At the same time, manual inspection of these results showed that a number of FP sequences in **Exact** graph fully comes from the backend SVM classifier.

While the trade-off between number of detection and false positives is inherently present in many machine learning algorithms including SVM, this error can be effectively reduced with guidance from normal specifications. In fact, populating **Exact** even with the small number of normal patterns reduces the number of false positives significantly (Table 4). Since the overall variability of sendmail behavior is small, even approximately 10% of normal sequences leads to recognition of majority of normal patterns. However, generally a greater variability in process behavior might require a larger set of normal patterns to improve the accuracy of classification. Note that, an **Exact** with partially populated normal specification does not affect the number of detected sequences. This is because abnormal, incoming sequences are still recognized as unknown and processed by SVM algorithm as they would be if **Exact** graph were empty.

## 8 Conclusion

In this paper, we present an ongoing work in the development of fast online classification of run-time behavioral patterns for software systems based on the combination of specification and anomaly-based ap-

proaches. We show that memorization of classification results from the SVM classifier can be effectively applied to generate (partial) specifications automatically. We introduce the data structure **Exact** for recording specifications and develop efficient algorithms for insertion and matching of sequences. Finally, our experimental results indicate that our methodology is practical and can be effectively applied in areas such as software testing (run-time monitoring) to complement traditional testing techniques and intrusion detection to provide fast online detection of variable-size anomalous patterns.

Recently, Bowring et al. [2] proposed a technique to model program executions as Markov models and merge similar/redundant models using clustering. They mainly focus on aggregating and predicting program behavior on the basis of small tractable set of features (branches and method calls). In the future work, we plan to investigate the effect of incorporating Bowring’s approach into the generation of **Exact** graph, specifically, by annotating stochastic information with each graph transition.

Another important avenue of future research is the classification of error/fault severity. Machine-learning algorithms have been effectively applied in this context of automatic classification of reported software failures with prioritization according to their relative complexity. Podgurski et al.[22] have applied clustering to group together failures that are likely to have the same or similar cause. The sequences in **Exact** specification for abnormal system-behavior can be effectively classified using these known techniques. Such classification might help to explain the cause of

an abnormality once it is detected using **Exact** specifications.

We also plan to further enhance the efficiency of this approach by combining the **Exact** graphs for normal and faulty specifications. Such a combination may support *early* identification of sequences as normal or abnormal. This may be especially applicable in the domain of adaptable software systems. The main aim of adaptability is to identify abnormalities, apply appropriate adaptation to avoid failures. Our technique would allow abnormal specifications to be annotated with corresponding adaptations paving the way to efficiently identify and apply adaptation automatically for similar/identical patterns of abnormal behavior.

## References

- [1] T. Ball and S. Rajamani. Slam, 2003. <http://research.microsoft.com/slam>.
- [2] J. F. Bowering, J. M. Rehg, and M. J. Harrold. Active learning for automatic classification of software behavior. *SIGSOFT Softw. Eng. Notes*, 29(4), 2004.
- [3] Y. Brun and M. Ernst. Finding latent code errors via machine learning over program executions. In *26th International Conference on Software Engineering*, 2004.
- [4] C.-C. Chang and C.-J. Lin. Libsvm: a library for support vector machines (version 2.31). Available from "<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>".
- [5] H. Debar, M. Dacier, M. Nassehi, and A. Wespi. Fixed vs. variable-length patterns for detecting suspicious process behavior. In *ESORICS '98: Proceedings of the 5th European Symposium on Research in Computer Security*, 1998.
- [6] M. D. Ernst, J. Cockrell, W. G. Griswold, and D. Notkin. Dynamically discovering likely program invariants to support program evolution. In *IEEE Transactions on Software Engineering*, volume 27, 2001.
- [7] E. Eskin, W. Lee, and S. J. Stolfo. Modeling system calls for intrusion detection with dynamic window sizes. In *Proceedings of DISCEX II*, 2001.
- [8] S. Fickas and M. S. Feather. Requirements monitoring in dynamic environments. In *RE '95: Proceedings of the Second IEEE International Symposium on Requirements Engineering*, 1995.
- [9] S. Forrest. Computer immune systems, data sets. Available from "<http://www.cs.unm.edu/~immsec/data/synthesm.html>".
- [10] S. Forrest, S. A. Hofmeyr, A. Somayaji, and T. A. Longstaff. A sense of self for Unix processes. In *Proceedings of the 1996 IEEE Symposium on Research in Security and Privacy*, 1996.
- [11] P. Godefroid. Model checking for programming languages using verisoft. In *Proceedings of POPL*, 1997.
- [12] R. Gopalakrishna, E. H. Spafford, and J. Vitek. Efficient intrusion detection using automaton inlining. In *SP '05: Proceedings of the 2005 IEEE Symposium on Security and Privacy*, pages 18–31, Washington, DC, USA, 2005.
- [13] S. Hangal and M. S. Lam. Tracking down software bugs using automatic anomaly detection. In *ICSE '02: Proceedings of the 24th International Conference on Software Engineering*, 2002.
- [14] J. Hatcliff and M. Dwyer. Using the bandera tool set to model-check properties of concurrent Java software. *LNCS*, 2154, 2001.
- [15] T. Henzinger, R. Jhala, R. Majumdar, and G. Sutre. Lazy abstraction. In *Proceedings of POPL*, 2002.
- [16] A. P. Kosoresow and S. A. Hofmeyr. Intrusion detection via system call traces. 14:24–42, 1997.
- [17] A. Lazarevich, L. Ertöz, A. Ozgur, J. Srivastava, and V. Kumar. A comparative study of anomaly detection schemes in network intrusion detection. In *Proceedings of SIAM Conference on Data Mining*, 2003.
- [18] Y. Liao and V. R. Vemuri. Use of k-nearest neighbor classifier for intrusion detection. *Computers & Security*, 21(5):439–448, 2002.
- [19] Y. Liu, S. Yerramalla, E. Fuller, B. Cukic, and S. Gururajan. Adaptive control software: Can we guarantee safety? In *28th International Computer Software and Applications Conference (COMPSAC) Workshop*, pages 100–103, 2004.
- [20] C. Marceau. Characterizing the behavior of a program using multiple-length n-grams. In *Proceedings of the 2000 workshop on New security paradigms*, pages 101–110, 2000.
- [21] A. Mili, G. Jiang, B. Cukic, Y. Liu, and R. Ayyed. Towards the verification and validation of online learning systems: General framework and applications. In *HICSS*, 2004.
- [22] A. Podgurski, D. Leon, P. Francis, W. Masri, M. Minch, J. Sun, and B. Wang. Automated support for classifying software failure reports. In *ICSE '03: Proceedings of the 25th International Conference on Software Engineering*, pages 465–475, 2003.
- [23] B. Scholkopf, J. Platt, J. Shawe-Taylor, A. Smola, and R. Williamson. Estimating the support of a high-dimensional distribution. Technical Report 99-87, Microsoft Research, 1999.
- [24] R. Sekar, A. Gupta, J. Frullo, T. Shanbhag, A. Tiwari, H. Yang, and S. Zhou. Specification-based anomaly detection: a new approach for detecting network intrusions. In *CCS '02: Proceedings of the 9th ACM conference on Computer and communications security*, 2002.
- [25] N. Shankar, S. Owre, and J. M. Rushby. A tutorial on specification and verification using PVS. Technical report, 1993.

- [26] V. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.
- [27] D. Wagner and D. Dean. Intrusion detection via static analysis. In *SP '01: Proceedings of the 2001 IEEE Symposium on Security and Privacy*, 2001.