


RESEARCH

Open Access

# Detecting web attacks with end-to-end deep learning



Yao Pan<sup>1\*</sup> , Fangzhou Sun<sup>1</sup>, Zhongwei Teng<sup>1</sup>, Jules White<sup>1</sup>, Douglas C. Schmidt<sup>1</sup>, Jacob Staples<sup>2</sup> and Lee Krause<sup>2</sup>

## Abstract

Web applications are popular targets for cyber-attacks because they are network-accessible and often contain vulnerabilities. An intrusion detection system monitors web applications and issues alerts when an attack attempt is detected. Existing implementations of intrusion detection systems usually extract features from network packets or string characteristics of input that are *manually selected* as relevant to attack analysis. Manually selecting features, however, is time-consuming and requires in-depth security domain knowledge. Moreover, large amounts of *labeled* legitimate and attack request data are needed by supervised learning algorithms to classify normal and abnormal behaviors, which is often expensive and impractical to obtain for production web applications.

This paper provides three contributions to the study of autonomic intrusion detection systems. First, we evaluate the feasibility of an unsupervised/semi-supervised approach for web attack detection based on the *Robust Software Modeling Tool* (RSMT), which autonomically monitors and characterizes the runtime behavior of web applications. Second, we describe how RSMT trains a stacked denoising autoencoder to encode and reconstruct the call graph for end-to-end deep learning, where a low-dimensional representation of the raw features with unlabeled request data is used to recognize anomalies by computing the reconstruction error of the request data. Third, we analyze the results of empirically testing RSMT on both synthetic datasets and production applications with intentional vulnerabilities. Our results show that the proposed approach can efficiently and accurately detect attacks, including SQL injection, cross-site scripting, and deserialization, with minimal domain knowledge and little labeled training data.

**Keywords:** Web security, Deep learning, Application instrumentation

## 1 Introduction

**Emerging trends and challenges.** Web applications are attractive targets for cyber attackers. SQL injection [1], cross site scripting (XSS) [2] and remote code execution are common attacks that can disable web services, steal sensitive user information, and cause significant financial loss to both service providers and users. Protecting web applications from attack is hard. Even though developers and researchers have developed many counter-measures (such as firewalls, intrusion detection systems (IDSs) [3] and defensive programming best practices [4]) to protect web applications, web attacks remain a major threat.

For example, researchers found that more than half of web applications during a 2015–2016 scan contained

significant security vulnerabilities, such as XSS or SQL Injection [5]. Moreover, hacking attacks cost the average American firm \$15.4 million per year [6]. The Equifax data breach in 2017 [7, 8] (which exploited a vulnerability in Apache Struts) exposed over 143 million American consumers' sensitive personal information. Although the vulnerability was disclosed and patched in March 2017, Equifax took no action until four months later, which led to an estimated insured loss of over 125 million dollars.

Conventional intrusion detection systems do not work as well as expected for a number of reasons, including the following:

- **Workforce limitations.** In-depth domain knowledge of web security is needed for web developers and network operators to deploy these systems [9]. An experienced security expert is often needed to determine what features are relevant to extract from network packages, binaries, or other inputs for intrusion detection systems. Due to the

\*Correspondence: [panyao98@gmail.com](mailto:panyao98@gmail.com)

<sup>1</sup>Department of EECS, Vanderbilt University, Nashville, TN, USA  
Full list of author information is available at the end of the article

large demand and relatively low barrier to entry into the software profession, however, many developers lack the necessary knowledge of secure coding practices.

- **Classification limitations.** Many intrusion detection systems rely on rule-based strategies or supervised machine learning algorithms to differentiate normal requests from attack requests, which requires large amounts of labeled training data to train the learning algorithms. It is hard and expensive, however, to obtain this training data for arbitrary custom applications. In addition, labeled training data is often heavily imbalanced since attack requests for custom systems are harder to get than normal requests, which poses challenges for classifiers [10]. Moreover, although rule-based or supervised learning approaches can distinguish existing known attacks, new types of attacks and vulnerabilities emerge continuously, so they may be misclassified.

- **False positive limitations.** Although prior work has applied unsupervised learning algorithms (such as PCA [11] and SVM [12]) to detect web attacks, these approaches require manual selection of attack-specific features. Moreover, while these approaches achieve acceptable performance they also incur false positive rates that are too high in practice, e.g., a 1% increase in false positives may cause an intrusion detection system to incorrectly flag thousands of legitimate users [13]. It is therefore essential to reduce the false positive rate of these systems.

Given these challenges with using conventional intrusion detection systems, an infrastructure that requires less expertise and labeled training data is needed.

**Solution approach ⇒ Applying end-to-end deep learning to detect cyber-attacks autonomically in real-time and adapt efficiently, scalably, and securely to thwart them.** This paper explores the potential of end-to-end deep learning [14] in intrusion detection systems. Our approach applies deep learning to the entire process from feature engineering to prediction, i.e., raw input is fed into the network and high-level output is generated directly. There is thus no need for users to select features and construct large labeled training sets manually.

We empirically evaluate how well an unsupervised-/semi-supervised learning approach based on end-to-end deep learning detects web attacks. Our work is motivated by the success deep learning has achieved in computer vision [15], speech recognition [16], and natural language processing [17]. In particular, deep learning is not only capable of classification, but also automatically extracting features from high dimensional raw input.

Our deep learning approach is based on the *Robust Software Modeling Tool* (RSMT) [18], which is a late-stage (i.e., post-compilation) instrumentation-based toolchain

that target languages designed to run on the *Java Virtual Machine* (JVM). RSMT is a general-purpose tool that extracts arbitrarily fine-grained traces of program execution from running software, which is applied in this paper to detect intrusions at runtime by extracting call traces in web applications. Our approach applies RSMT in the following steps:

1. During an unsupervised training epoch, traces generated by test suites are used to learn a model of correct program execution with a stacked denoising autoencoder, which is a symmetric deep neural network trained to have target value equal to a given input value [19].

2. A small amount of labeled data is then used to calculate reconstruction error and establish a threshold to distinguish normal and abnormal behaviors.

3. During a subsequent validation epoch, traces extracted from a live application are classified using previously learned models to determine whether each trace is indicative of normal or abnormal behavior.

A key contribution of this paper is the integration of autonomic runtime behavior monitoring and characterization of web applications with end-to-end deep learning mechanisms, which generate high-level output directly from raw feature input.

This paper extends our prior work [18] by focusing on attack detection using stacked denoising autoencoders. This improved approach significantly improves upon our past approaches that relied on other machine learning techniques and typically required labeled training sets. A key benefit of the approaches presented in this paper versus our prior work is that they do not require manual feature engineering, which is needed for our past detection techniques. Moreover, the approaches work well with standard software engineering artifacts, the execution data from tests, which can be gleaned from many application-types.

The remainder of this paper is organized as follows: Section 2 summarizes the key research challenges we are addressing in our work; Section 3 describes the structure and functionality of the *Robust Software Modeling Tool* (RSMT); Section 4 explains our approach for web attack detection using unsupervised/semi-supervised end-to-end deep learning and the stacked denoising autoencoder; Section 5 empirically evaluates the performance of our RSMT-based intrusion detection system on representative web applications; Section 6 compares our work with related web attack detection techniques; and Section 7 presents concluding remarks.

## 2 Research challenges

This section describes the key research challenges we address and provides cross-references to later portions of the paper that show how we applied RSMT to detect web attacks by applying end-to-end deep learning.

**Challenge 1: Attacks can have significantly different characteristics.** Different types of web attacks, such as SQL injection, cross site scripting, remote code execution and file inclusion vulnerabilities, use different forms of attack vector and exploit different vulnerabilities inside web applications. These attacks therefore often exhibit completely different characteristics. For example, SQL injection targets databases, whereas remote code execution targets file systems. Conventional intrusion detection systems [2, 20], however, are often designed to detect only one type of attack. For instance, a grammar-based analysis that works on SQL injection detection will not work on XSS. Section 3 describes how we applied RSMT to characterize the normal behaviors and detect different types of attacks comprehensively.

**Challenge 2: Monitoring can have a significant performance cost.** Static analysis approaches that analyze source code and search for potential flaws incur various drawbacks, including vulnerability to unknown attacks and the need for source code access. An alternative is to apply dynamic analysis by instrumenting applications. Instrumentation invariably incurs monitoring overhead [21], however, which may degrade web application throughput and latency, as described in Section 5.3. Section 3.2 explores techniques applied by RSMT to minimize the overhead of monitoring and characterizing application runtime behavior.

**Challenge 3: Collecting labeled attack training data.** Machine learning-based intrusion detection systems rely on labeled training data to learn what should be considered normal and abnormal behaviors. Collecting this labeled training data can be hard and expensive in large-scale production web applications since labeling data requires extensive human effort and it is hard to cover all the possible cases. For example, normal request training data can be generated with load testing tools, web crawlers, or unit tests. If the application has vulnerabilities, however, the generated data may also contain some abnormal requests, which can undermine the performance of supervised learning approaches.

Abnormal training data is even harder to obtain [22], e.g., it is hard to know what types of vulnerabilities a system has and what attacks it will face. Even manually creating attack requests targeted for a particular application may not cover all scenarios. Moreover, different types of attacks have different characteristics, which makes it hard for supervised learning methods to capture what attack requests should look like. Although supervised learning approaches often distinguish known attacks effectively, they may miss new attacks and vulnerabilities that emerge continuously, especially when web applications frequently

depend on many third-party packages [8]. Section 4.3 describes how we applied an autoencoder-based unsupervised learning approach in RSMT to resolve the labeled training data problem.

**Challenge 4: Developing intrusion detection systems without requiring users to have extensive web security domain knowledge.** Traditional intrusion detection systems apply rule-based approach where users must have domain-specific knowledge in web security. Experienced security experts are thus needed to determine what feature(s) are relevant to extract from network packages, binaries, or other input for intrusion detection systems. This feature selection process can be tedious, error-prone, and time-consuming, such that even experienced engineers often rely on repetitive trial-and-error processes. Moreover, even web security experts may struggle to keep pace with the latest vulnerabilities due to quick technology refresh cycles and the continuous release of new tools and packages. Sections 4.1 and 4.2 describe how we applied RSMT to build intrusion detection systems with “featureless” approaches that eliminated the feature engineering step and directly used high-dimensional request traces data as input.

### 3 The structure and functionality of the robust software modeling tool (RSMT)

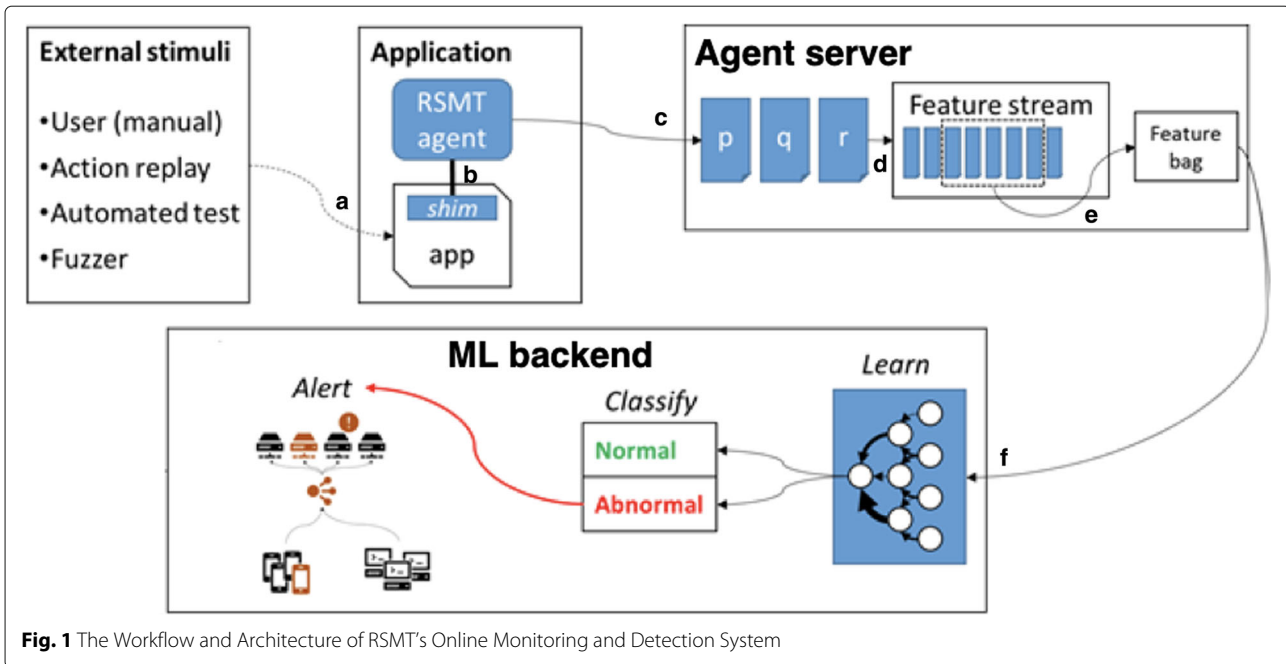
This section describes the structure and functionality of the *Robust Software Modeling Tool* (RSMT), which we developed to autonomously monitor and characterize the runtime behavior of web applications, as shown in Fig. 1.

