

Lom: Discovering Logic Flaws Within MongoDB-based Web Applications

Shuo Wen¹ Yuan Xue² Jing Xu¹ Li-Ying Yuan¹ Wen-Li Song¹ Hong-Ji Yang³ Guan-Nan Si⁴

¹Institute of Machine Intelligence, College of Computer and Control Engineering, Nankai University, Tianjin 300350, China

²Department of Electrical Engineering and Computer Science, Vanderbilt University, Nashville, Tennessee 37212, USA

³Centre for Creative Computing, Bath Spa University, Bath, BA2 9BN, UK

⁴School of Information Science and Electrical Engineering, Shandong Jiaotong University, Ji'nan 250357, China

Abstract: Logic flaws within web applications will allow malicious operations to be triggered towards back-end database. Existing approaches to identifying logic flaws of database accesses are strongly tied to structured query language (SQL) statement construction and cannot be applied to the new generation of web applications that use not only structured query language (NoSQL) databases as the storage tier. In this paper, we present Lom, a black-box approach for discovering many categories of logic flaws within MongoDB-based web applications. Our approach introduces a MongoDB operation model to support new features of MongoDB and models the application logic as a mealy finite state machine. During the testing phase, test inputs which emulate state violation attacks are constructed for identifying logic flaws at each application state. We apply Lom to several MongoDB-based web applications and demonstrate its effectiveness.

Keywords: Logic flaw, web application security, not only structured query language (NoSQL) database, black-box, MongoDB.

1 Introduction

Web applications have become a major information access portal these years. These applications interact with back-end databases on behalf of their users. The back-end database executes all the operations requested by the web application with its privileges, and therefore the application is indispensable for ensuring security checks effective before the database accepts an operation. So web applications become one of the primary targets for maliciously acquiring or manipulating the sensitive information in back-end databases. One category of attacks exploits the application's input validation mechanisms that allow malformed user inputs to be used for constructing database operations, e.g., structured query language (SQL) queries. The case of notorious SQL injections belongs to this type. Another category of attacks, which is referred to as state violation attacks^[1], exploits logic flaws within the application. This type of attacks misleads the application into sending database operation at incorrect application states.

In contrast to input validation vulnerabilities which have

received considerable attention, only limited works have been presented to address logic flaws. The key challenge comes from the fact that logic vulnerabilities are specific to intended functionality of a particular web application, hence general approaches that can be applied to all web applications require an automated way of deriving the application's intended logic or specification.

On the other hand, not only structured query language (NoSQL) databases are increasingly being employed as alternative to traditional SQL databases. Their notable characteristics, such as flexible data models, scalable data storage, nicely support the need of web applications where the workloads are massive and data sources may not have a pre-defined structure. Such flexibility also brings higher risk of logic vulnerabilities into the web applications. However, to the best of our knowledge, no previous work has made efforts to address logic flaws in web applications with NoSQL database as a backend.

In this paper, we present Lom, the first systematic black-box approach which discovers logic flaws of database access within MongoDB-based web applications. The reason why we choose MongoDB is two folds: 1) According to the DB-engines ranking^[2], the popularity of MongoDB is the top 1 among all the NoSQL databases. 2) As far as data modeling is concerned, MongoDB, which has a complicated hierarchical data model, is a representative NoSQL database.

Although a few existing solutions aim to address logic vulnerabilities within web applications, the characteristics of MongoDB make their approaches not applicable for MongoDB-based web applications: 1) Identical MongoDB operations represented in distinct programming languages

Research Article
Manuscript received December 16, 2015; accepted August 8, 2016;
published online December 29, 2016

This work was supported by China Scholarship Council, Tianjin Science and Technology Committee (No. 12JCZDJC20800), Science and Technology Planning Project of Tianjin (No. 13ZCZDZX01098), NSF TRUST (The Team for Research in Ubiquitous Secure Technology) Science and Technology Center (No. CCF-0424422), National High Technology Research and Development Program of China (863 Program) (No. 2013BAH01B05), and National Natural Science Foundation of China (No. 61402264).

Recommended by Associate Editor Xun Chen
© Institute of Automation, Chinese Academy of Sciences and Springer-Verlag Berlin Heidelberg 2017

have various appearances. However, previous static analysis approaches^[3–5], which can only address pattern-unchanged SQL queries or specific languages, cannot handle the diversified MongoDB operation appearances of multiple programming languages. 2) Some black-box approaches^[6, 7], which can only target the flat data model of relational databases, are not appropriate for the hierarchical and flexible data model of MongoDB. 3) Many static techniques^[3, 5] require the source code of applications for analyzing, or can only be applied to specific web development languages and platforms^[5, 8]. 4) A few approaches^[6] need to access to server-side session information. 5) Some previous approaches^[8–10] can address only one specific vulnerability and cannot be easily extended to handle other forms of logic flaws.

By contrast, our approach supports the features of MongoDB. We explore the protocol layer to extract the MongoDB operation regardless of programming languages and introduce MPath, which is an XPath-like representation to locate each value in the hierarchical model within the MongoDB operation. In addition, our technique is designed to be general and cover many kinds of logic vulnerabilities.

The logic of a web application is modeled by a mealy finite state machine^[11] (mealy FSM). To discover logic vulnerabilities, the intended state machine is built as a partial state machine over the expected user inputs (MongoDB operations) when users follow the navigation paths within the web application first. After that, on basis of the inferred intended mealy machine, we generate unexpected test inputs to exploit logic vulnerabilities within the application. These test inputs are related to three categories of attacks. After producing test inputs, we send the test web requests to web applications and evaluate the outputs to discover potential logic flaws.

Our contributions are summarized as follows:

1) We present a novel black-box approach for discovering logic vulnerabilities within MongoDB-based web applications. In particular, by observing the messages in the protocol layer, our approach introduces a MongoDB operation model to represent the MongoDB actions triggered within the web application. We characterize the logic flaws over the mealy FSM, systematically utilize the observed user inputs for deriving the specification and generate test inputs to exploit vulnerabilities.

2) Our approach is able to cover numerous categories of logic flaws without the need of application source code and server-side session information, therefore it can support different coding languages and environments.

3) We implemented a prototype system Lom and demonstrate that Lom can be used to identify logic flaws in today's MongoDB-based web applications.

The rest of this paper is organized as follows. Section 2 discusses related work. Our problem formulation is presented in Section 3. Our approach and implementation are illustrated in details in Sections 4 and 5, respectively. Sec-

tion 6 presents our experimental results and finally the paper is concluded in Section 7.

2 Related work

To the best of our knowledge, only two existing researches make efforts on NoSQL database security. Okman et al.^[12] analyze the main functionality and security features of two popular NoSQL databases: MongoDB and Cassandra. Aniello et al.^[13] analyze the vulnerabilities of the gossip-based membership protocol used by Cassandra. Nonetheless, none of these approaches concentrates on the flaw within NoSQL database based web applications, while our approach detects logic flaws within modern MongoDB-based web applications.

Most previous researches^[10, 14–16] endeavor to exploit various vulnerabilities within web applications. For instance, SecuBat^[17] is used to identify input validation vulnerabilities.

