



BLProM: A Black-Box Approach for Detecting Business-Layer Processes in the Web Applications

Mitra Alidoosti^a, Alireza Nowroozi^{b,*}, Ahmad Nickabadi^c

^aMalek-Ashtar University of technology, Tehran, Iran.

^bMalek-Ashtar University of technology, Tehran, Iran.

^cAmirkabir University of Tehran, Tehran, Iran.

ARTICLE INFO.

Article history:

Received: 24 May 2019

Revised: 22 December 2019

Accepted: 9 February 2020

Published Online: 3 May 2020

Keywords:

Business Layer, Business Process, Navigation Graph.

ABSTRACT

Web application vulnerability scanners cannot detect business logic vulnerabilities (vulnerabilities related to logic) because they are not able to understand the business logic of the web application. To identify the business logic of the web application, this paper presents BLProM, Business-Layer Process Miner, the black-box approach that identifies business processes of the web application. Detecting business processes of the web applications can be used in dynamic security testing to identify business logic vulnerabilities in web applications. BLProM first extracts the navigation graph of the web application then identifies business processes from the navigation graph. The evaluation conducted on three well-known open-source web applications shows that BLProM can detect business logic processes. Experimental results show that BLProM improves web application scanning because it clusters web application pages and prevents scanning similar pages. The proposed approach is compared to OWASP ZAP, an open-source web scanner. We show that BLProM improves web application scanning about %96.

© 2019 JComSec. All rights reserved.

1 Introduction

Most of the vulnerabilities reported in the Common Vulnerabilities and Exposures database [1] are related to the web application vulnerabilities. The number of security breaches increased by %35.5 in 2015 compared to last year [2]. Business logic vulnerabilities affect web application security as the most potent vulnerabilities.

So far, there are subtle vulnerabilities related to the web application logic that are still discovered

manually. Automated scanners cannot detect business logic flaws in applications because scanners are not able to understand the context [3]. Such vulnerabilities can only be detected through manual testing and be relied on tester creativity and skills [4].

There is no formal definition for business logic vulnerabilities [3]. It is very difficult to detect business logic vulnerabilities and this type of vulnerability causes serious damage in case of misuse [3]. Understanding context is difficult for automated tools, so penetration testers are responsible for detecting business logic vulnerabilities. Since business logic vulnerabilities are application-specific, it's hard to detect these types of vulnerabilities.

Web applications do not have any formal documentation describing their internal states and expected

* Corresponding author.

Email addresses: Alidoosti@mut.ac.ir (M. Alidoosti), Nowroozi@mut.ac.ir (A. Nowroozi), Nickabadi@aut.ac.ir (A. Nickabadi)

<https://dx.doi.org/10.22108/jcs.2020.117223.1028>

ISSN: 2322-4460 © 2019 JComSec. All rights reserved.



user behavior. The lack of such a document makes it difficult to detect business logic vulnerabilities. For example, adding a specific item several times in a shopping cart is a common feature but repeated usage of discount code is a kind of business logic vulnerabilities. A human easily understands the difference between these two scenarios, while a scanner without an appropriate model of the web application cannot distinguish between these two scenarios [4]. Researches [5–7] have been conducted to automatically detect business logic vulnerabilities but they are used for small applications. Also, the application source code is required to generate the appropriate model of the application [4]. So nowadays; an automatic tool that can find business logic vulnerabilities is required.

In this paper, we propose BLProM, a black-box technique for detecting web application business layer. BLProM By identifying business processes in the web application provides the ability to identify business layer vulnerabilities. In other words, to dynamic security testing of the web application in the business layer, it is necessary first to identify the business processes of the web application (detecting business logic of the web applications). Then, by analyzing the processes, the business layer vulnerabilities are identified. The BLProM output is the Web application business processes that are used as input in dynamic security testing in the business layer. The proposed approach is independent of the technology used in web applications and automatically finds business processes. Also, we will show that BLProM improves web application scanning, because it detects similar pages in the web application, and prevents scanning of similar pages. Comparing the results of the web application scanning, between BLProM and OWASP ZAP (web application open source scanner), shows that BLProM improves web application scanning by about 96%. In summary, this paper makes the following contributions:

- We present BLProM, a black-box technique for detecting web application business-layer.
- We present a new black-box approach for clustering web pages.
- We show that web application scanning improved about %96 by identifying web application business processes.

2 BACKGROUND and RELATED WORK

The business layer determines the business logic of the web applications. The business layer is responsible for data processing and data management and specifies business logic policies and rules. Besides, this layer validates the input data. Figure 1 shows the three-layer

architecture of a web application and the position of the business layer in the web application. The presentation layer is a user interface that displays data to the user and receives inputs from the user. In a web application, this is the part that receives the HTTP request and returns the HTML response.

The business layer handles data validation and business rules. The data access layer communicates with the database by constructing SQL queries.

After receiving data from the user, the data is available to the business layer. The web application uses the data to run business processes. Every business process has several steps that should be implemented respectively and processes may interact with each other.