This section first gives an overview of RSMT, then focuses on RSMT’s agent and agent server components, and finally explains how these components address *Challenge 1* (detection of different types of attacks) and *Challenge 2* (minimizing instrumentation overhead) summarized in Section 2. Section 4 later describes RSMT’s learning backend components and examines the challenges from Section 2 that they address.

#### 3.1 Overview of rSMT

As discussed in Section 2, different attacks have different characteristics and traditional feature engineering approaches lack a unified solution for all types of attacks. RSMT bypasses these attack vectors and instead captures the low-level call graph. It assumes that no matter what the attack type is (1) some methods in the server that should not be accessed are invoked and/or (2) the access pattern is statistically different than the legitimate traffic.

RSMT operates as a late-stage (post-compilation) instrumentation-based toolchain targeting languages that run on the *Java Virtual Machine* (JVM). It extracts



arbitrarily fine-grained traces of program execution from running software and constructs its models of behavior by first injecting lightweight shim instructions directly into an application binary or bytecode. These shim instructions enable the RSMT runtime to extract features representative of control and data flow from a program as it executes, but do not otherwise affect application functionality.

Figure 1 shows the high-level workflow of RSMT's web attack monitoring and detection system. This system is driven by one or more environmental stimuli (a), which are actions transcending process boundaries that can be broadly categorized as either manual (e.g., human interaction-driven) or automated (e.g., test suites and fuzzers) inputs. The manifestation of one or more stimuli results in the execution of various application behaviors. RSMT attaches an agent and embeds lightweight shims into an application (b). These shims do not affect the functionality of the software, but instead serve as probes that allow efficient examination of the inner workings of software applications. The events tracked by RSMT are typically control flow-oriented, though dataflow-based analysis is also possible.

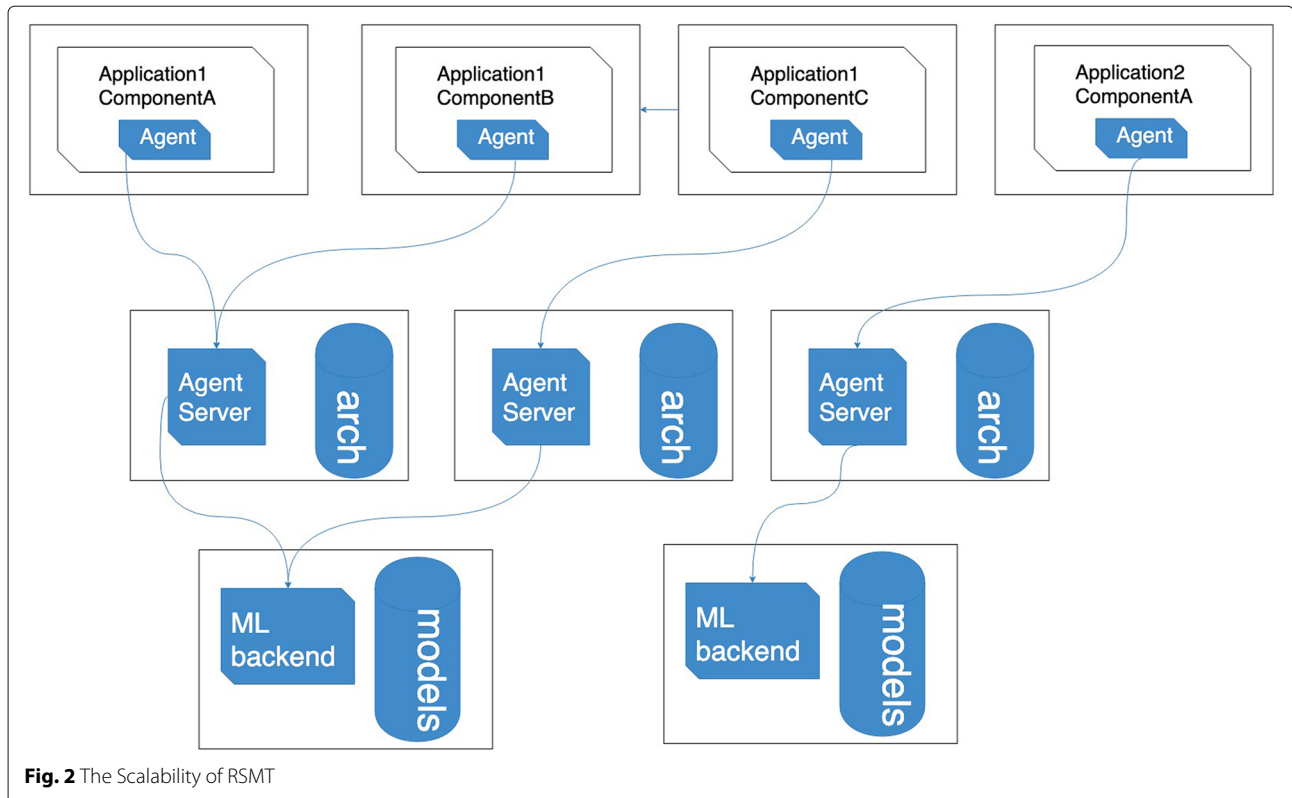
As the stimuli drive the system, the RSMT agent intercepts event notifications issued by shim instructions. These notifications are used to construct traces of behavior that are subsequently transmitted to a separate trace management process (c). This process aggregates traces over a sliding window of time (d) and converts these traces into "bags" of features (e). RSMT uses feature bags to enact online strategies (f), which involve the following two epochs:

- during a training epoch, where RSMT uses the traces generated by test suites to learn a model of correct program execution, and
- during a subsequent validation epoch, where RSMT classifies traces extracted from a live application using previously learned models to determine whether each trace is indicative of normal or abnormal behavior.

Figure 1 also shows the three core components of RSMT's architecture, which include (1) an *application*, to which the RSMT agent is attached, (2) an *agent server*, which is responsible for managing data gathered from various agents, and (3) a *machine learning backend*, which is used to train various machine learning models and validating traces. This architecture is scalable to accommodate arbitrarily large and complex applications, as shown in Fig. 2.

For example, a large web application may contain multiple components, where each component can be attached with a different agent. When the number of agents increases, a single agent server may be overwhelmed by requests from agents. Multiple agent servers can therefore be added and agent requests can then be directed to different agent servers using various partitioning rules.

It is also possible to scale the machine learning backend, e.g., by deploying machine learning training and testing engine on multiple servers. An application generally comprises multiple tasks. For example, the tasks in a web forum service might be *init*, *registerNewUser*, *createThread*, and *createPost*. Machine learning models are built at the task granularity. Different machine learning backends store and process different models.



### 3.2 The rSMT agent

**Problem to resolve.** To monitor and characterize web application runtime behavior, a plugin program is needed to instrument the web application and record necessary runtime information. This plugin program should require minimum human intervention to avoid burdening developers with low-level application behavior details. Likewise, instrumentation invariably incurs performance overhead that should be minimized to avoid unduly degrading web application throughput and latency.

**Solution approach.** To address the problem of instrumentation with minimum developer burden and performance overhead, the RSMT agent captures features that are representative of application behavior. This agent defines a class transformation system that creates events to generalize and characterize program behavior at runtime. This transformation system is plugin-based and thus extensible, e.g., it includes a range of transformation plugins providing instrumentation support for extracting timing, coarse-grained (method) control flow, fine-grained (branch) control flow, exception flow, and annotation-driven information capture.

For example, a profiling transformer can inject ultra-lightweight instructions to store the timestamps when methods are invoked. A trace transformer could add **methodEnter()** and **methodExit()** calls to construct a control flow model. Each transformation plugin conforms

to a common API. This common API can be used to determine whether the plugin can transform a given class, whether it can transform individual methods in that class, and whether it should actually perform those transformations if it is able.

We leverage RSMT's publish-subscribe (pub/sub) framework to (1) rapidly disseminate events by instrumented code and (2) subsequently capture these events via event listeners that can be registered dynamically at runtime. RSMT's pub-sub framework is exposed to instrumented bytecode via a proxy class that contains various static methods<sup>1</sup>. In turn, this proxy class calls various listeners that have been registered with it. The following event types are routed to event listeners:

- *Registration events* are typically executed once per method in each class as its `< clinit >` (class initializer) method is executed. These events are typically consumed (not propagated) by the listener proxy.
- *Control flow events* are issued just before or just after a program encounters various control flow structures. These events typically propagate through the entire listener delegation tree.
- *Annotation-driven events* are issued when annotated methods are executed. These events propagate to the offline event processing listener children.

<sup>1</sup>We use static methods since calling a Java static method is up to 2x faster than calling a Java instance method.

The root listener proxy is called directly from instrumented bytecode and delegates event notifications to an error handler, which gracefully processes exceptions generated by child nodes. Specifically, the error handler ensures that all child nodes receive a notification regardless of whether that notification results in the generation of an exception (as is the case when a model validator detects unsafe behavior). The error handler delegates to the following model construction/validation subtrees:

- The *online model construction/validation subtree* performs model construction and verification in the current thread of execution *i.e.*, on the critical path.
- The *offline model construction/validation subtree* converts events into a form that can be stored asynchronously with a (possibly remote) instance of Elasticsearch [23], which is an open-source search and analytics engine that provides a distributed real-time document store.

To address *Challenge 1* (minimizing the overhead of monitoring and charactering application runtime behavior) described in Section 2, RSMT provides a dynamic filtering mechanism. We analyzed the method call patterns and observed that most method calls are lightweight and occur in a small subset of nodes in the call graph. By identifying a method as being called frequently and having a significantly larger performance impact, we can disable events issued from it entirely or reduce the number of events it produces (thereby improving performance). These observations, along with a desire for improved performance, motivated the design of RSMT’s dynamic filtering mechanism.

To enable filtering, each method in each class is associated with a new static field added to that class during the instrumentation process. The value of the field is an object used to filter methods before they make calls to the runtime trace API. This field is initialized in the constructor and is checked just before any event would normally be issued to determine if the event should actually occur.

To characterize feature vector abilities to reflect application behaviors, we added an online model builder and

model validator to RSMT. The model builder constructs two views of software behavior: a *call graph*, which is used to quickly determine whether a transition is valid, and a *call tree*, which is used to determine whether a sequence of transitions is valid. The model validator is a closely related component that compares current system behavior to an instance of a model assumed to represent correct behavior.

When tracking calls, we can store them in a buffer that retains up to *N* past entries. The fastest and simplest tracking strategy is to use *N=1*, in which case we essentially are performing a reachability analysis. More involved is *N=2*, which yields a traditional call graph. If we do not restrict *N*, we allow the capture of arbitrarily large call histories (*i.e.*, full call stacks).

There is a tradeoff here—tracking a small, fixed-number of entries (*N=2*) can be done much faster than tracking a complete call history. However, a full call history gives a better model of program behavior.

Figure 3 demonstrates the complexity of the graphs we have created by applying RSMT on various SQL statements.

Each directed edge in a call graph connects a parent method (source) to a method called by the parent (destination). Call graph edges are not restricted with respect to forming cycles. Suppose the graph in Fig. 4 represented correct behavior. If we observed a call sequence *e,a,x* at runtime, we could easily tell this was not a valid execution path since no *a,x* edge is present in the call graph.

Although the call graph is fast and simple to construct, it has shortcomings. For example, suppose a transition sequence *e,a,d,c,a* is observed. Using the call graph, none of these transition edges violated expected behavior. If we account for past behavior, however, there is no *c,a* transition occurring after *e,a,d*. To handle these complex cases, a more robust structure is needed. This structure is known as the *call tree*, as shown in the right-hand side of Fig. 4.