Nevertheless, very few techniques address logic flaws within modern web applications. There are two categories of approaches researched for securing legacy web applications from logic flaws:

1) Vulnerability analysis: It tries to identify and fix the logic vulnerabilities within the applications.

2) Attack detection: It tries to detect and block logic attacks launched against the vulnerable applications.

The key issue that is common in both approaches is how to derive the application logic specification. Then, the logic specification is used for either attack detection or vulnerability analysis.

The logic specification that is general to a number of web applications can be manually pre-specified. Dalton et al.^[18] aim at providing reliable authentication and authorization mechanisms for web applications. By modifying the language runtime, it can track users' credentials and enforce pre-specified security policies over resources, such as files, database objects, etc. Parno et al.^[19] employs virtualization technology to isolate the application components for different users, so that the current user can only access his/her own data. However, more commonly, the logic specification is specific to each application and not available as a priori. SWADDLER^[1], BLOCK^[20] and SENTINEL^[6] establishes application-specific behavioral models and identify the runtime deviation from the established model as potential logic attacks. In particular, SENTINEL focuses on securing the database access triggered by the web application based on a set of invariants extracted from execution traces. The objective of these work is to detect whether a given user input violates the application specification, while our objective is to effectively identify concrete inputs to the web application which can violate the specification, which is much more challenging.

Our work shares the same objective of identifying logic flaws within web applications as a number of existing works. Cova et al.^[1], Bisht et al.^[3], Son et al.^[4], Felmetzger et

al.^[5], Doupé et al.^[8], Sun et al.^[21], Balzarotti et al.^[22] and Son et al.^[23] infer the logic specification from application source code, through either static analysis or instrumentation. However, these techniques are language-dependent and limited in the spectrum of logic flaws they can deal with by their capability of handling language details. For example, Waler can only identify violations of value-related invariants in JSP web applications, which are inferred from dynamic executions^[5]. Sun et al.^[21] assume a strong role lattice model for identifying access control flaws within PHP web applications. WAPTEC collects the set of constraints along the paths leading to sensitive operations and constructs exploits to circumvent the security checks^[3]. Doupé et al.^[8] specifically focus on execution after redirection (EAR) vulnerabilities in ruby web applications by analyzing control flows from application source code.

In contrast, our approach extracts the MongoDB operation from the protocol layer without source code requirement, and can be utilized for all programming languages supported by MongoDB. Moreover, most of the above approaches target only one specific vulnerability and cannot be easily extended to handle other categories of logic flaws. Our technique is designed to be general and covers many kinds of logic vulnerabilities.

Techniques are also designed to discover logic flaws within web applications without source code. For example, Doupé et al.^[8] and Bisht et al.^[9] can address EAR vulnerability and parameter tampering respectively. In comparison, our approaches can cover not only these two attacks, but also forceful browsing attack. InteGuard^[24] and EURECOM^[25] attempt to secure multi-party web applications. LogicScope^[26], SENTINEL^[6] and BLOCK^[20] make use of session information to construct application specifications. In comparison, our work does not require server-side session information from the application developers. Li et al.^[7] propose an automated black-box technique for identifying access control vulnerabilities. Though SENTINEL^[6] and the work of [7] can be applied to traditional relational database, they cannot handle the hierarchical and schema-less data model of MongoDB, which brings in new challenges. Our technique supports these new features of MongoDB back-end web applications^[27].

Web applications are more and more built with third-party web services through application programming interfaces and split at both client-side and server-side, where logic vulnerabilities might arise. Wang et al.^[28] discover logic vulnerabilities within the checkout procedures, which can be exploited by the attackers to shop for free. Its further research^[29] also identified logic vulnerabilities within web-based single-sign-on services. Xing et al.^[24] perform security checks over a set of invariant relations among HTTP interactions to defeat logic attack at runtime. Wu et al.^[30] employs hybrid analysis to infer the dependency constraints on parameters for web services. Guha et al.^[31] extract event graphs from client-side web applications and detect malicious client behaviors at runtime. Capsules^[32] can be

used to build secure web applications, where security policies specified by developers can be automatically verified and enforced. Our technique focuses on logic vulnerabilities within server-side web applications and has the potential to be extended to handle the above scenarios.

A number of testing tools, both open-source, e.g., Spike, Burp, IBM AppScan, have been proposed for identifying input validation vulnerabilities within web applications^[15]. They feed random inputs from a library of known attack patterns into applications. To improve the testing coverage and efficiency, random fuzzing can be enhanced by guided test input generation^[16, 33, 34]. None of these techniques can effectively handle logic vulnerabilities within web applications.

3 Problem description

3.1 Background of MongoDB

3.1.1 Data model of MongoDB

Document. In MongoDB, the basic unit of data is document whose structure is hierarchical and non-relational. A document includes a set of field/value pairs where the value of a field can even be a document or an array which is a list of values. Array values can be all the supported values for normal field/value pairs in MongoDB, even nested arrays and embedded documents. Fig. 1 shows a document which employs embedded documents and array values.

```
{
  "$or" : [
    {
      "Number" : {
        "$lt" : 10
      }
    },
    {
      "selected" : true
    }
  ]
}
```

Fig. 1 A document / A MongoDB request variable

Collection. MongoDB documents are grouped as one or more collections in a MongoDB database. The schema of a collection does not need to be defined while the collection is created, which means users have more data-modeling flexibilities to match the design and performance requirements of an application.

3.1.2 MongoDB wire protocol

MongoDB offers many additional drivers for users to work with their proficient programming languages. The same operations represented in distinct drivers may have different appearances. To avoid this difference, we focus on internals of how drivers access the MongoDB server. The drivers use MongoDB wire protocol, which is a simple socket-based, request-response style and lightweight transmission control protocol/internet protocol (TCP/IP) wire protocol, to make clients communicate with the MongoDB server through MongoDB request messages. A message

defines the concrete data which an operation can access and the type of the operation. With these messages, update, delete, insert and read operations can be performed on MongoDB.

3.1.3 MongoDB request variable

Fig. 2 shows the structure of a category of MongoDB request message (update message). As can be seen from Fig. 2, the data structures of most useful variables in MongoDB wire protocol are documents, such as the query selector and the update definition. This structure is able to support complex commands. For instance, Fig. 1 is also a MongoDB request variable (query selector). Here “\$lt” is a

```
struct OP_UPDATE {
    MsgHeader header;           // standard header
    int32  ZERO;                // reserved for future use
    cstring fullCollectionName; // "databaseName.collectionName"
    int32  flags;                // 0 - upsert; 1 - multiupdate
    document querySelector;      // to select the document
    document updateDefinition;   // to specify the update to perform
}
```

Fig. 2 Structure of the update message

comparison operator corresponding to “less than”. Each of these document structure variables is denoted as MongoDB request variable in this paper. Apparently, all the operation parameters are placed in these hierarchical and non-relational variables.

3.2 Logic flaws within MongoDB-based web applications

Fig. 3 shows a simple vulnerable application to illustrate the logic vulnerabilities we concentrate on in our research. A logged in user will be redirected to the “index.php” at first. If the current user is an administrator, he is allowed to achieve links for adding new users, editing and deleting any of the registered users. If the current user is a regular user, he can only browse the page for editing his personal information.