The business layer specifies the logic of the web application. Business logic vulnerability is a defect in the business layer. A business logic attack vector is a legitimate request (usually multiple requests) and has legitimate input values that abuse a module's functionality to inflict damage and direct damage to the business.

2.1 Business Logic Attacks

There are two approaches to prevent business logic attacks: 1) identifying attacks at runtime (*defense approach*) and 2) identifying logic vulnerabilities in the *web* applications (*prevention approach*). In the *defense approach*, the behavior of the web application is monitored and an attack is reported whenever the web application exits from the normal state. In the *prevention approach*, attack vectors are used to identify business logic vulnerabilities.

BLOCK [9] and Swaddler [10] use a defense approach to prevent business logic attacks. BLOCK First obtains the behavioral model of the web application by observing the interaction of the clients and the web application. It extracts a set of constants from the request/response sequence and session variables. BLOCK Identifies any request or response that violates identified constants as an attack.

Swaddler provides an anomaly detection method for detecting attacks. It uses anomaly in the internal state of the web application to detect vulnerabilities. In other words, the web application's internal state is monitored in the learning phase and the normal values of the web application state are extracted which define the web application profile. Then in the detection phase, abnormal states are identified.

MiMoSA [5] uses a *prevention approach* in the form of a *white-box* approach to identify business logic vulnerabilities. MiMoSA first provides a web application model based on the web application's state and work-



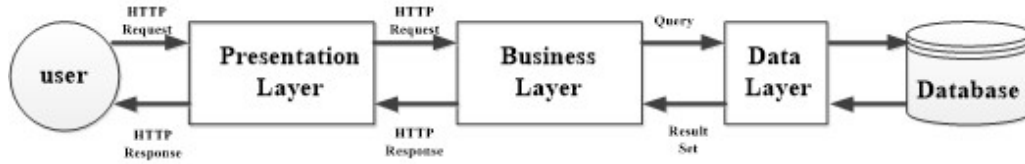


Figure 1. The Three-Layer Architecture of a Web Application [8].

flow. MiMoSA detects multi-step attacks by analyzing the relationship between the web application and the database, as well as the connections in the web application.

SENTINEL [11] and Pellegrino [3] use a *prevention approach* to identify business logic vulnerabilities in the form of a *black-box* approach.

SENTINEL [11] is a *black-box* approach for detecting logical weaknesses in database access. SENTINEL generates a state machine of the web and extracts a set of invariants from observed SQL queries and responses and session variables as the application specification. Any SQL query that violates defined invariants is identified as an attack.

Pellegrino et. al. [3] propose a *black-box* technique to detect logic vulnerabilities in web applications. This technique extracts behavioral patterns from network traces in which the user interacts with a certain application's functionality. First, the web application is modeled and then attack vectors are applied to the model.

Our previous works, BLDAST [12, 13] and BLTOCTTOU [14] use a prevention approach to identify business logic vulnerabilities in the form of a black-box approach. BLProM [15] can be used as an input for BLTOCTTOU and BLDAST.

BLDAST [12, 13] is a dynamic and a black-box vulnerability analysis approach that identifies business logic vulnerabilities of a web application against flooding DoS attacks. BLDAST assesses web application resiliency against flooding DoS attacks. It can take into account the business processes of a web application. BLDAST selects critical pages in business processes. A critical page has considerable response time. Therefore a critical process can enforce heavy load into the target and lead the web server to become unresponsive. The goal of the BLDAST is to find these critical processes within the web applications.

BLTOCTTOU [14] is a black-box dynamic application security tester for detecting business logic vulnerabilities against race condition attacks. BLTOCTTOU identifies vulnerabilities with the help of finding the business processes of the web application. BLTOCTTOU detects business processes that interact with each other; one process should set the value of a vari-

able and the other should read or write that variable. To identify the race condition, BLTOCTTOU first executes identified processes sequentially and then executes them in reverse order. At last, it evaluates the outputs of these two modes. If they are different, the web application is vulnerable to a race condition.

2.2 Clustering Web Pages

Crescenzi [16] presented an approach to cluster web pages based on the page structure. The structural similarity between web pages is defined by DOM trees of their hyperlinks. The final clusters are used to build a model that describes the structure of the site according to classes of pages and their connectivity.

3 BUSINESS-LAYER PROCESS MINER

In this paper, the BLProM is proposed to identify business processes of the web application and we use its outputs as the input in the web application security testing in the business layer. Then by analyzing the interaction between business processes, business layer vulnerabilities can be detected.

BLProM first preprocesses normal user HTTP traffic. Then extracts web application pages in the traffic. BLProM clusters similar pages to prevent the infinite growth of the user navigation graph. Detected clusters are graph nodes and the graph edges are the relations between detected clusters. Based on detected nodes and edges, the user navigation graph is extracted. Then BLProM extracts business processes from the navigation graph. BLProM has two main steps:

1. Extracting user navigation graph.
2. Detecting business processes in the web application.

Figure 2 shows the proposed steps to identify the business process in the web application. In the following, we will explain each of these steps in detail.