Whereas the call graph falsely represents it as a valid sequence, there is no path along sequence *e,a,d,c,a* in the call tree (this requires two backtracking operations), so we determine that this behavior is incorrect. The call tree is not a tree in the structural sense. Instead, it is

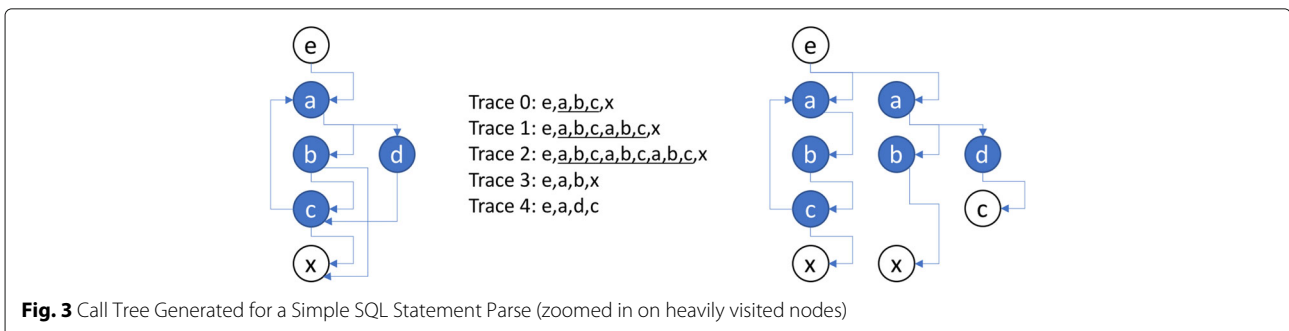
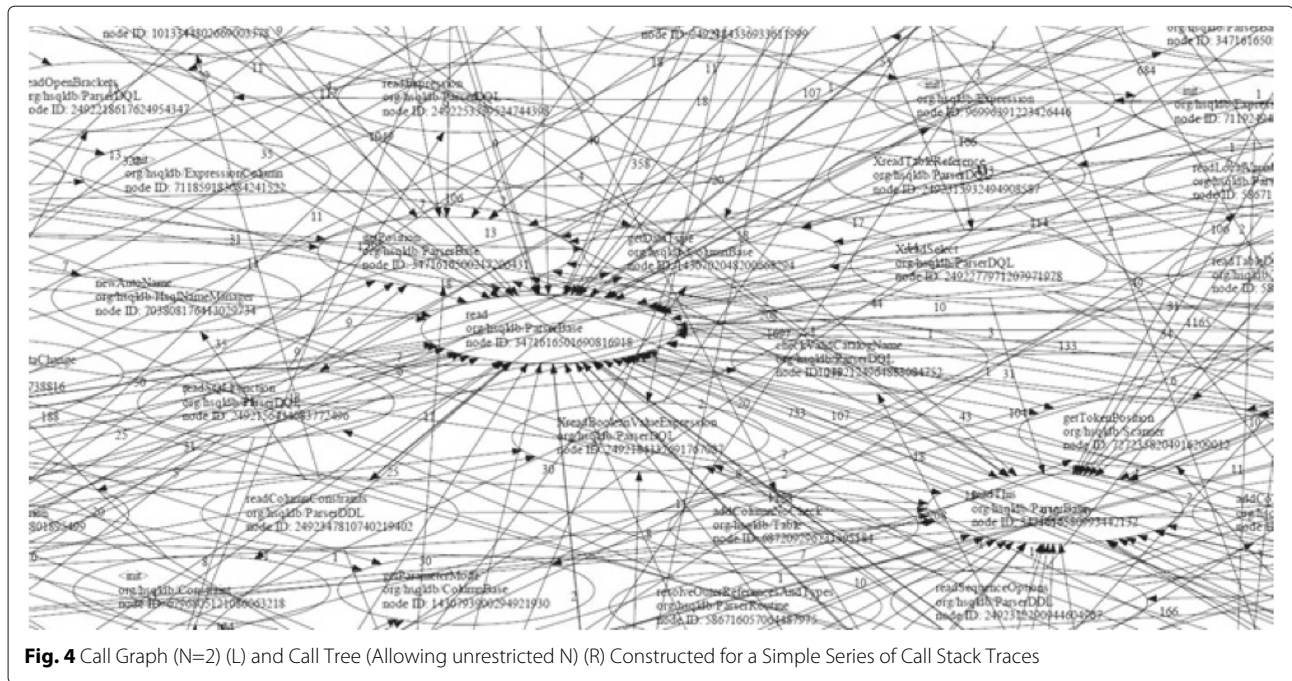


Fig. 3 Call Tree Generated for a Simple SQL Statement Parse (zoomed in on heavily visited nodes)



**Fig. 4** Call Graph (N=2) (L) and Call Tree (Allowing unrestricted N) (R) Constructed for a Simple Series of Call Stack Traces

a tree where each branch represents a possible execution path. If we follow the current execution trace to any node in the call tree, the current behavior matches the expectation.

Unlike a pure tree, the call tree does have self-referential edges (e.g., the *c,a* edge in Fig. 4) if recursion is observed. Using this structure is obviously more processing-intensive than tracking behavior using a call graph. Section 5.3 presents empirical evaluation of the performance overhead of the RSMT agent.

### 3.3 The rSMT agent server

**Problem to resolve.** A web application may comprise multiple components where multiple agents are attached. Likewise, multiple instances of the application may run on different physical hardware for scalability. It is important for agents to communicate effectively with our machine learning backend to process collected traces, which requires some means of mapping the task- and application-level abstractions onto physical computing resources.

**Solution approach.** RSMT defines an agent server component to address the problem of mapping task/application-level abstractions to physical computing resources. This component receives traces from various agents, aligns them to an application architecture, maps application components to models of behavior, and pushes the trace to the correct model in a remote machine learning system that is architecture agnostic.

The agent server exposes three different REST APIs, which are described below:

- **A trace API** that RSMT agents use to transmit execution traces. This API allows an agent to (1) register a recently launched JVM as a component in a previously defined architecture and (2) push execution trace(s).
- **An application management API** for defining and maintaining applications by (1) defining/deleting/modifying an application, (2) retrieving a list of applications, and (3) transitioning components in an application from one state to another. This design affects how traces received from monitoring agents are handled, e.g., in the *IDLE* state, traces are discarded whereas in the *TRAIN* state they are conveyed to a machine learning backend that applies them incrementally to build a model of expected behavior. In the *VALIDATE* state, traces are compared against existing models and classified as normal or abnormal.
- **A classification API** that monitors the health of applications. This API can be used to query the status of application components over a sliding window of time, whose width determines how far back in time traces are retrieved during the health check and which rolls up into a summary of all classified traces for an application's operation. This API can also be used to retrieve a JSON representation of the current health of an application.

## 4 Unsupervised web attack detection with end-to-End deep learning

This section describes how our unsupervised/semi-supervised web attack detection system augments the RSMT architectural components described in Section 3 with *end-to-end deep learning* mechanisms [16, 24], which generate high-level output directly from raw feature input. The RSMT components covered in Section 3 provide feature input for the end-to-end deep learning mechanisms described in this section, whose output indicates whether a given web request is legitimate or an attack. This capability addresses *Challenge 4* (developing intrusion detection systems without domain knowledge) summarized in Section 2.

### 4.1 Traces collection with unit tests

The RSMT agent is responsible for collecting application runtime traces, as described in Section 3.2. These collected traces include the program's execution path information, which is then used as the feature input for our end-to-end deep learning system. Below we discuss how the raw input data is represented.

When a client sends a request to a web application the RSMT agent records a *trace*, which is a histogram of directed f-calls-g edges observed beginning after the execution of a method. In particular, from a starting entry method A, we record call traces up to depth  $d$ . We record the number of times each trace triggers each method to fulfill a request from a client.

For example, A calls B one time and A calls B and B calls C one time will be represented as: A-B: 2; B-C: 1; A-B-C: 1. Each trace can be represented as a  $1*N$  vector [2,1,1] where N is the number of different method call sequences. Unlike sequence-base approaches, we do not capture every order of method call, but instead use the frequency count as features. The order information, however, is still partially preserved by recording the frequency of call sequences.

We also pad the  $1*N$  histogram feature with an additional dimension to represent un-seen method calls. If a test dataset contains method calls that never appear in the training dataset, its count will be recorded in this bit. Our goal is to determine if the request is an attack request when given the trace signature  $T_i = \{c_1, c_2, \dots, c_n\}$  produced in response to a client request  $P_i$ .

### 4.2 Anomaly detection with deep learning

Machine learning approaches for detecting web attacks can be categorized into the following two types

- **Supervised learning** approaches (such as Naive Bayes [25] and SVM [26]) work by calibrating a classifier with a training dataset that consists of data labeled as either normal traffic or attack traffic. The

classifier then classifies the incoming traffic as either normal data or an attack request. Two general types of problems arise when applying supervised approaches to detect web attacks: (1) classifiers cannot handle new types of attacks that are not included in the training dataset, as described in *Challenge 3* (hard to obtain labeled training data) in Section 2 and (2) it is hard to get a large amount of labeled training data, as described in *Challenge 3* in Section 2.

- **Unsupervised learning** approaches (such as Principal Component Analysis (PCA) [27] and autoencoder [19]) do not require labeled training datasets. Instead, they rely on the assumption that data can be embedded into a lower dimensional subspace in which normal instances and anomalies appear significantly different. The idea is to apply dimension reduction techniques (such as PCA or autoencoders) for anomaly detection. PCA or autoencoders try to learn a function  $h(X) = X$  that maps input to itself.

The input traces to web attack detection can have a very high dimension (thousands or more). If no constraint is enforced, an identity function will be learned, which is not useful. We therefore force some information loss during the process. For example, in PCA we only select a subset of eigenvalues. In autoencoder, the hidden layers will have smaller dimension than the input.

For PCA, the original input  $X$  will be projected to  $Z = XV$ .  $V$  contains the eigenvectors and we can choose  $k$  eigenvectors with the largest eigenvalues. To reconstruct the original input,  $x = XVV^T$ . If all eigenvectors are used, then  $VV^T$  is an identity matrix, no dimensionality reduction is performed, the reconstruction is perfect. If only a subset of eigenvectors are used, the reconstruction is not perfect, the reconstruction error is given by  $E = \|x - X\|^2$ .

If a test input shares similar structure or characteristics with training data, the reconstruction error should be small. To apply the same principle to web attack detection, if a test trace is similar to the ones in the training set, the reconstruction error should be small and it is likely to be a legitimate request. If the reconstruction error is large, it implies the trace is statistically different, thereby suggesting it has a higher probability of being an attack request.

### 4.3 End-to-end deep learning with stacked denoising autoencoders

The transformation performed by PCA is linear, so it cannot capture the true underlying input and output relationships if the modeled relationship is non-linear. *Deep neural networks* (DNNs) [28] have achieved success in computer vision, speech recognition, natural language processing, etc. With non-linear activation functions and



multiple hidden layers, DNNs can model complex non-linear functions.

The decision functions for anomaly detection in web attacks are often complex since no simple threshold can be used to determine if the request is an attack. Complicated interactions, such as co-occurrence and order of method calls, are all involved in the decision making. These complexities make DNNs ideal candidates for anomaly detection in web attacks. In particular, we use a special case of neural network called an *autoencoder* [19], which is a neural network with a symmetric structure.

An autoencoder consists of two parts: (1) an encoder that maps the original input to a hidden layer  $h$  with an encoder function  $h = f(x) = s(Wx + b)$ , where  $s$  is the activation function and (2) a decoder that produce a reconstruction  $r = g(h)$ . The goal of normal neural networks is to learn a function  $h(x) = y$  where the target variable  $y$  can be used for classification or regression. An autoencoder is trained to have target value equal to input value, *i.e.*, to minimize the difference between target value and input value, *e.g.*,  $L(x, g(f(x)))$  where  $L$  is the loss function. In this case, the autoencoder penalizes  $g(f(x))$  for being dissimilar from  $x$ .

If no constraint is enforced, an autoencoder will likely learn an identity function by just copying the input to the output, which is not useful. The hidden layers in autoencoders are therefore usually constrained to have smaller dimensions than the input  $x$ . This dimensionality constraint forces autoencoders to capture the underlying structure of the training data.

Figure 5 shows a visualization of normal and abnormal requests using the compressed representation learned from an autoencoder via a t-Distributed Stochastic Neighbor Embedding (t-SNE) [29]. Blue dots in this figure represent normal requests and red dots represent abnormal

requests, which can thus be easily distinguished in the low-dimensional subspace learned with the autoencoder.

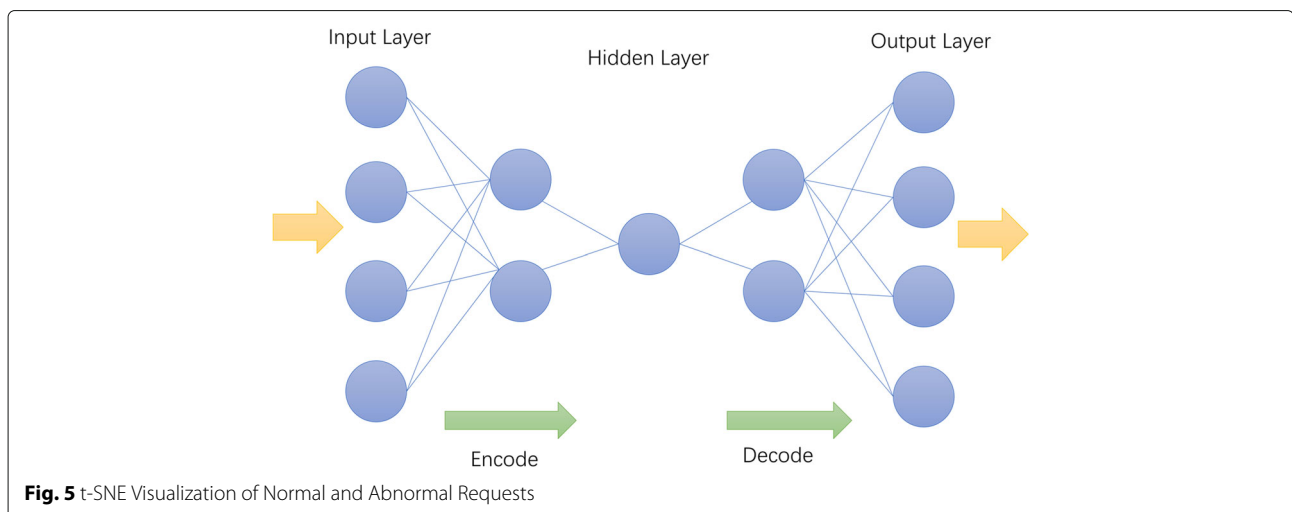
To address *Challenge 2* (detecting different types of attacks) described in Section 2, the autoencoder performs feature extraction automatically. The input  $x$  is mapped to a low dimensional representation and reconstructed trying to restore input. When the reconstruction  $g(f(x))$  is different from  $x$ , the reconstruction error  $e = ||g(f(x) - x)||^2$  can be used as an indicator for abnormality.