We model a web application using a mealy finite-state machine (mealy FSM) model $(S, s_0, \Sigma, \Lambda, T, G)$, where S is the set of states, $s_0 \in S$ is the initial state, Σ is the set of input symbols, Λ is the set of output symbols, $T : S \times \Sigma \rightarrow S$ is the set of transition functions mapping pairs of a state and an input symbol to the corresponding next state, $G : S \times \Sigma \rightarrow \Lambda$ is the set of output functions mapping pairs of a state and an input symbol to the corresponding output symbol.

To find out the logic flaws within a web application, we are required to analyze its two categories of mealy FSMs:

- 1) Intended FSM (denoted as F_i), which models the behavior of the originally planned web application without any logic flaws.
- 2) Realistic FSM (denoted as F_r), which models the be-

index.php	delUser.php
<pre><?php if (\$_SESSION['privilege'] == "admin") { \$alluser = getAllUsers(); foreach (\$alluser as \$eachuser) { echo "Delete "; echo "Edit "; } echo "Add "; } else if (\$_SESSION['privilege'] == "commonuser") { echo "Edit"; } ... ?></pre>	<pre><?php if (\$_SESSION['privilege'] != "admin") echo("Forbidden access."); //op2 require('dbconnection.php'); \$mongo = DBConnection::instantiate(); \$collection = \$mongo->getCollection('users'); \$id=\$_GET['userid']; \$collection->remove(array('id' => new MongoClient(\$id))); ... ?></pre>
<pre>editUser.php //op1 require("dbconnection.php"); \$mongo = DBConnection::instantiate(); \$collection = \$mongo-> getCollection("users"); \$id=\$_GET["userid"]; \$user = array(); \$user["name"] = \$_POST["name"]; \$user["password"] = \$_POST["pwd"]; \$collection-> update(array('id' => new MongoClient(\$id)), \$user); ... ?></pre>	<pre>addUser.php //op3 require('dbconnection.php'); \$mongo = DBConnection::instantiate(); \$collection = \$mongo->getCollection('users'); \$user =array('name' => \$_POST['name'], 'password' => \$_POST['pwd']); \$collection->insert(\$user); echo 'User created successfully'; ... ?></pre>

Fig. 3 Example application

havior of the actual web application implemented by the developer.

If F_r is equivalent to F_i , the web application is regarded as secure. Once disparities which involve sensitive operations exist between F_r and F_i , we affirm the application has logic flaws.

As illustrated in Fig. 4, the example application has three states: the guest user who is not logged in (s_0), regular user (s_1) and administrator (s_2). Each input symbol $I \in \Sigma$ is an abstract representation of the triggered operation on back-end MongoDB (e.g., op1, op2 and op3 in Fig. 3), which consists of two parts:

1) Operation contour (denoted by C), which represents the contour of the operation (refer to Section 4.3.1 for details).

2) Transmitted parameter mapping (denoted by P), which represents both the parameter which can be transmitted from web request to the operation and its related value set (refer to Section 4.3.2 for details).

Each output symbol in Λ is the acceptance of the operation by back-end MongoDB.

The intended FSM (F_i) for the application works as follows: At state s_1 , since it is intended that the regular user can only edit his personal information, when the regular user sends an input symbol $I_1 = C_1 \times P_1$, where the “userid” parameter is equal to the current user id, back-end MongoDB will accept this operation (output symbol O_1). When this user attempts to edit other users’ information, delete

or add a user, i.e., send I_2 (I_2 is different from I_1 due to the diverse parameter mappings.), I_3 or I_4 , MongoDB will not accept or trigger the operation (output symbol O_2).

Nonetheless, in this application, there are three logic flaws which are reflected as the discrepancies between F_i and F_r . First, the “editUser.php” fails to check whether the “userid” parameter is the same as current user’s information. Second, despite the “delUser.php” checks whether current user is an administrator and seems to reject the operation from web response, it does not end the application execution, thus the back-end MongoDB operation is still triggered. Third, the “addUser.php” does not check whether the current user has the admin privilege. These vulnerabilities allow three types of attacks:

1) Parameter manipulation attack: When I_2 is sent to the application at state s_1 , O_1 is returned, which means a regular user can edit other users’ information.

2) Execution after redirection (EAR) attack^[8]: When I_3 is sent to the application at state s_1 , O_3 is returned, which means a regular user can still successfully make back-end MongoDB delete other users’ information although O_2 appears to be returned from web response.

3) Forceful browsing attack: When I_4 is sent to the application at state s_1 , O_4 is returned, which means a regular user can add new users.

All the attacks mentioned above are common attacks targeting different kinds of logic vulnerabilities within database based web applications. EAR attack is especially

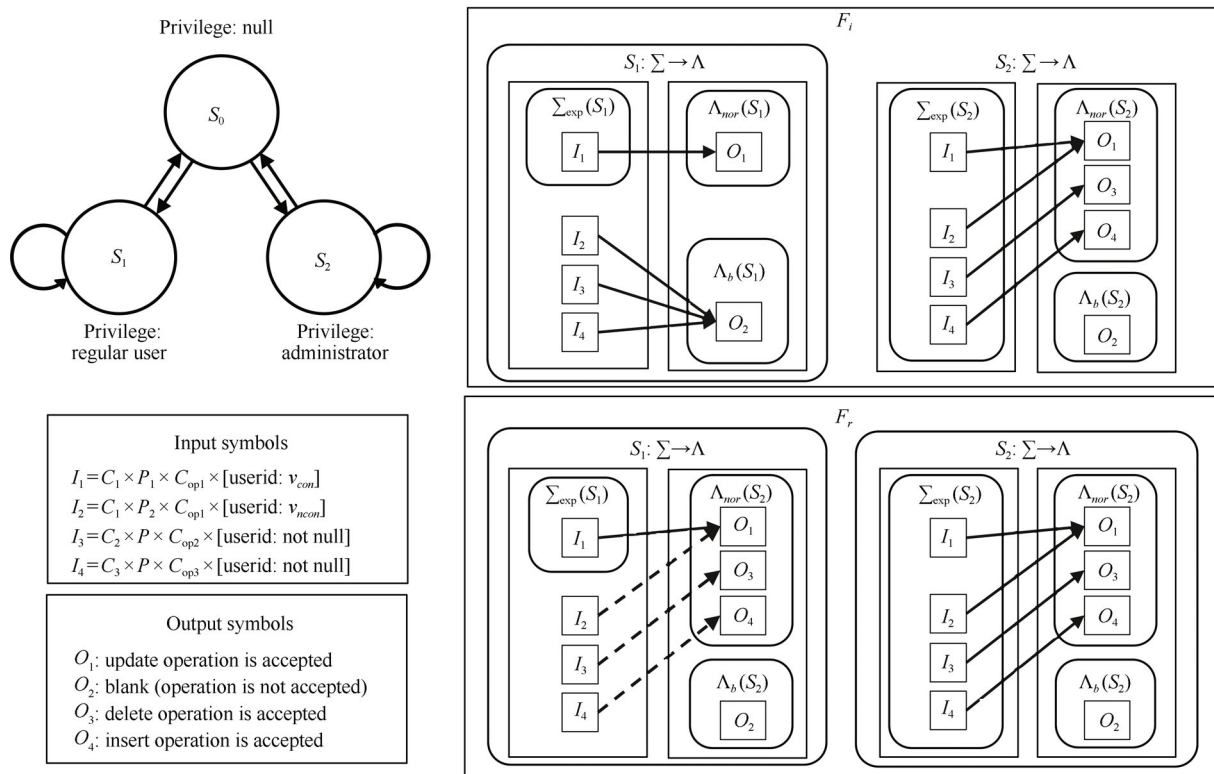


Fig. 4 FSM representation of Fig. 3

challenging. This attack seems to be defended from web response. However, the back-end database still triggers the database operation which is not designed to run.