If the training data share similar structure or characteristics, the reconstruction error should be small. An outlier is a datum that has significantly different underlying structure or characteristic. It is therefore hard to represent the outlier with the feature we extract. As a result, the reconstruction error will be larger. We can use the reconstruction error as a standard to distinguish abnormal traffic and legitimate traffic.

Compared to PCA, autoencoders are more powerful because the encoder and decoder functions can be chosen to be non-linear, thereby capturing non-linear manifolds. In contrast, PCA just does linear transformations, so it can only create linear decision boundaries, which may not work for complex attack detection problems. Moreover, non-linearity allows the network to stack to multiple layers, which increases the modeling capacity of the network. While the combination of multiple linear transformation is still one linear layer deep, it may lack sufficient capacity to model the attack detection decision.

*Challenge 4* (developing intrusion detection systems without domain knowledge) in Section 2 is also addressed by applying the following two extensions to conventional autoencoders:

**1. Stacked autoencoders**, which may contain more than one hidden layer [19], have been applied in network security systems along with deep neural networks [30, 31]



to detect and differentiate web attacks, allowing models to learn more abstract features [32] to improve performance. Stacking increases the expressing capacity of the model, which enables the autoencoders to differentiate attacks and legitimate traffic from high dimensional input without web security domain knowledge. The output of each preceding layer is fed as the input to the successive layer. For the encoder:  $h_1 = f(x)$ ,  $h_i = f(h_{i-1})$ , whereas for the decoder:  $g_1 = g(h_i)$ ,  $g_i = g(g_{i-1})$ . Deep neural networks have shown promising applications in a variety of fields such as computer vision, natural language processing due to its representation power. These advantages also apply to deep autoencoders.

To train our stacked autoencoder we use a pretraining step involving greedy layer-wise training. The first layer of encoder is trained on raw input. After a set of parameters are obtained, this layer is used to transform the raw input to a vector represented as the hidden units in the first layer. We then train the second layer on this vector to obtain the parameters of second layers. This process is repeated by training the parameters of each layer individually, while keep the parameters of other layers unchanged.

**2. Denoising**, which prevents the autoencoder from over-fitting. Our system must be able to generalize to cases that are not presented in the training set, rather than only memorizing the training data. Otherwise, our system would not work for unknown or new types of attacks. Denoising works by corrupting the original input with some form of noise. The autoencoder now needs to reconstruct the input from a corrupted version of it, which forces the hidden layer to capture the statistical dependencies between the inputs. More detailed explanation of why denoising autoencoder works can be found in [33]. In our experiment (outlined here and described further

in Section 5) we implemented the corruption process by randomly setting 20% of the entries for each input to 0.

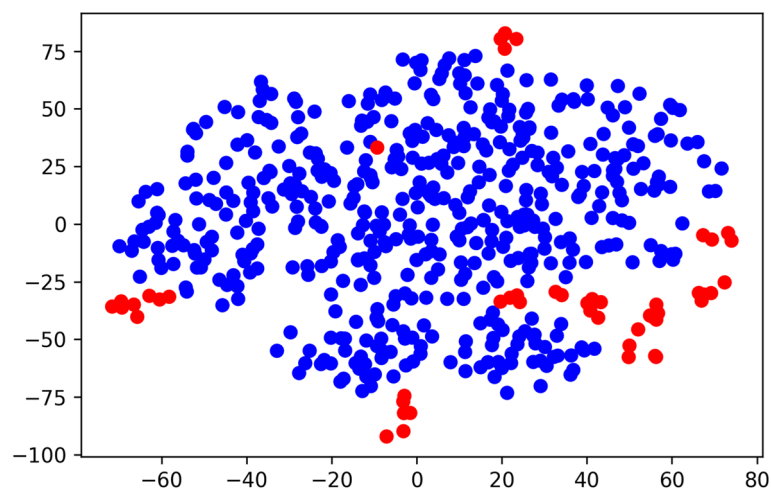
We chose a denoising autoencoder with three hidden layers for our experiments in Section 5. The structure of the autoencoder is shown in Fig. 6. The hidden layer contains  $n/2$ ,  $n/4$ ,  $n/2$  dimensions respectively. Adding more hidden layers does not improve the performance and can easily overfit. Relu [34] was chosen as the non-linear activation function in the hidden layer. Section 5.5 presents the results of experiments that evaluate the performance of a stacked denoising autoencoder in web attack detection.

The architecture of our unsupervised/semi-supervised web attack detection system is shown in Fig. 7 and described below (each numbered bullet corresponds to a numbered portion of the figure):

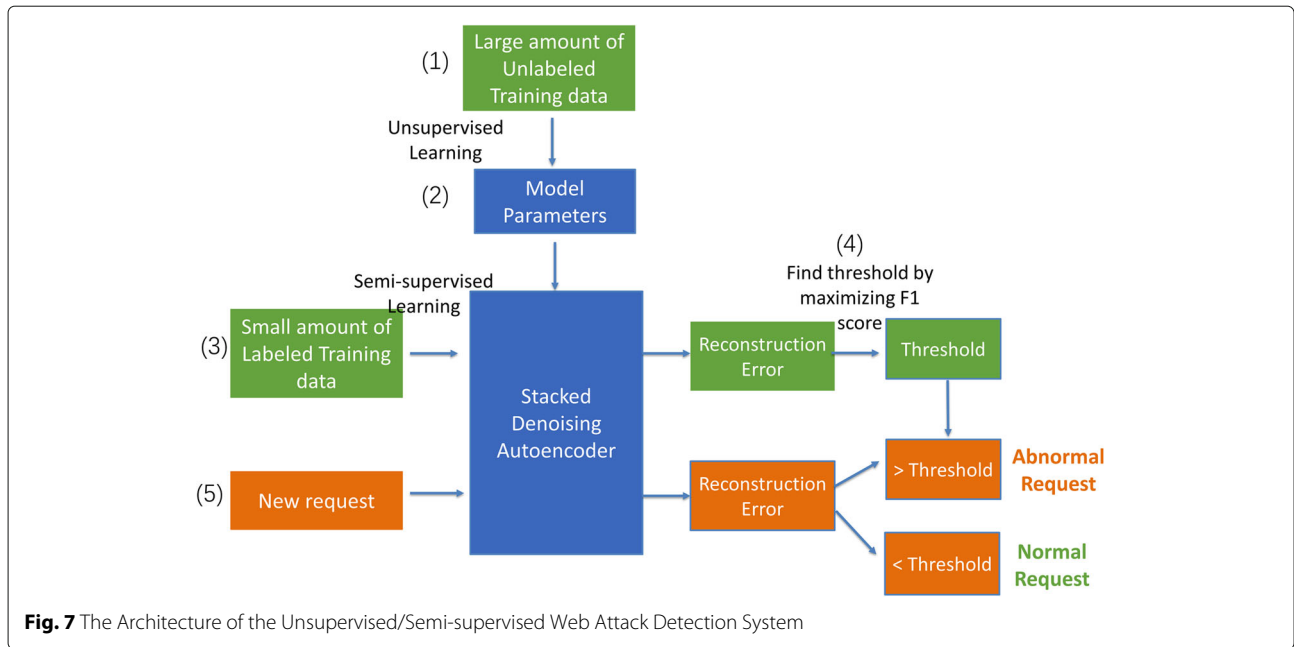
1. RSMT collected a large number of unlabeled training traces by simulating normal user requests. These unlabeled training traces should contain mostly normal requests, although a few abnormal requests may slip in.

2. A stacked denoising autoencoder is used to train on the unlabeled training traces. By minimizing the reconstruction error, the autoencoder learns an embedded low dimensional subspace that can represent the normal requests with low reconstruction error.

3. A semi-supervised learning step can optionally be performed, where a small amount of labeled normal and abnormal request data is collected. Normal request data can be collected by running repetitive unit tests or web traffic simulators, such as Apache JMeter (<http://jmeter.apache.org/>). Abnormal request data can be collected by manually creating attack requests, such as SQL injection and Cross-Site Scripting (XSS) attacks against the system. The transformation learned in unsupervised learning is applied to both normal and abnormal requests and



**Fig. 6** Structure of Stacked Autoencoder



their average reconstruction error is calculated respectively. A threshold for reconstruction error is chosen to maximize a metric, such as the F1 score, which measures the harmonic average of the precision and recall.

4. If no semi-supervised learning is conducted, the highest reconstruction error for unlabeled training data is recorded and the threshold is set to a value that is higher than this maximum by an adjustable percentage.

5. When a new test request arrived, the trained autoencoder will encode and decode the request vector and calculate reconstruction error  $E$ . If  $E$  is larger than the learned threshold  $\theta$ , it will be classified as attack request. If  $E$  is smaller than  $\theta$ , it will be considered as a normal request.

### 5 Analysis of experimental results

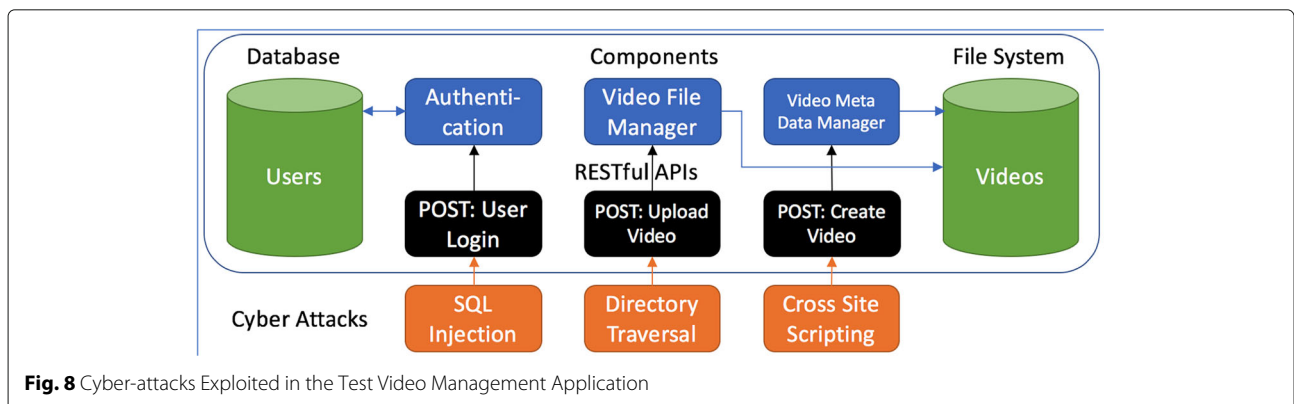
This section presents the results of experiments that empirically evaluate our deep learning-based intrusion detection system. We first describe the test environment

and evaluation metrics. We then compare the performance of our end-to-end deep learning approach with alternative methods. Our experiments were conducted using the machine learning library for Python Keras 2.0 [35], scikit-learn 0.19 [36], and Weka 3.7 [37].

#### 5.1 Testbed

We used the following two web applications as the basis for the testbed in our experiments: (1) a **video management application** built on Apache Spring framework using an embedded HSQL database and which handles HTTP requests for uploading, downloading, and viewing video files, and (2) a **compression service application** built upon the Apache Commons Compress library and which takes a file as input and outputs a compressed file in the chosen compression format.

Figure 8 shows how the test video management application provides several RESTful APIs, including: (1) *user authentication*, where a GET API allows clients to send



usernames and passwords to the server and then checks the SQL database in the back-end for authentication, (2) *video creation*, where a POST API allows clients to create or modify video metadata, and (3) *video uploading/downloading*, where POST/GET APIs allow users to upload or download videos from the server's back-end file system using the video IDs.

Our test web applications (webapps) were engineered in a manner that intentionally left them susceptible to several widely-exploited vulnerabilities. The test emulated the behavior of both normal (good) and abnormal (malicious) clients by issuing service requests directly to the test webapp's REST API. For example, the test harness might register a user with the name "Alice" to emulate a good client's behavior or "Alice 'OR true'" to emulate a malicious client attempting a SQL injection attack.

To evaluate the system's attack detection performance, we exploited three attacks from OWASP's top ten cybersecurity vulnerabilities list [38] and used them against the test webapp. These attacks included SQL injection, Cross-Site Scripting (XSS), and object deserialization vulnerabilities, as described below.

**SQL injection.** The SQL injection attack was constructed by creating queries with permutations/combinations of keywords INSERT, UPDATE, DELETE, UNION, WHERE, AND, OR, etc. The following types of SQL injections were examined:

- **Type1: Tautology-based.** Statements like OR '1' = '1' and OR '1' < '2' were added at the end of the query to make the preceding statement always true. For example, SELECT \* FROM user WHERE username = 'user1' OR '1' = '1'.
- **Type2: Comment-based.** A comment was used to ignore the succeeding statements, e.g., SELECT \* FROM user WHERE username = 'user1' # AND password = '123'.
- **Type3: Use semicolon to add additional statement,** e.g., SELECT \* FROM user WHERE username = 'user1'; DROP TABLE users; AND password = '123'.

**Cross-Site Scripting (XSS).** For the XSS attack, we added a new method with a `@RequestMapping`<sup>2</sup> in a controller that was never called in the "normal" set. We then called this method in the abnormal set to simulate an XSS attack that accessed code blocks a client should not be able to access. We also modified an existing controller method with `@RequestMapping` so a special value of one request path called a completely different code path to

execute. This alternate code path was triggered only in the abnormal set.

**Object deserialization.** Object deserialization vulnerabilities [39] can be exploited by crafting serialized objects that will invoke reflective methods that result in unsafe behaviors during the deserialization process. For example, we could store `ReflectionTransformer` items in an `ArrayList` that result in `Runtime.exec` being reflectively invoked with arguments of our choice (effectively enabling us to execute arbitrary commands at the privilege level of the JVM process). To generate such serialized objects targeting the Commons-Collections library, we used the `ysoersial` tool [40].

We collected 1000 traces for the compression service application. All runs compressed 64 MB of randomly generated data using a different method of random data generation for each run. For each of  $x \in \{1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024, 2048, 4096\}$ , a single chunk of size 64 MB/x was generated and duplicated times (with  $x = 4096$  the data is repetitive, whereas with  $x = 1$ , the data is not repetitive at all). This test shows the input dependency of compression algorithm control flow, so it was not feasible to create inputs/test cases that would exercise all possible control flow paths.

## 5.2 Evaluation metrics

An ideal intrusion detection system should classify the legitimate traffic as normal and classify attack traffic as abnormal. Two types of errors therefore exist: (1) **A false positive (FP)** or false alarm, which refers to classifying benign traffic as an attack, and (2) **A false negative (FN)**, which refers to classifying attack traffic as benign traffic. A key goal of an intrusion detection system is to minimize both the FP rate and FN rate. A tradeoff exists, however, since a more strict algorithm will tend to reduce the FN rate at the cost of classifying benign traffic as attack traffic.

Anomaly detection is an imbalanced classification problem, *i.e.*, the attack test cases appear much less frequently than the normal test cases. Accuracy is therefore not a good metric because simply predicting every request as normal will give very high accuracy. To address this issue, we use the following metrics to evaluate our approaches: (1) **Precision** =  $TP/(TP+FP)$ , which penalizes false positives, (2) **Recall** =  $TP/(TP+FN)$ , which penalizes false negatives, and (3) **F1 score** =  $2 * precision * recall / (precision + recall)$ , which evenly weights precision and recall.

## 5.3 Overhead observations

To examine the performance overhead of the RSMT agent described in Section 3, we conducted experiments that evaluated the runtime overhead in average cases and worst cases, as well as assessed how "real-time"

<sup>2</sup>`@RequestMapping` is an annotation used in Spring framework for mapping web requests onto specific handler classes or handler methods.

application execution monitoring and abnormal detection could be. As discussed in Section 3, RSMT modifies bytecode and subsequently executes it, which incurs two key sources of overhead: (1) the cost of the instrumentation itself and (2) the performance cost of executing the new instructions injected into the original bytecode.

Such instruction-level tracing can significantly increase execution time in the worst case. For example, consider a while loop that iterates 100,000 times and contains 5 instructions. If a `visitInstruction()` method call is added to each static instruction in the loop, roughly 500,000 dynamic invocations of the `visitInstruction()` method will be incurred, which is a two-fold increase in the number of dynamic instructions encountered. Moreover, this overhead can be even greater when considering the number of instructions needed to initialize fields and make the appropriate calls to `visitMethodEnter()` or handle exceptions.

RSMT has low overhead for non-computationally constrained applications. For example, a Tomcat web server that starts up in 10 seconds takes roughly 20 seconds to start up with RSMT enabled. This startup delay is introduced since RSMT examines and instruments every class loaded by the JVM. This startup cost typically occurs only once, however, since class loading usually happens just once per class.

In addition to startup delays, RSMT incurs runtime overhead every time instrumented code is invoked. We tested several web services and found RSMT had an overhead ranging from 5% to 20%. The factors most strongly impacting its overhead are the number of methods called (more frequent invocation results in higher overhead) and the ratio of computation to communication (more computation per invocation results in lower overhead).

To evaluate worst-case performance, we used RSMT to monitor the execution of an application that uses Apache's Commons Compress library to "bz2 compress" randomly-generated files of varying sizes ranging from 1x64 byte blocks to 1024x64 byte blocks, which is a control-flow intensive task. Moreover, the Apache Commons implementation of bz2 is "method heavy," (*i.e.*, there are a significant number of setter and getter calls), which are typically optimized by the JVM's hotspot compiler and converted into direct variable accesses. The instrumentation performed by RSMT prevents this optimization from occurring, however, since these lightweight methods are wrapped in calls to the model construction and validation logic. As a result, our bz2 benchmark represents the worst case for RSMT performance.

Figure 9 shows that registration adds a negligible overhead to performance (0.5 to 1%), which is expected since registration events only ever occur once per class, at class initialization. Adding call graph tracking incurs a significant performance penalty, particularly as the number of randomly generated blocks increases. Call graph tracking ranges from 1.5x to over 10x slower than the original application, whereas call tree tracking results in a 2-5x slowdown. Similarly, fine-grained control flow tracking results in a 4-6x slowdown. With full fine-grained tracking enabled, therefore, an application might run at 1% its original speed. By filtering getters and setters, however, it is possible to reduce this overhead by several orders of magnitude, as described later.

To further quantify RSMT's performance overhead, we used SPECjvm2008 [41], which is a suite comprising various integer and floating point benchmarks that quantitatively compare the performance of JVM implementations (*e.g.*, to determine whether one implementation's JIT

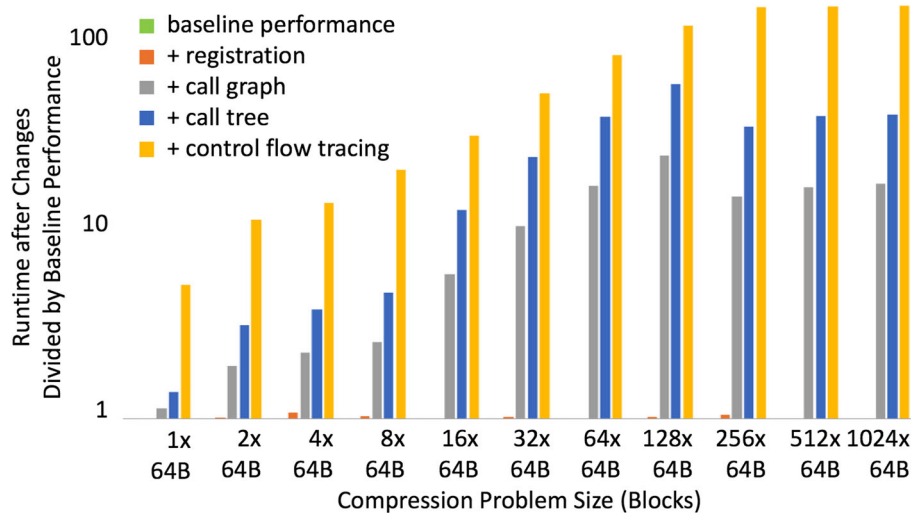


Fig. 9 Analysis of RSMT Performance Overhead

compiler is superior to another for a certain type of workload). We used the same JVM implementation across our tests, but varied the configuration of our instrumentation agents to measure the performance tradeoffs.

We evaluated the following configurations: (a) no instrumentation (no RSMT features emitted), (b) reachability instrumentation only (disabled after first access to a code region), (c) call tracing but all events passed into a null implementation, and (d) reachability + call tracing (null). We executed each configuration on a virtualized Ubuntu 14 instance provisioned with two cores and 8 GB of memory. The results of this experiment are shown below in Fig. 10. We would expect a properly tuned RSMT system to perform somewhere between configurations 3 and 4.

Although we observed that the overhead incurred by naively instrumenting all control flows within an application could be quite large (see Fig. 9), a well-configured agent should extract useful traces with overheads ranging from nearly 0% (for computation-bound applications) to 40% (for control-bound applications). Most production applications contain a blend of control-bound and computation-bound regions. Under this assumption we anticipate an overhead of 15–20% based on the composite score impact shown in Fig. 10.

#### 5.4 Supervised attack detection with manually extracted features

Before evaluating the performance of our deep learning approach, we present several supervised learning methods as benchmarks for comparison. We also describe the manually extracted features we used.

##### 5.4.1 Experiment benchmarks

Datasets and feature vectors are crucial for cyber-attack detection systems. The following feature attributes were chosen as the input for our supervised learning algorithms:

- 1 *Method execution time*. Attack behaviors can result in abnormal method execution times, e.g., SQL injection attacks may execute faster than normal database queries.
- 2 *User Principal Name (UPN)*. UPN is the name of a system user in an e-mail format, such as my\_name@my\_domain\_name. When attackers log into the test application using fake user principal names, the machine learning system can use this feature to detect it.
- 3 *The number of characters of an argument*, e.g., XSS attacks might input some abnormally large argument lengths.
- 4 *Number of domains*, which is the number of domains found in the arguments. The arguments can be inserted with malicious URLs by attackers to redirect the client “victim” to access malicious web sources.
- 5 *Duplicate special characters*. Many web browsers ignore and correct duplicated characters, so attackers can insert duplicated characters into requests to fool validators.
- 6 *N-gram*. Feature vector was built using the n-gram [42] model. The original contents of the arguments and return values are filtered by Weka’s StringToWordVector tool (which converts plain word into a set of attributes representing word occurrence) and the results are then applied to make the feature vectors.

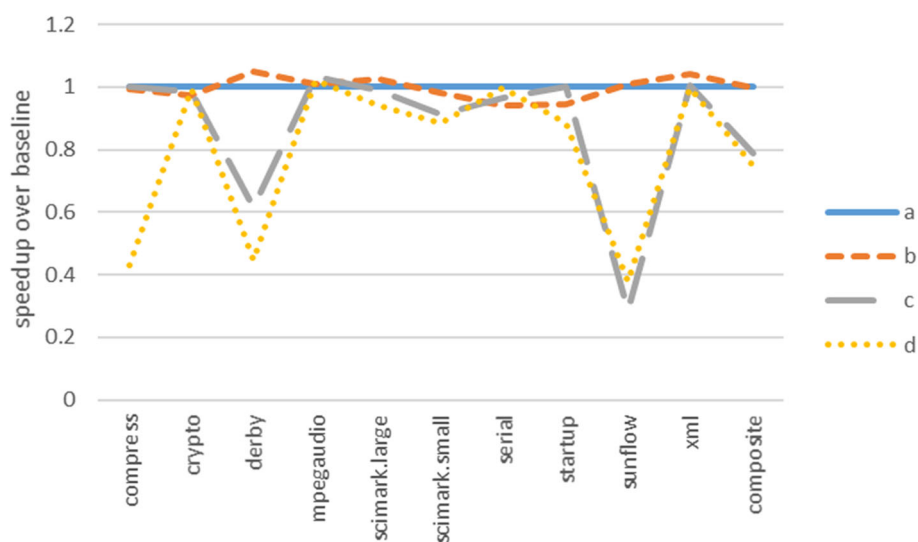


Fig. 10 SPECjvm2008 Performance Impact for Various Benchmarks and Test Configurations

After instrumenting the runtime system to generate measurements of the system when it is executing correctly or incorrectly, supervised approaches use these measurements to build a training data set. In this data set the measurements are viewed as features that can characterize the correct and incorrect system operation. Machine learning algorithms use these features to derive models that classify the correctness of the execution state of the system based on a set of measurements of its execution. When new execution measurements are given to the machine-learned model, algorithms can be applied to predict whether the previously unseen trace represents a valid execution of the system.