At a given state s , only a subset of input symbols are expected by the application (denoted as $\Sigma_{\text{exp}}(s)$) and processed to produce normal output symbols, i.e., $\Lambda_{\text{nor}}(s) = G(s, \Sigma_{\text{exp}}(s))$. The expected input symbols are the triggered MongoDB operations when the user follows the navigation links of the web application. The normal output symbols mean that MongoDB accepts the expected MongoDB operations. All the other input symbols, which are not expected at state s , should not be triggered by MongoDB, resulting in blank output symbols. A blank output symbol means that the application refuses to accept the operation and therefore back-end MongoDB does not execute anything. As shown in Fig. 4, for state s_1 , the expected input set is $\{I_1\}$, the normal output set and the blank output set is $\{O_1\}$ and $\{O_2\}$, respectively. For state s_2 , the expected input set is $\{I_1, I_2, I_3, I_4\}$ and the normal output set is $\{O_1, O_3, O_4\}$. The behaviors of F_i and F_r over the expected input symbols should be consistent because the web application aims at implementing all the intended functionalities. Nevertheless, there may be unexpected inputs which are accepted by F_r . Therefore, if an input symbol, which is not expected at state s , can be transmitted into the application and triggered by MongoDB, MongoDB then generates an output symbol that falls beyond the blank output set, we recognize this web application has a logic vulnerability at state s . The related input symbol is defined as a malicious input symbol (I_{mal}).

4 Proposed approach

4.1 Approach overview

As mentioned in Section 3.2, we need to construct malicious inputs to verify their outputs for each state. It is a challenging task because we do not possess anything about the entire input symbol set and unexpected input symbol set at each state. Some malicious inputs, e.g., EAR attacks, can modify the data in back-end MongoDB secretly without affecting intended web responses. To symbolize the input symbol, we need to learn the operation over MongoDB (Section 4.2). The characteristics of MongoDB make the understanding more sophisticated:

1) As illustrated in Section 3.1.2, the same MongoDB operation may have dissimilar expression in different programming language and furthermore, an operation may be characterized by several statements in the source code (such as op1, op2 or op3 in Fig. 3). Hence, we utilize dynamic analysis but not static analysis to make our approach not constrained to specific programming language or driver. We look into the protocol layer, which is the underlying unification of distinct drivers, to extract the MongoDB operation no matter which programming language the application is written in.

2) As Section 3.1.3 shows, the basic data model of MongoDB, which is also utilized in the MongoDB request variable, is hierarchical and non-relational. MongoDB request variables are the most important components of MongoDB request messages. Thus, we need to locate each field/value pair in the hierarchical data model. We present MPath to support this nested data structure.

Our approach first builds a partial mealy FSM over the expected input domain by leveraging the collected traces. For each application, we identify user privileges and construct each privilege as a State. Normal users' traces are collected for different users at each state. The traces we collect include web requests/responses and MongoDB requests/responses from protocol layer. The traces are symbolized as following:

1) Input symbolization (Section 4.3), in which we abstract concrete MongoDB operations into input symbols to profile the expected input domain at each state, i.e., $\Sigma_{\text{exp}}(s), \forall s \in S$.

2) Output symbolization, in which we observe whether the MongoDB accepts the operations or not for generating output symbols and the mappings between the expected inputs and normal output symbols, i.e., $G(s, \Sigma_{\text{exp}}(s)) \rightarrow \Lambda_{\text{nor}}(s), \forall s \in S$. Application state transitions and the corresponding input symbols that trigger the transitions are also observed in this phase, i.e., $T : S \times \Sigma \rightarrow S$.

After the inference of partial FSM, we will leverage this inferred FSM to construct unexpected inputs at each application state (Section 4.4) and test the application. Output symbols will be evaluated to discover potential logic flaws (Section 4.5).

4.2 MongoDB operation analysis

A MongoDB operation is related to the read, delete, update or insert message in MongoDB wire protocol. It can read or modify the records in MongoDB. We extract the kernel information (message/operation type, collection name and the MongoDB request variables) of a message as its MongoDB operation, which represents its execution on MongoDB performed by the user through the web application.

4.3 Input symbolization

Given a set of MongoDB operations, we need to represent them with a finite number of input symbols. We symbolize each MongoDB operation with a two-part structure, i.e., the operation contour and the transmitted parameter mapping.

4.3.1 Variable/Operation contour

MPath and variable contour. Since MongoDB request variables are included in MongoDB operations, all the MongoDB request variables in the operations need to be stored reasonably. So the main challenge is how to model all of these variables in a more efficient way for convenient comparison, i.e., locating each parameter easily.

To locate each parameter, we introduce MPath which is an analogue of XPath. As an example, the “\$lt” parameter of Fig. 1 can be expressed as “\$or/Number/\$lt”. With this kind of effective representation, we can express the original hierarchical MongoDB request variable.

We then define the contour of a MongoDB request variable as the variable without any parameter values. Each original variable is represented as its extracted contour and its parameter value set. For instance, the contour of the Fig. 1 can be represented as Fig. 5, where “p1” and “p2” represent the values of related parameters. The top of Fig. 5 is the document view of the contour and the bottom is the MPath view which is the implementation. Both the contour and the parameter set are derived from its original variable.

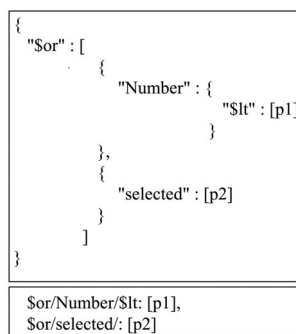


Fig. 5 A variable contour (2 appearances)

Operation contour. The operation type, collection name and the variable contours of a MongoDB operation are denoted as operation contour. Similarly, each operation is represented as its contour and its parameter value set.

4.3.2 Transmitted parameter mapping

We group all MongoDB operations based on their contours as well as the kernels of their respective web requests. A web request kernel includes HTTP method and request URL path without URL parameters. Each group is denoted as an operation group. For a MongoDB operation mo and its related web request wr , we denote a web request parameter of wr is p_{wr} and its value is $v_{p_{wr}}$, a MongoDB operation parameter of mo is p_{mo} , and its value is $v_{p_{mo}}$. If $\exists p_{wr}, p_{mo} \wedge v_{p_{wr}} = v_{p_{mo}}$ holds for all MongoDB operations and web requests within the same operation group, we define there is a parameter transmission path from p_{wr} to p_{mo} , and denote p_{wr} and its related value set $V_{p_{wr}}$ as a transmitted parameter mapping.