To provide an environment for classification, regression, and clustering we used the following three supervised machine learning algorithms from the Weka workbench:

- *Naive Bayes*, whose classification decisions calculate the probabilities/costs for each decision and are widely used in cyber-attack detection [43].
- *Random forests*, which is an ensemble learning method for classification that train decision trees on sub-samples of the dataset and then improve classification accuracy via averaging. A key parameter for random forest is the number of attributes to consider in each split point, which are selected automatically by Weka.
- *Support vector machine (SVM)*, which is an efficient supervised learning model that draws an optimal hyperplane in the feature space and divides separate categories as widely as possible. RSMT uses Weka's *Sequential Minimal Optimization* algorithm to train the SVM.

Likewise, to reduce variance and avoid overfitting [44], we also used the following two aggregate models:

- *Aggregate\_vote*, which returns ATTACK if a majority of classifiers detect attacks and returns NOT\_ATTACK otherwise.
- *Aggregate\_any*, which returns attack if any classifier detects attacks and NOT\_ATTACK otherwise.

#### 5.4.2 Experiment results

Tables 1 and 2 show the performance comparison of different machine learning algorithms on testbed web applications. For the SQL injection attacks, the training dataset contains 160 safe unit tests and 80 attack unit tests, while the test dataset contains 40 safe unit tests and 20 attack unit tests. The SQL injection attack samples bypass the test application's user authentication and include the most common SQL injection attack types.

There are three different machine learning models (Naive Bayes, Random Forest and SVM) along with two derived approaches (AGGREGATE VOTE and

**Table 1** Machine Learning Models' Experimental Results for SQL Injection Attacks

	Precision	Recall	F-score
Naive bayes	0.941	0.800	0.865
Random forest	1.000	0.800	0.889
SVM	0.933	0.800	0.889
AGGREGATE_VOTE	1.000	0.800	0.889
AGGREGATE_ANY	0.941	0.800	0.865

AGGREGATE ANY). For the SQL injection results, all three models misclassified 4 attack queries as benign with the remaining 16 samples as true positives. Since these 3 models all misclassified the 4 attack vectors, they have the same recall. The derived approaches also have the same results. Additionally, considering there are only 20 positive samples in the testing dataset, additional test data will be needed in future work to observe differences between these machine learning approaches.

The XSS training dataset contains 1000 safe unit tests and 500 attack unit tests, while the test dataset contains 150 safe unit tests and 75 attack unit tests (XSS attack samples were obtained from [www.xssed.com](http://www.xssed.com)). All three classifiers are similar in detecting XSS attacks.

## 5.5 Unsupervised attack detection with deep learning

### 5.5.1 Experiment benchmarks

Several techniques can be applied to differentiate benign traffic and attack traffic. The first is the naive approach, which learns a set of method calls from a training set (obtained by unit test or simulated legitimate requests). When a new trace is encountered, the naive approach checks if the trace contain any method call that is never seen from the training set. If there is such method, the trace will be treated as attack trace, otherwise it is considered safe.

The naive approach can detect attack traces easily since attack traces usually contains some dangerous method calls that will not be used in legitimate operation. The naive approach, however, also suffer from high false positive rate since it may not be possible to iterate through all the legitimate request scenarios. A legitimate request may

**Table 2** Machine Learning Models' Experimental Results for Cross-site Scripting Attacks

	Precision	Recall	F-score
Naive bayes	0.721	1.000	0.838
Random forest	0.721	1.000	0.838
SVM	0.728	1.000	0.843
AGGREGATE_VOTE	0.724	1.000	0.840
AGGREGATE_ANY	0.710	1.000	0.831

thus contain some method call(s) that do not exist in the training set, which results in blocking benign traffic.

A more advanced technique is one-class SVM [45]. Traditional SVM solves the two or multi-class situation. While the goal of a one-class SVM is to test new data and find out whether it is similar to the training data or not. By just providing the normal training data, one-class classification creates a representational model of this data. If newly encountered data is too different (e.g., outliers in the projected high-dimensional space), it is labeled as out-of-class.

### 5.5.2 Experiment results

Tables 3 and 4 compare the performance of different machine learning algorithms on our two testbed web applications. For the video upload application, the attack threat is SQL injection and XSS. The results in these tables show that autoencoder outperforms the other algorithms. For the compression application, we evaluate the detection performance in terms of a deserialization attack.

The naive approach, not to be confused with "naïve bayes," can detect attack traces in some circumstances since attack traces often contain some unusual method calls that are never seen in legitimate operation. The naive approach, however, also suffers from a high false positive rate since it may not be possible to iterate through all the legitimate request scenarios. A legitimate request may thus contain some method call(s) that do not exist in the training set, which results in flagging benign traffic. This is why the naive approach has high recall but low precision.

Figure 11 plots the precision/recall/F-score curve along with threshold value. This figure shows a tradeoff between precision and recall. If a threshold is chosen that is too low, many normal request will be classified as abnormal, resulting in higher false negative and low recall score. In contrast, if a threshold is chosen that is too high, many abnormal requests will be classified as normal, leading to higher false positive and low precision score. To balance precision and recall in our experiments, we choose a threshold that maximizes the F-score in the labeled training data.

To understand how various parameters (such as training data size, input feature dimension, and test coverage ratio) affect the performance of machine learning algorithms,

**Table 3** Performance Comparison of Different Machine Learning Algorithms on Video Management Application

	Precision	Recall	F-score
Naive	0.722	0.985	0.831
PCA	0.827	0.926	0.874
One-class SVM	0.809	0.909	0.858
Autoencoder	0.898	0.942	0.914

**Table 4** Performance Comparison of Different Machine Learning Algorithms on Compression Application

	Precision	Recall	F-score
Naive	0.421	1.000	0.596
PCA	0.737	0.856	0.796
One-class SVM	0.669	0.740	0.702
Autoencoder	0.906	0.928	0.918

we manually created a synthetic dataset to simulate web application requests.

Figure 12 compares the performance of machine learning algorithms with different unlabeled training data sizes. Since the test case contains method calls that were not presented in the training data, the naive approach simply treats every request as abnormal, resulting 100% recall, but 0% precision. Both PCA and autoencoder's performance improves since we have more training data.

PCA performs better, however, when there is limited training data (below 1000). The autoencoder needs more training data to converge, but outperforms the other machine learning algorithms after it is given enough training data. Our results show the autoencoder generally needs 5000 unlabeled training data to achieve good performance.

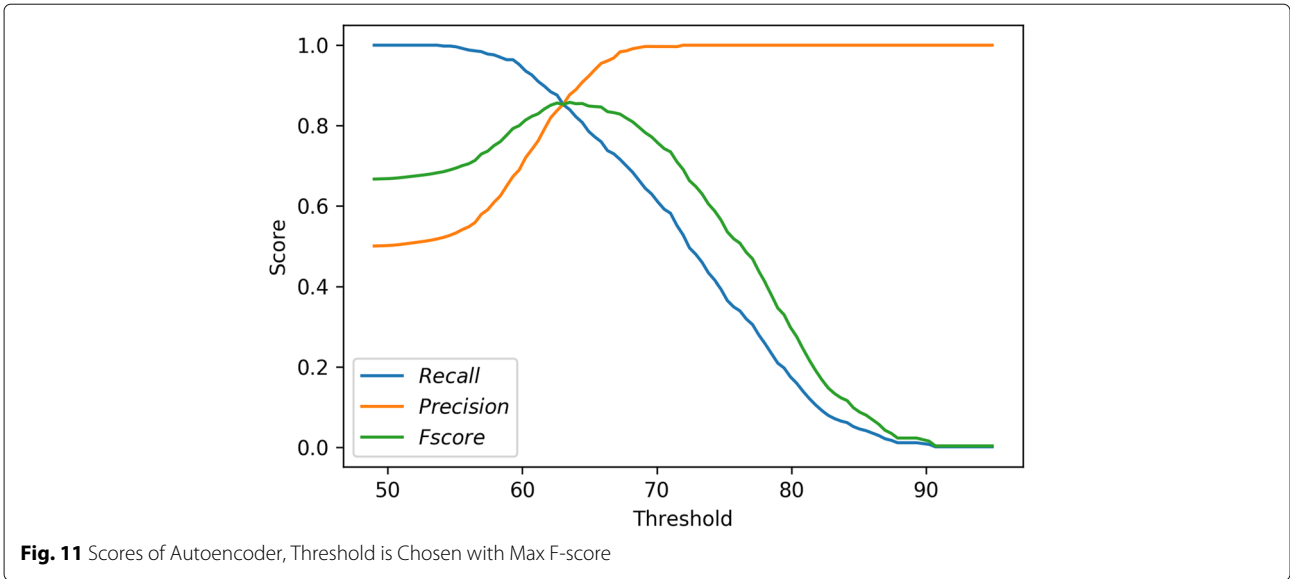
Figure 13 shows the performance of machine learning algorithms under different test coverage ratios. The test coverage ratio is the percentage of method calls covered in the training dataset. For large-scale web applications, it is impossible to traverse every execution path and method calls due to the "path explosion problem" [46], where the number of feasible paths in a program grows exponentially with an increase in program size.

If only a subset of method calls are present in the training dataset, the naive approach or other supervised learning approaches may classify the legitimate test request with uncovered method calls as abnormal. In contrast, PCA and autoencoder algorithms can still learn a hidden manifold by finding the similarity in structure instead of exact method calls. They can thus perform well even given only a subset of coverage for all the method calls.

Figure 14 shows the performance of machine learning algorithms under different input feature dimensions (the unique feature ratio is kept constant). This figure shows the gap between autoencoder and other ML techniques increases as the number of features increases. As the number of feature increases, however, this gap becomes larger. The autoencoder shows robust performance even with complicated high dimension input data.

Figure 15 compares the performance of machine learning algorithms under different unique feature ratios. This figure shows that the performance of the machine learning algorithms improves as the unique feature ratio increases.





This result is not surprising because the statistical difference between normal and abnormal requests is larger and easier to capture. For the autoencoder algorithm at least 2% of unique features are needed in the abnormal requests for acceptable performance.

In our synthetic dataset, the ratio of trace features existing in training but not test is constant. For the naive approach, the precision and recall remain constant regardless of the number of trace features.

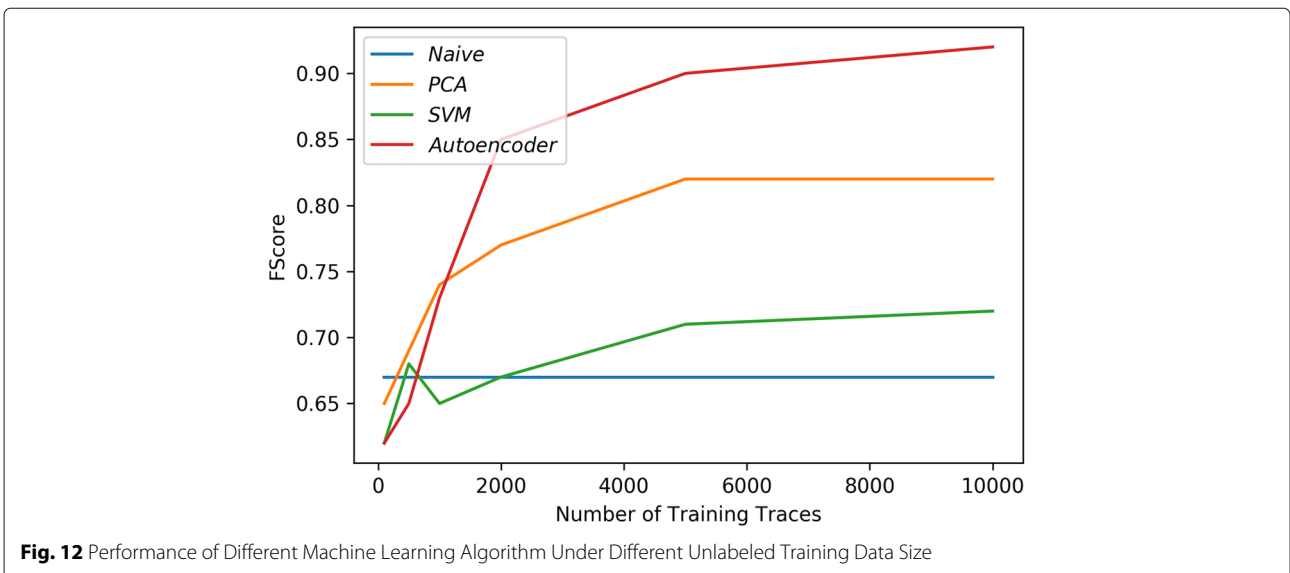
The experiment was conducted on a desktop with Intel i5 3570 and GTX 960 GPU running Windows 10. The autoencoder was implemented using Keras 2.0 with a TensorFlow backend.

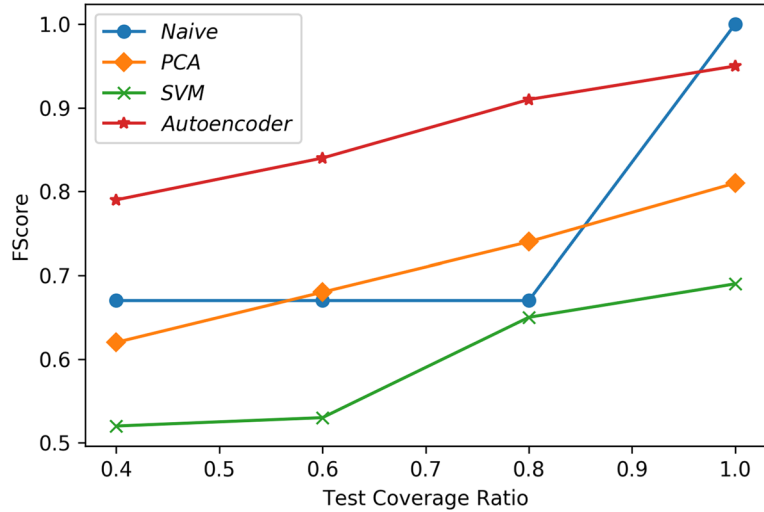
Table 5 compares the training/classification time for different algorithms. The training was performed with the

same set of 5000 traces with default parameters specified in Section 4.3. The classification time is the average time to classify one trace over 1000 test traces.

The results in Table 5 show that the training time of the deep autoencoder is significant longer than other approaches. This training need not be performed frequently, however, and can also be done offline. Moreover, existing deep learning frameworks (such as TensorFlow) support powerful GPUs, which can also significantly accelerate training time.

For the classification time, all machine learning algorithms can perform classification in an acceptable short period of time with the trained model. Moreover, hardware advances (such as the Tensor Processing Unit [47]) are bringing high performance and low cost computing





**Fig. 13** Performance of Different Machine Learning Algorithms Under Different Test Coverage Ratios

resources in the future. Computation cost should thus not be a bottleneck for future deployments of deep learning to detect web attacks.

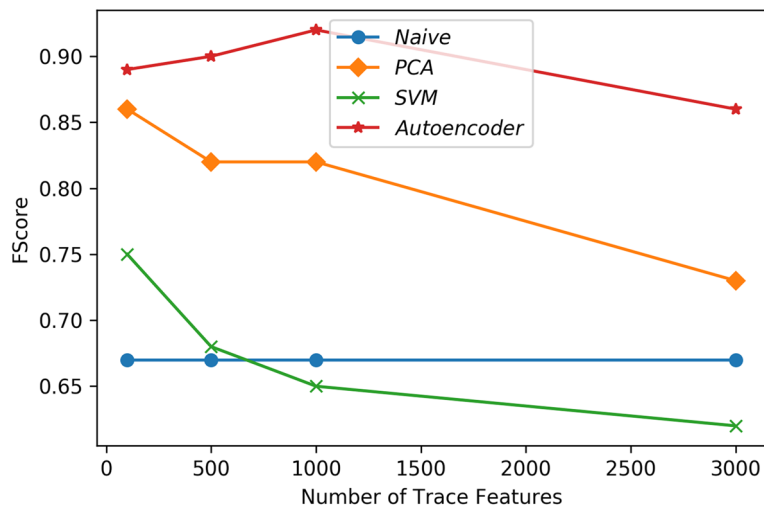
### 6 Related work

Intrusion detection systems monitor a network and/or system for malicious activity or policy violations [3]. These types of systems have been studied extensively in the literature based on various approaches, including static analysis [2, 20], sequence-based [48, 49], manual modeling [50, 51], and machine learning [52, 53]. This section describes prior work and compares/contrasts it to our research on RSMT presented in this paper.

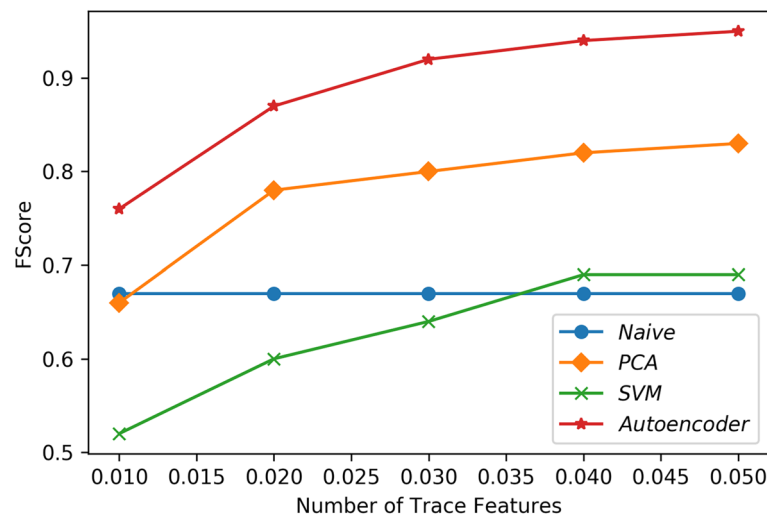
### 6.1 Static analysis

Static analysis approaches examine an application’s source code and search for potential flaws in its construction and expected execution that could lead to attack. For example, Fu et al. [20] statically analyzed SQL queries and built grammars representing expected parameterization. Wassermann et al. [2] presented a static analysis for detecting XSS vulnerabilities using tainted information flow with string analysis. Kolosnjaji et al. [54] proposed an analysis on system call sequences for malware classification.

Statically derived models can also be used at runtime to detect parameterizations of the SQL queries that do not



**Fig. 14** Performance of Different Machine Learning Algorithms Under Different Input Feature Dimensions



**Fig. 15** Performance of Different Machine Learning Algorithm Under Different Unique Feature Ratios

fit the grammar and indicate possible attack. Static analysis approaches, however, typically focus on specific types of attacks that are known a priori. In contrast, RSMT bypasses these various attack vectors and captures the low-level call graph under the assumption that the access pattern of attack requests will be statistically different than legitimate requests, as shown in Section 3.

Moreover, many static analysis techniques require access to application source code, which may not be available for many production systems. Employing attack-specific detection approaches requires building a corpus of known attacks and combining detection techniques to secure an application. A significant drawback of this approach, however, is that it does not protect against unknown attacks for which no detection techniques have been defined. In contrast, RSMT models correct program execution behaviors and uses these models to detect abnormality, which works even if attacks are unknown, as shown in Section 4.

## 6.2 Sequence-based

Sequence-based anomaly detection approaches [48, 49] either try to model the call sequences as a finite-state automaton (FSA), Hidden Markov Models (HMM) [55]

or N-gram [56]. FSA can capture common program structures such as loops or branches and predict future behaviors from past behaviors. Although sequence-based approaches achieved early success, their time complexity is high.

It has also been shown that there is no polynomial time algorithm for learning an optimal FSA [57]. N-gram [56] breaks a system call sequence into subsequences of fixed length  $N$ . The limitation of N-gram, however, is the number of N-gram grow exponentially with  $N$ .  $N$  must therefore be small, though a small  $N$  makes the algorithm ineffective at capturing long-term correlations. Moreover, the false alarm rate is high for N-gram because it cannot generalize to any N-gram that are not present in the training dataset.

## 6.3 Manual modeling

Manual modeling relies on designers to annotate code or build auxiliary textual or graphical models to describe expected system behavior. For example, SysML [50] is a language that allows users to define parametric constraint relationships between different parameters of the system to indicate how changes in one parameter should propagate or affect other parameters. Scott [51] proposed a Bayesian model-based design for intrusion detection systems. Ilgun et al. [58] used state transitions to model the intrusion process and build a rule-based intrusion detection system.

Manual modeling is highly effective when analysis can be performed on models to simulate or verify that error states are not reached. Although expert modelers can manually make errors, many errors can be detected via model simulation and verification. A key challenge of using manual modeling alone for detecting

**Table 5** Comparison of Training/Classification Time for Different Algorithms

	Training Time	Classification Time
Naive	51s	0.05s
PCA	2min 12s	0.2s
One-class SVM	2min 6s	0.2s
Autoencoder	8min 24s	0.4s

cyber-attacks, however, is that it may not fully express or capture all characteristics needed to identify the attacks. Since manual models typically use abstractions to simplify their usage and specification of system properties, these abstractions may not provide sufficient expressiveness to describe properties needed to detect unknown cyber-attacks. Our deep learning approach uses RSMT to analyze raw request trace data without making any assumption of the relationships or constraints of the system, thereby overcoming limitations with manual modeling, as shown in Section 3.

#### 6.4 Web application firewalls

Web Application Firewalls [59, 60] are a related approach for detecting and thwarting attacks that are complementary to the proposed approach. With web application firewalls, the firewall itself looks for abnormal interactions with the application that should be filtered or blocked. The proposed work demonstrates potential feasibility of denoising autoencoders to learn expected application behavior and identify attacks. The proposed approaches could be used in concert with existing web application firewall techniques.

#### 6.5 Machine learning

Machine learning approaches require instrumenting a running system to measure various properties (such as execution time, resource consumption, and input characteristics) to determine when the system is executing correctly or incorrectly due to cyber-attacks, implementation bugs, or performance bottlenecks. For example, Farid et al. [52] proposed an adaptive intrusion detection system by combining naive bayes and decision tree. Zolotukhin et al. [53] analyzed HTTP request with PCA, SVDD, and DBSCAN for unsupervised anomaly detection. Likewise, Shar et al. [61] used random forest and co-forest on hybrid program features to predict web application vulnerabilities.

Anomaly detection is another machine learning [22] application that addresses cases where traditional classification algorithms work poorly, such as when labeled training data is imbalanced. Common anomaly detection algorithms include mixture Gaussian models, support vector machines, and cluster-based models [62]. Likewise, autoencoder techniques have shown promising results in many anomaly detection tasks [63–65].

Our RSMT-baed approach described in this paper uses a stacked autoencoder to build an end-to-end deep learning system for the intrusion detection domain. The accuracy of conventional machine learning algorithms [52, 61] rely heavily on the quality of manually selected features, as well as the labeled training data. In contrast, our deep learning approach uses RSMT to extract features from high-dimensional raw input automatically without relying

on domain knowledge, which enables it to achieve better detection accuracy with large training data, as shown in Section 5.5.

## 7 Concluding remarks

This paper describes the architecture and results of applying a unsupervised end-to-end deep learning approach to automatically detect attacks on web applications. We instrumented and analyzed web applications using the Robust Software Modeling Tool (RSMT), which autonomically monitors and characterizes the runtime behavior of web applications. We then applied a denoising autoencoder to learn a low-dimensional representation of the call traces extracted from application runtime. To validate our intrusion detection system, we created several test applications and synthetic trace datasets and then evaluated the performance of unsupervised learning against these datasets.

While cross validation is widely used in traditional machine learning, it is often not used for evaluating deep learning models because of the great computational cost. We needed to compare autoencoder approaches with other machine learning methods. To enable a fair comparison, we didn't use cross validation in our experiments.

The following are key lessons learned from the work presented in this paper:

- **Autoencoders can learn descriptive representations from web application stack trace data.** Normal and anomalous requests are significantly different in terms of reconstruction error with representations learned by autoencoders. The learned representation reveals important features, but shields application developers from irrelevant details. The results of our experiments in Section 5.5 suggest the representation learned by our autoencoder is sufficiently descriptive to distinguish web request call traces.

- **Unsupervised deep learning can achieve over 0.91 F1-score in web attack detection without using domain knowledge.** By modeling the correct behavior of the web applications, unsupervised deep learning can detect different types of attacks, including SQL injection, XSS or deserialization with high precision and recall. Moreover, less expertise and effort is needed since the training requires minimum domain knowledge and labeled training data.

- **End-to-end deep learning can be applied to detect web attacks.** The accuracy of the end-to-end deep learning can usually outperform systems built with specific human knowledge. The results of our experiments in Section 5.5 suggest end-to-end deep learning can be successfully applied to detect web attacks. The end-to-end deep learning approach using autoencoders achieves better performance than supervised methods in web attack

detection without requiring any application-specific prior knowledge.

In future work, we plan to investigate more complex network structures such as LSTM autoencoders or autoencoders with CNNs. We would like to see whether these structures can provide better accuracy for web attack detection tasks. Although the results show that the autoencoder outperforms other approaches, there is still significant research needed to be done to show performance on zero-day attacks. The results show promise that autoencoders will be able to potentially detect zero-day attacks, but more research in this area is still needed. A fundamental challenge of this work will be assessing efficacy of autoencoders against attacks that are unknown.

Determining the frequency that models should be retrained is also an open research question that will need to be analyzed in future work. Retraining using online data from real-world usage opens the possibility of incorporating attack data into the normal behavior data set. It is possible that the overwhelming valid usage of the application will outweigh this issue and lead to correct detection, but this hypothesis requires additional exploration. Also, we plan to develop mechanisms to distribute the machine learning analysis workload across remote machines and support coordinated distributed detection across hosts.