4.3.3 Symbolization

We first profile each transmitted parameter mapping and construct this part based on its related value set, i.e., the values of all transmitted parameters. The characterization of each value domain is a two step process. The constraints between the parameter value set and the specific state, i.e., privilege, are extracted first by profiling each parameter at each state. For each state, the value set collected for each parameter within the same operation group is utilized for grouping the parameter into three categories:

1) Random parameter (denoted as $para_{ur}$): The value set of this type of parameter has no limitation. Its value domain is represented with two values: null and not null.

2) Unbounded constrained parameter (denoted as $para_{uc}$): The value set of this type of parameter is affected by certain constraints though it is infinite. Single privilege-related constraint is our focus in this paper, which means the parameter value is always specific for each user under this state (e.g., the value of “userid” of “editUser.php” at s_1 in Fig. 3 is particular for each user under s_1). Its value domain is represented with three kinds of values: null, v_{con} and v_{ncon} , where v_{con} denotes the value satisfying a constraint linked to a specific user under this state and v_{ncon} denotes other values.

3) Bounded parameter (denoted as $para_b$): We represent its value domain with the value set and two kinds of values: null and v_{outb} , where v_{outb} denotes the values out of the bounded set.

We aggregate all the state views of the parameter value domains into a macroscopic view afterwards. If the value domain types of a parameter is consistent for all states, its domain type will not be changed and it will not be recomputed. For $para_{uc}$, the updated value domain is value set divided by constraints. For $para_b$, its value domain adds additional values. If the value domain types of the parameter over different states are disparate, the more restrictive type (the restrictiveness order is defined as $para_{ur} < para_{uc} < para_b$) is adopted and its value domain will be divided. For instance, the parameter “userid” of “editUser.php” in Fig. 3 is constrained by specific user at s_1 , but inferred as $para_{ur}$ at s_2 . So its macroscopic type will be $para_{uc}$ and two input symbols are produced at s_2 , i.e., $C_1 \times P_1$ and $C_1 \times P_2$.

4.4 Test input symbol generation

As Fig. 6 illustrates, there are two methods designed for generating test input symbols at a given state s .

4.4.1 Parameter manipulation

For an expected input symbol $I = C \times P \in \Sigma_{exp}(s)$ at state s , we manipulate P directly, i.e., values of one or more parameters will be changed so as to make the tampered input symbol not included by the expected input set at state s . For an unbounded constrained parameter, we modify its value from v_{con} to v_{ncon} . For a bounded parameter, its value is changed to another value in the bounded set or v_{outb} . The left of Fig. 6 shows an example, P_1 of input symbol I_1 is manipulated so as to generate a test input I_{mal} for s_1 . This method exhibits parameter manipulation attacks, where parameter values are manipulated for violating constraints between operations and the current state.

4.4.2 Forceful browsing/execution after redirection

We observe another state s' which has one or more expected input symbols excluded from the expected input set of current state s . Input symbols at s' with operation contours which are not included by the expected input symbol

set of state s are chosen as test input symbols for s , i.e., $I_{mal} \in \Sigma_{exp}(s') - \Sigma_{exp}(s') \cap \Sigma_{exp}(s)$. The right of Fig. 6 shows an example, the input symbols at state s_2 with C_3 and C_4 are selected as test inputs for state s_1 since they are not included by s_1 . This method exhibits two types of attacks:

- 1) Forceful browsing attacks: One or more hidden sensitive links which should not be accessible at current state can be forcefully browsed.
- 2) Execution after redirection attacks: The attacker seems to be blocked by the application from the web response of the page, but the sensitive MongoDB operations related to the page can still be successfully run on back-end MongoDB. These EAR attacks, which only manipulate the data stored in MongoDB, violate the state secretly.

4.5 Output evaluation

We denote the output symbol generated after the test input I_{mal} being delivered into the application at state s as O_{test} . The output evaluation will determine whether O_{test} belongs to the blank output set. Since the blank output symbol means that the application refuses to accept the operation thus if back-end MongoDB triggers the operation of the test input, O_{test} falls out of the blank output set.

We collect the traces during the testing, after all the test inputs are delivered and the traces are gathered, we analyze each interaction in the traces to examine whether each test operation has been triggered or not. If a test operation is performed, we recognize its related test input as a potential logic flaw.

5 Implementation

We implemented a prototype system Lom for discovering logic vulnerabilities within MongoDB-based web applications. As Fig. 7 shows, Lom has three major components, including trace collector, specification analyzer and testing engine. These components are corresponding to three phases: trace collection, specification inference and testing.

5.1 Phase 1: Trace collection

Trace collector, which collects the communication between the web application/MongoDB and the client when users navigate through the application during attack-free sessions, is implemented in our research by utilizing the open source network protocol analyzer wireshark.

5.2 Phase 2: Specification inference

Specification analyzer is executed in Phase 2 to derive both the partial mealy FSM and the testing specification. Symbolizer first transforms collected traces (in Phase 1) into symbolized session logs. Then, session logs are used by the mealy FSM analyzer module to derive the partial FSM, resulting in two files: StateProfile, which characterizes the mapping between input/output symbols at each state and DriverSpec, which records the transitions between the set of application states, as well as the input symbols that trigger the transitions. Finally, StateProfile is analyzed by TestSpec generator to generate the testing specification, which includes both a set of test input symbols for each state and their related output symbols for evaluation.

5.3 Phase 3: Testing

Testing engine is executed in Phase 3 to test whether the application has logic vulnerabilities, based on the above derived profiles and specifications. It will produce test web requests from test inputs (by web request generator), deliver them into the application and evaluate the test traces for logic flaw identification (by output evaluator).

Testing controller is the core module that takes charge of the entire testing procedure. It first loads TestSpec and other profiles and checks the current application state. If the test of the current state is not completed, it retrieves the next available test input symbol, delegates request generator to generate a concrete web request and submit to the application. After it receives the web response, it will wrap up all the necessary information and send it to output evaluator for evaluation, where logic vulnerabilities, if exist,