#### Abbreviations

DNNS: Deep neural networks; FSA: Finite-state automaton; HMM: Hidden Markov models; JVM: Java virtual machine; LSTM: Long short-term memory; PCA: Principle component analysis; RSMT: Robust software modeling tool; SVM: Support vector machine; XSS: Cross site scripting

#### Acknowledgements

We would like to thank the anonymous reviewers for their careful reading of our manuscript and their many insightful comments and suggestions.

#### Authors' contributions

YP conducted experiments on unsupervised web attack detection and implemented stacked autoencoder. FS conducted experiments on supervised attack detection. ZT helped conduct experiments and revised the manuscript. JW participated in the design of the experiments and provided critical revision to the manuscript. DS provided critical revision to the manuscript. JS implemented RSMT and contributed to Section 3. LK contributed to Section 3. All authors read and approved the final manuscript.

#### Funding

Not applicable.

#### Availability of data and materials

The datasets used and/or analysed during the current study are available from the corresponding author on reasonable request.

#### Competing interests

The authors declare that they have no competing interests.

#### Author details

<sup>1</sup>Department of EECS, Vanderbilt University, Nashville, TN, USA. <sup>2</sup>Securboracion Inc., Melbourne, FL, USA.

Received: 18 February 2019 Accepted: 25 July 2019

Published online: 27 August 2019

#### References

- Halfond WG, Viegas J, Orso A. A classification of sql-injection attacks and countermeasures. In: Proceedings of the IEEE International Symposium on Secure Software Engineering. IEEE; 2006. p. 13–5.
- Wassermann G, Su Z. Static detection of cross-site scripting vulnerabilities. In: Proceedings of the 30th International Conference on Software Engineering. ACM; 2008. p. 171–80.
- Di Pietro R, Mancini LV. Intrusion Detection Systems vol. 38: Springer; 2008.
- Qie X, Pang R, Peterson L. Defensive programming: Using an annotation toolkit to build dos-resistant software. *ACM SIGOPS Oper Syst Rev.* 2002;36(5):45–60.
- <https://doi.org/https://www.acunetix.com/acunetix-web-application-vulnerability-report-2016>. Accessed 16 Aug 2017.
- <https://doi.org/http://money.cnn.com/2015/10/08/technology/cybercrime-cost-business/index.html>. Accessed 16 Aug 2017.
- <https://doi.org/https://www.consumer.ftc.gov/blog/2017/09/equifax-data-breach-what-do>. Accessed 16-August-2017.
- <https://doi.org/https://theconversation.com/why-dont-big-companies-keep-their-computer-systems-up-to-date-84250>. Accessed 16 Aug 2017.
- Ben-Asher N, Gonzalez C. Effects of cyber security knowledge on attack detection. *Comput Hum Behav.* 2015;48:51–61.
- Japkowicz N, Stephen S. The class imbalance problem: A systematic study. *Intell Data Anal.* 2002;6(5):429–49.
- Liu G, Yi Z, Yang S. A hierarchical intrusion detection model based on the pca neural networks. *Neurocomputing.* 2007;70(7):1561–8.
- Xu X, Wang X. An adaptive network intrusion detection method based on pca and support vector machines. *Advanced Data Mining and Applications.* 2005;3584:696–703.
- Pietraszek T. Using adaptive alert classification to reduce false positives in intrusion detection. In: *Recent Advances in Intrusion Detection*. Springer; 2004. p. 102–24.
- Goodfellow I, Bengio Y, Courville A. *Deep Learning*: MIT press; 2016.
- Krizhevsky A, Sutskever I, Hinton GE. Imagenet classification with deep convolutional neural networks. In: *Advances in Neural Information Processing Systems*. Curran Associates, Inc.; 2012. p. 1097–105.
- Amodei D, Ananthanarayanan S, Anubhai R, Bai J, Battenberg E, Case C, Casper J, Catanzaro B, Cheng Q, Chen G, et al. Deep speech 2: End-to-end speech recognition in english and mandarin. In: *International Conference on Machine Learning*. New York: PMLR; 2016. p. 173–82.
- Sutskever I, Vinyals O, Le QV. Sequence to sequence learning with neural networks. In: *Advances in Neural Information Processing Systems*. Curran Associates, Inc.; 2014. p. 3104–12.
- Sun F, Zhang P, White J, Schmidt D, Staples J, Krause L. A feasibility study of autonomously detecting in-process cyber-attacks. In: *Cybernetics (CYBCON), 2017 3rd IEEE International Conference On*. IEEE; 2017. p. 1–8.
- Vincent P, Larochelle H, Lajoie I, Bengio Y, Manzagol P-A. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *J Mach Learn Res.* 2010;11(Dec): 3371–408.
- Fu X, Lu X, Peltsverger B, Chen S, Qian K, Tao L. A static analysis framework for detecting sql injection vulnerabilities. In: *Computer Software and Applications Conference, 2007. COMPSAC 2007. 31st Annual International*. IEEE; 2007. p. 87–96.
- Waddington DG, Roy N, Schmidt DC. Dynamic analysis and profiling of multi-threaded systems. *IGI Glob.* 2009;156–99.
- Chandola V, Banerjee A, Kumar V. Anomaly detection: A survey. *ACM Comput Surv (CSUR).* 2009;41(3):15.
- Elasticsearch. <https://www.elastic.co/products/elasticsearch>. Accessed 13 Aug 2019.
- Graves A, Jaitly N. Towards end-to-end speech recognition with recurrent neural networks. In: *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*. JMLR.org; 2014. p. 1764–72.
- Russell S, Norvig P. *Intelligence A. A modern approach*. Artif Intell Prentice-Hall, Egnlewood Cliffs. 1995;25:27.
- Cortes C, Vapnik V. Support vector machine. *Mach Learn.* 1995;20(3): 273–97.
- Wold S, Esbensen K, Geladi P. Principal component analysis. *Chemometr Intell Lab Syst.* 1987;2(1–3):37–52.
- LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature.* 2015;521(7553): 436–44.

29. Maaten Lvd, Hinton G. Visualizing data using t-sne. *J Mach Learn Res*. 2008;9(Nov):2579–605.
30. Ma T, Wang F, Cheng J, Yu Y, Chen X. A hybrid spectral clustering and deep neural network ensemble algorithm for intrusion detection in sensor networks. *Sensors*. 2016;16(10):1701.
31. Vartouni AM, Kashi SS, Teshnehlab M. An anomaly detection method to detect web attacks using stacked auto-encoder. In: 2018 6th Iranian Joint Congress on Fuzzy and Intelligent Systems (CFIS). IEEE; 2018. p. 131–4.
32. Yadav S, Subramanian S. Detection of application layer ddos attack by feature learning using stacked autoencoders. In: 2016 International Conference on Computational Techniques in Information and Communication Technologies (ICCTICT). IEEE; 2016. p. 361–6.
33. Vincent P, Larochelle H, Bengio Y, Manzagol P.-A. Extracting and composing robust features with denoising autoencoders. In: Proceedings of the 25th International Conference on Machine Learning. ACM; 2008. p. 1096–103.
34. Nair V, Hinton GE. Rectified linear units improve restricted boltzmann machines. In: Proceedings of the 27th International Conference on Machine Learning (ICML-10). USA: Omnipress; 2010. p. 807–14.
35. Chollet F, et al. Keras: GitHub; 2015.
36. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E. Scikit-learn: Machine learning in Python. *J Mach Learn Res*. 2011;12:2825–30.
37. Witten IH, Frank E, Hall MA, Pal CJ. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann; 2016.
38. 2013 OWASP Top 10 Most Dangerous Web Vulnerabilities. [https://www.owasp.org/index.php/Top\\_10\\_2013-Top\\_10](https://www.owasp.org/index.php/Top_10_2013-Top_10). Accessed 13 Aug 2019.
39. <https://www.owasp.org/index.php/Deserialization-of-untrusted-data>. Accessed 16 Aug 2017.
40. yoserial. <https://github.com/frohoff/yoserial>. Accessed 13 Aug 2019.
41. SPECjvm2008. <https://www.spec.org/jvm2008/>. Accessed 13 Aug 2019.
42. Powers DM. Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation. *J Mach Learn Technol*. 2011;2(1):37–63. Bioinfo Publications.
43. Buczak AL, Guven E. A survey of data mining and machine learning methods for cyber security intrusion detection. *IEEE Commun Surv Tutor*. 2015;18(2):1153–76.
44. Opitz DW, Maclin R. Popular ensemble methods: An empirical study. *J Artif Intell Res (JAIR)*. 1999;11:169–98.
45. Wang Y, Wong J, Miner A. Anomaly intrusion detection using one class svm. In: Information Assurance Workshop, 2004. Proceedings from the Fifth Annual IEEE SMC. IEEE; 2004. p. 358–64.
46. Boonstoppel P, Cadar C, Engler D. Rwsset: Attacking path explosion in constraint-based test generation. *Tools Algorithm Constr Anal Syst*. 2008;351–66. Springer.
47. Schneider D. Deeper and cheaper machine learning [top tech 2017]. *IEEE Spectr*. 2017;54(1):42–3.
48. Sekar R, Bendre M, Dhurjati D, Bollineni P. A fast automaton-based method for detecting anomalous program behaviors. In: Security and Privacy, 2001. S&P 2001. Proceedings. 2001 IEEE Symposium On. IEEE; 2001. p. 144–55.
49. Wagner D, Dean R. Intrusion detection via static analysis. In: Security and Privacy, 2001. S&P 2001. Proceedings. 2001 IEEE Symposium On. IEEE; 2001. p. 156–68.
50. Friedenthal S, Moore A, Steiner R. *A Practical Guide to SysML: the Systems Modeling Language*. Morgan Kaufmann; 2014.
51. Scott SL. A bayesian paradigm for designing intrusion detection systems. *Comput Stat Data Anal*. 2004;45(1):69–83.
52. Farid DM, Harbi N, Rahman MZ. Combining naive bayes and decision tree for adaptive intrusion detection. arXiv preprint. 2010. arXiv:1005.4496.
53. Zolotukhin M, Hämäläinen T, Kokkonen T, Siltanen J. Analysis of http requests for anomaly detection of web attacks. In: Dependable, Autonomic and Secure Computing (DASC), 2014 IEEE 12th International Conference On. IEEE; 2014. p. 406–11.
54. Kolosnjaji B, Zarras A, Webster G, Eckert C. Deep learning for classification of malware system call sequences. In: Australasian Joint Conference on Artificial Intelligence. Springer; 2016. p. 137–49.
55. Warrender C, Forrest S, Pearlmutter B. Detecting intrusions using system calls: Alternative data models. In: Security and Privacy, 1999. Proceedings of the 1999 IEEE Symposium On. IEEE; 1999. p. 133–45.
56. Hofmeyr SA, Forrest S, Somayaji A. Intrusion detection using sequences of system calls. *J Comput Secur*. 1998;6(3):151–80.
57. Kearns M, Valiant L. Cryptographic limitations on learning boolean formulae and finite automata. *J ACM (JACM)*. 1994;41(1):67–95.
58. Ilgun K, Kemmerer RA, Porras PA. State transition analysis: A rule-based intrusion detection approach. *IEEE Trans Softw Eng*. 1995;21(3):181–99.
59. Becher M. *Web Application Firewalls*. VDM Verlag; 2007.
60. Desmet L, Piessens F, Joosen W, Verbaeten P. Bridging the gap between web application firewalls and web applications. In: Proceedings of the Fourth ACM Workshop on Formal Methods in Security. ACM; 2006. p. 67–77.
61. Shar LK, Briand LC, Tan HBK. Web application vulnerability prediction using hybrid program analysis and machine learning. *IEEE Trans Dependable Secure Comput*. 2015;12(6):688–707.
62. Leung K, Leckie C. Unsupervised anomaly detection in network intrusion detection using clusters. In: Proceedings of the Twenty-eighth Australasian Conference on Computer Science—Volume 38. Australian Computer Society, Inc.; 2005. p. 333–42.
63. Erfani SM, Rajasegarar S, Karunasekera S, Leckie C. High-dimensional and large-scale anomaly detection using a linear one-class svm with deep learning. *Pattern Recogn*. 2016;58:121–34.
64. Xiong Y, Zuo R. Recognition of geochemical anomalies using a deep autoencoder network. *Comput Geosci*. 2016;86:75–82.
65. Sakurada M, Yairi T. Anomaly detection using autoencoders with nonlinear dimensionality reduction. In: Proceedings of the MLSDA 2014 2nd Workshop on Machine Learning for Sensory Data Analysis. ACM; 2014. p. 4.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Submit your manuscript to a SpringerOpen<sup>®</sup> journal and benefit from:**

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

---

Submit your next manuscript at ► [springeropen.com](https://www.springeropen.com)