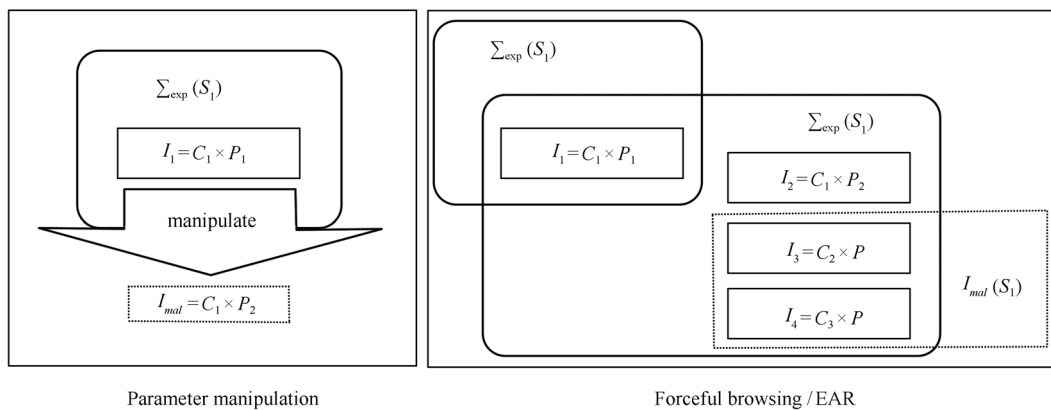


Fig. 6 Test input generation

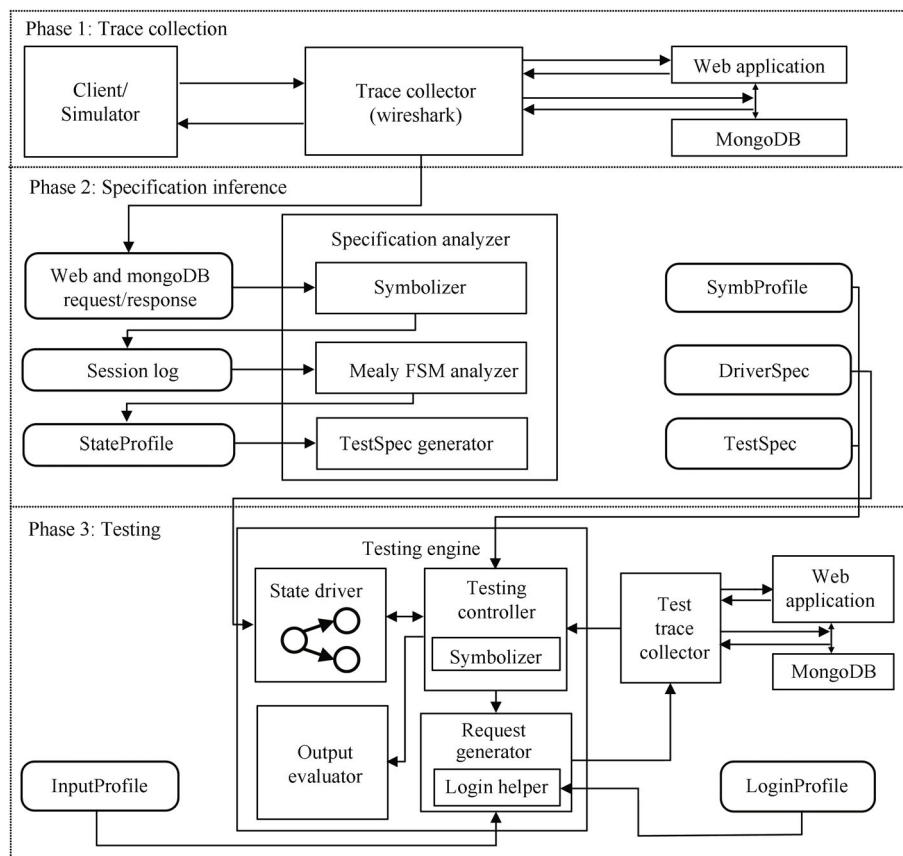


Fig. 7 Prototype system architecture

will be reported. If the test of the current state is completed, i.e., no test inputs are left, Testing controller will move to the next available test state. It will consult state driver, which loads DriverSpec and keeps track of the transition graph of the application, to get the path leading to the next test state. The path computed by state driver is essentially the shortest path from the current state to the target state (i.e., a sequence with minimum number of input symbols), which will be instantiated by request generator and trigger the state transition step by step to the target state. This mechanism is desirable, since we cannot directly drive the application into our desired abstract state. For the example application, after we have tested state s_1 for the regular user, we have to first log out (i.e., move to state s_0) and log in as an administrator to test state s_2 . If all the states have been fully tested, the testing procedure is finished.

One key challenge we need to address is how to instantiate abstract input symbols into concrete web requests with meaningful parameters. In Phase 2, when we profile web requests, we also infer the value type (e.g., number, literal string) of each parameter. When request generator tries to generate the concrete value for a parameter, it checks its value type and randomly generates a value of that type or retrieve a value from a pre-loaded value store (i.e., InputProfile). In particular, request generator includes login

helper module, which helps testing engine successfully log into the application. For each type of user, e.g., regular user, administrator, login helper requires the user to provide a LoginProfile file, which specifies the input symbol that represents the login request and at least one set of legitimate user credentials, e.g., username and password.

6 Evaluation

We choose a set of interactive MongoDB-based web applications for evaluating our prototype system Lom. We deploy all web applications on a 3.30 GHz Intel core i3-2120 Linux server with 4GB RAM. To facilitate trace collection, we build user simulators for each application based on Selenium WebDriver. We first identify user privileges and their corresponding atomic actions by following navigation links. All of the atomic actions can be recognized as intended functions by the web application designers due to each of them follows the navigation paths under normal situation implemented by the designer. Therefore, the correctness can be guaranteed. Then, the simulator performs a random sequence of atomic actions automatically with different privileges and users, each user will run all the atomic actions under his state at least once.

Our inference is performed through dynamic analysis, where the web application is executed under the constraint

of navigation paths. This constraint has been applied in several existing approaches^[5, 21] and shown to be effective and general to cover a large number of web applications.

6.1 Analysis of results

Lom first runs in Phases 1 and 2 to collect traces and infer the application logic specification. The statistics of collected traces and inferred FSMs are shown in Table 1, including the number of files, collected web requests, MongoDB requests, states, input symbols. Then, Lom generates the testing specification and launches the testing procedure against each web application. It constructs test web requests and sends them to the application. Testing evaluator then evaluates the test inputs based on the collected test traces. One feature of Lom is that it also gives concrete attack vectors and evidences for further inspection.

Table 2 shows the testing results, including the number of test inputs generated by each method, flagged attacks and false positives. We also report the sum of real attacks (true positives) and vulnerable web pages. Note that these two numbers can be different, because a web page may have one or more unexpected operations which can be triggered under different states. In the following, we describe the details of logic flaws we identify from each web application. As Table 2 shows, 31 vulnerable web pages are discovered without false positive.

6.1.1 MongoBlog

There are three states in this web application: guest, regular user and admin user. Regular users can post new articles, add comments under articles, edit or delete the articles or comments created by the user. Admin users can manage all the articles and comments. Either a regular user or an admin user can mark his favorite articles. Several logic vulnerabilities are identified within this application. First, forceful browsing attacks can be applied on the application, guest users can publish and manage articles and comments as other types of users. Second, the application can be attacked by parameter manipulation, a regular user can view other user's summary page which shows articles, comments or favorite articles of corresponding user by manipulating a parameter.

6.1.2 QuickBlog

This application also has three states: guest, regular user and administrator. Only the administrator is allowed to modify all of the posts. The regular user can edit or delete his own posts. Logic flaws exist within administrative or regular users' pages which fail to check the current application state before any database operations. Thus, an attacker can forcefully browse those pages and trigger sensitive operations, a regular user can perform parameter manipulation attacks to view other regular user's pages.

6.1.3 SimpleNote

There are three states in SimpleNote: regular user, user manager and super administrator. Each regular user can

Table 1 Summary of traces and inferred FSM

Web application	File	Web request	MongoDB request	State	Input
MongoBlog	41	371	1 165	2	24
QuickBlog	15	336	346	3	11
SimpleNote	21	437	493	3	10
ProductShow	8	65	25	2	2

Table 2 Summary of testing results

Web application	Method	Test inputs	Flagged attacks	False positives	True positives	Vulnerable web pages
MongoBlog	FE	11	10	0	14	13
	PM	4	4	0		
QuickBlog	FE	21	13	0	14	10
	PM	1	1	0		
SimpleNote	FE	34	18	0	18	7
	PM	0	0	0		
ProductShow	FE	1	1	0	1	1
	PM	0	0	0		
Summary		72	47	0	47	31

FE: Forceful browsing and execution after redirection.

PM: Parameter manipulation.

True positives: i.e., the sum of real attacks.

only view, edit and delete his own notes. User managers can manage the profile of regular users. Super administrators have the highest privilege, they are allowed to handle all users and notes. We identify logic vulnerabilities within user managers' and super administrators' pages which miss the examination of current application state. These vulnerabilities allow an attacker to browse vulnerable pages directly for managing other users' notes or profiles.

6.1.4 ProductShow

ProductShow has two states: the administrator which can add new product to MongoDB from his own page, the common user which can read products' information. An attacker can forcefully browse the administrative page because the application does not check the current application state.

7 Conclusions

In this paper, we presented the first systematic black-box approach to identify logic flaws within MongoDB-based web applications. A prototype system Lom, which introduces a MongoDB operation model to support new features of MongoDB and models the application logic as a mealy finite state machine, was implemented and evaluated to demonstrate the practical utility of our approach.

With the development of web application technology, based on the method of this paper, there are several related areas, which will be the concentration of our further research, could be extended:

1) The logic flaw within other NoSQL database based web application: The database we concentrate on now is MongoDB which is a representation of NoSQL database. Nevertheless, there are various categories of NoSQL database. The features of other NoSQL database may bring new challenges which may be worth studying.

2) Other kinds of vulnerabilities within NoSQL database based web application: The approach we present in this paper targets logic flaws within NoSQL database based web applications. It is worth considering whether NoSQL database will bring challenges to other security problems, such as input validation vulnerabilities.

In summary, this paper makes progress on discovering logic flaws within MongoDB-based web applications and makes progress and has the value of practical application. Our work also has some reference value on further research of web application security.

References

- [1] M. Cova, D. Balzarotti, V. Felmetsger, G. Vigna. Swaddler: An approach for the anomaly-based detection of state violations in web applications. In *Proceedings of the 10th International Conference on Recent Advances in Intrusion Detection*, Springer, Queensland, Australia, pp. 63–86, 2007.
- [2] *DB-Engines Ranking*, July 2014, [Online], Available: <http://dbengines.com/en/ranking>.
- [3] P. Bisht, T. Hinrichs, N. Skrupsky, V. N. Venkatakrishnan. WAPTEC: Whitebox analysis of web applications for parameter tampering exploit construction. In *Proceedings of the 18th ACM Conference on Computer and Communications Security*, ACM, Chicago, USA, pp. 575–586, 2011.
- [4] S. Son, K. S. McKinley, V. Shmatikov. RoleCast: Finding missing security checks when you do not know what checks are. In *Proceedings of ACM International Conference on Object Oriented Programming Systems Languages and Applications*, ACM, Portland, USA, pp. 1069–1084, 2011.
- [5] V. Felmetsger, L. Cavedon, C. Kruegel, G. Vigna. Toward automated detection of logic vulnerabilities in web applications. In *Proceedings of the 19th USENIX Conference on Security*, USENIX Association, Berkeley, USA, pp. 10, 2010.
- [6] X. W. Li, W. Yan, Y. Xue. SENTINEL: Securing database from logic flaws in web applications. In *Proceedings of the 2nd ACM Conference on Data and Application Security and Privacy*, ACM, San Antonio, USA, pp. 25–36, 2012.
- [7] X. W. Li, X. J. Si, Y. Xue. Automated black-box detection of access control vulnerabilities in web applications. In *Proceedings of the 4th ACM Conference on Data and Application Security and Privacy*, ACM, San Antonio, USA, pp. 49–60, 2014.
- [8] A. Doupé, B. Boe, C. Kruegel, G. Vigna. Fear the EAR: Discovering and mitigating execution after redirect vulnerabilities. In *Proceedings of the 18th ACM Conference on Computer and Communications Security*, ACM, Chicago, USA, pp. 251–262, 2011.
- [9] P. Bisht, T. Hinrichs, N. Skrupsky, R. Bobrowicz, V. N. Venkatakrishnan. NoTammer: Automatic blackbox detection of parameter tampering opportunities in web applications. In *Proceedings of the 17th ACM Conference on Computer and Communications Security*, ACM, Chicago, USA, pp. 607–618, 2010.
- [10] A. Doupé, L. Cavedon, C. Kruegel, G. Vigna. Enemy of the state: A state-aware black-box web vulnerability scanner. In *Proceedings of the 21st USENIX Conference on Security Symposium*, USENIX Association, Berkeley, USA, pp. 26, 2012.
- [11] G. H. Mealy. A method for synthesizing sequential circuits. *Bell System Technical Journal*, vol. 34, no. 5, pp. 1045–1079, 1955.
- [12] L. Okman, N. Gal-Oz, Y. Gonen, E. Gudes, J. Abramov. Security issues in NoSQL databases. In *Proceedings of the 10th International Conference on Trust, Security and Privacy in Computing and Communications*, IEEE, Changsha, China, pp. 541–547, 2011.
- [13] L. Aniello, S. Bonomi, M. Breno, R. Baldoni. Assessing data availability of Cassandra in the presence of non-accurate membership. In *Proceedings of the 2nd International Workshop on Dependability Issues in Cloud Computing*, ACM, Braga, Portugal, Article number 2, 2013.

- [14] P. Chapman, D. Evans. Automated black-box detection of side-channel vulnerabilities in web applications. In *Proceedings of the 18th ACM Conference on Computer and Communications Security*, ACM, Chicago, USA, pp.263–274, 2011.
- [15] Y. W. Huang, S. K. Huang, T. P. Lin, C. H. Tsai. Web application security assessment by fault injection and behavior monitoring. In *Proceedings of the 12th International Conference on World Wide Web*, ACM, Budapest, Hungary, pp. 148–159, 2003.
- [16] M. Martin, M. S. Lam. Automatic generation of XSS and SQL injection attacks with goal-directed model checking. In *Proceedings of the 17th Conference on USENIX Security Symposium*, USENIX Association, Berkeley, USA, pp. 31–43, 2008.
- [17] S. Kals, E. Kirda, C. Kruegel, N. Jovanovic. Secubat: A web vulnerability scanner. In *Proceedings of the 15th International Conference on World Wide Web*, ACM, Edinburgh, UK, pp. 247–256, 2006.
- [18] M. Dalton, C. Kozyrakis, N. Zeldovich. Nemesis: Preventing authentication & access control vulnerabilities in web applications. In *Proceedings of the 18th Conference on USENIX Security Symposium*, USENIX Association, Berkeley, USA, pp. 267–282, 2009.
- [19] B. Parno, J. M. McCune, D. Wendlandt, D. G. Andersen, A. Perrig. CLAMP: Practical prevention of large-scale data leaks. In *Proceedings of the 30th IEEE Symposium on Security and Privacy*, IEEE, Oakland, USA, pp. 154–169, 2009.
- [20] X. W. Li, Y. Xue. BLOCK: A black-box approach for detection of state violation attacks towards web applications. In *Proceedings of the 27th Annual Computer Security Applications Conference*, ACM, Orlando, USA, pp. 247–256, 2011.
- [21] F. Q. Sun, L. Xu, Z. D. Su. Static detection of access control vulnerabilities in web applications. In *Proceedings of the 20th USENIX Conference on Security*, USENIX Association, Berkeley, USA, pp. 11, 2011.
- [22] D. Balzarotti, M. Cova, V. V. Felmetsger, G. Vigna. Multi-module vulnerability analysis of web-based applications. In *Proceedings of the 14th ACM Conference on Computer and Communications Security*, ACM, Alexandria, USA, pp. 25–35, 2007.
- [23] S. Son, K. S. McKinley, V. Shmatikov. Fix me up: Repairing access-control bugs in web applications. In *Proceedings of 20th Network and Distributed System Security Symposium*, Internet Society, San Diego, USA, 2013.
- [24] L. Y. Xing, Y. Y. Chen, X. F. Wang, S. Chen. InteGuard: Toward automatic protection of third-party web service integrations. In *Proceedings of 20th Annual Network and Distributed System Security Symposium*, Internet Society, San Diego, USA, 2013.
- [25] G. Pellegrino, D. Balzarotti. Toward black-box detection of logic flaws in web applications. In *Network and Distributed System Security Symposium*, Internet Society, San Diego, USA, 2014.
- [26] X. W. Li, Y. Xue. LogicScope: Automatic discovery of logic vulnerabilities within web applications. In *Proceedings of the 8th ACM SIGSAC Symposium on Information, Computer and Communications Security*, ACM, Hangzhou, China, pp. 481–486, 2013.
- [27] S. Wen, Y. Xue, J. Xu, H. J. Yang, X. H. Li, W. L. Song, G. N. Si. Toward exploiting access control vulnerabilities within mongodb backend web applications. In *Proceedings of the 40th Annual IEEE Computer Software and Applications Conference*, IEEE, Atlanta, USA, 2016.
- [28] R. Wang, S. Chen, X. F. Wang, S. Qadeer. How to shop for free online – security analysis of cashier-as-a- service based web stores. In *Proceedings of the 32nd IEEE Symposium on Security and Privacy*, IEEE, Berkeley, USA, pp. 465–480, 2011.
- [29] R. Wang, S. Chen, X. F. Wang. Signing me onto your accounts through facebook and google: A traffic-guided security study of commercially deployed single-sign-on web services. In *Proceedings of IEEE Symposium on Security and Privacy*, IEEE, Washington DC, USA, pp. 365–379, 2012.
- [30] Q. Wu, L. Wu, G. T. Liang, Q. X. Wang, T. Xie, H. Mei. Inferring dependency constraints on parameters for web services. In *Proceedings of the 22nd International Conference on World Wide Web*, ACM, Rio de Janeiro, Brazil, pp. 1421–1432, 2013.
- [31] A. Guha, S. Krishnamurthi, T. Jim. Using static analysis for Ajax intrusion detection. In *Proceedings of the 18th International Conference on World Wide Web*, ACM, Madrid, Spain, pp. 561–570, 2009.
- [32] A. Krishnamurthy, A. Mettler, D. Wagner. Fine-grained privilege separation for web applications. In *Proceedings of the 19th International Conference on World Wide Web*, ACM, Raleigh, USA, pp. 551–560, 2010.
- [33] A. Kieyzun, P. J. Guo, K. Jayaraman, M. D. Ernst. Automatic creation of SQL injection and cross-site scripting attacks. In *Proceedings of the 31st International Conference on Software Engineering*, IEEE, Washington DC, USA, pp. 199–209, 2009.
- [34] P. Saxena, S. Hanna, P. Poosankam, D. Song. Flax: Systematic discovery of client-side validation vulnerabilities in rich web applications. In *Proceedings of the 17th Annual Network and Distributed System Security Symposium*, NDSS, San Diego, USA, 2010.



Shuo Wen received the B.Sc. degree in computer science and technology from Nankai University, China in 2009. He is currently a Ph.D. degree candidate at the Institute of Machine Intelligence, College of Computer and Control Engineering, Nankai University, China.

His research interests include networking and distributed systems with a focus

on web applications and services and cloud computing.

E-mail: wenshuo@mail.nankai.edu.cn

ORCID iD: 0000-0001-6750-3735



Yuan Xue received the B.Sc. degree in computer science from Harbin Institute of Technology, China in 1998, and the M.Sc. and Ph.D. degrees in computer science from the University of Illinois at Urbana-Champaign, UK in 2002 and 2005. Currently, she is an assistant professor at the Department of Electrical Engineering and Computer Science of Vanderbilt University, USA. She is a NSF CAREER Award winner.

Her research interests include networking and distributed systems with a focus on wireless and mobile systems, web applications and services, clinical information system and cloud computing.

E-mail: yuan.xue@vanderbilt.edu



Jing Xu has been a professor of Nankai University in the Institute of Machine Intelligence, College of Computer and Control Engineering, Nankai University, China since 2006. She is a member of China Computer Federation, Software Engineering Technical Committee.

Her research interests include software engineering, software testing and information technology security evaluation.

E-mail: xujing@nankai.edu.cn (Corresponding author)

ORCID iD: 0000-0001-8532-2241



Li-Ying Yuan received the B.Sc. degree in computer science and technology from Nankai University, China in 2014. Currently, she is a master student at the Institute of Machine Intelligence, College of Computer and Control Engineering, Nankai University, China.

Her research interest is software analysis.

E-mail: yuanliying@mail.nankai.edu.cn



Wen-Li Song received the B.Sc. degree in computer science and technology from Nankai University, China in 2013. Currently, she is a master student at the Institute of Machine Intelligence, College of Computer and Control Engineering, Nankai University, China.

Her research interest include software analysis.

E-mail: wenli.song@foxmail.com



Hong-Ji Yang is a professor in Centre for Creative Computing, Bath Spa University, Bath, UK. He has taken part in many important international conferences, such as International Conference on Software Maintenance, the 8th IEEE Workshop on Future Trends of Distributed Computing Systems, the 26th Annual International Computer Software and Applications Conference, etc. He is also the leader of Software Evolution and Reengineering Group at the Software Technology Research Laboratory. He has become IEEE Computer Society Golden Core Member since 2010. Also, he is a member of EPSRC Peer Review College since 2003.

His research interests include software evolution, software engineering and creative computing.

E-mail: h.yang@bathspa.ac.uk



Guan-Nan Si received the Ph.D. degree from Nankai University, China in 2011. He is currently an assistant professor of Shandong Jiaotong University, China.

His research interests include software engineering and software evaluating technology.

E-mail: siguannan@163.com