3.1 Extracting the User Navigation Graph

First, the normal user starts to crawl the web application. The traffic of a normal user is captured and stored. It should be noted that the user permission



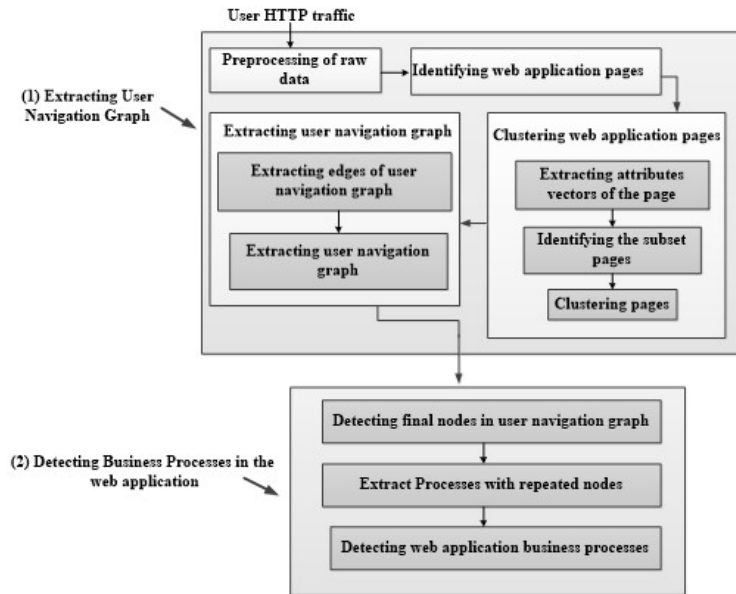


Figure 2. Black-Box Approach to Detect Web Application Business Process.

level can be different according to the role of the user and consequently the user navigation graph varies depending on the permission level. In this paper, the normal user has the user-level permission and searches through related permissible pages. The normal user crawls all different parts of the application that are allowed to search.

BLProM initially extracts the user navigation graph from the stored traffic; this is performed through steps as follows:

1. Preprocessing of raw input data
2. Identifying existing web application pages in the stored traffic
3. Clustering the web application pages
4. Extracting the user navigation graph

1) *Preprocessing of raw input data*

In the preprocessing step, the BLProM cleans the data and removes irrelevant data samples. In this paper, only HTTP requests and responses are used. For the responses, only those with successful status codes 200 and 209 are employed. The BLProM removes responses with failure status codes as well as their corresponding requests. Additionally, in this paper, only GET and POST requests are needed and the remaining ones are discarded.

2) *Identifying the web application pages in the stored traffic*

2) *Identifying the web application pages in the stored traffic*

Each page of the application can be represented as a pair (main request, corresponding response to the main request). To identify the web applica-

tion pages, it is first necessary to detect the main requests in the traffic. After identifying the main requests, the corresponding responses must also be identified.

Identifying the main HTTP requests in the traffic: In the user's stored traffic, there are both the main HTTP requests that lead to loading the web application pages and the secondary HTTP requests that are responsible to load a file, image, etc. of the page. The BLProM must distinguish between the main and secondary requests. In other words, when the main requests are identified, the remaining ones are considered as secondary requests. In this way, any request that its Referer header is different from the Referer of the previous request is considered as the secondary requests and the previous one is the main request. Additionally, the first request in the traffic is considered the main request because the first request does not include the Referer. It should be noted that the Referer header field in the HTTP request indicates the URL of the previous page the user visited.

Identifying the corresponding responses to the main requests: To identify the corresponding responses to the main requests, it is only needed to select the responses that their content-type field is text/html. Because the corresponding response to the secondary requests is often a file, photo, etc. while the corresponding response to the main requests is in the form of text and HTML. Such responses are the main responses in the HTTP traffic. After identifying the main requests and responses in the traffic, each pair (main request, corresponding responses to the main request) indicates a page of the web application. It



# Request URI	Referer	Response Content-Type
1 /osCommerce	/	text/html
2 /index.php?cPath=1	/osCommerce	text/html
3 /images/cdrm_drives.gif	/index.php?cPath=1	image/jpeg
4 /images/graphic_card.gif	/index.php?cPath=1	image/jpeg
5 /images/keyboards.gif	/index.php?cPath=1	image/jpeg
6 /index.php?cPath=1_17	/index.php?cPath=1	text/html
7 /css/stylings.php	/index.php?cPath=1_17	text/css
8 /index.php?cPath=1_4	/index.php?cPath=1_17	text/html

Figure 3. Number of Consecutive Requests.

should be noted that in responses of the main requests, all secondary requests exist that lead to load the web page.

Identifying whether the last request in the traffic is the main request or not: If the last response in the traffic is the main response, the last HTTP request is the main request as well.

For example, in Figure 3, requests 1, 2, 6, and 8 are the main requests that are shown in red.

The pseudocode of the algorithm used by the BLProM to identify web application pages is shown in Algorithm 1. As mentioned before, each web application page can be indicated by a pair (main request, corresponding responses to the main request). The pseudocode in Algorithm 1 can be divided into four main parts:

1. Extracting main requests in the traffic (line 13-26)
2. Extracting corresponding responses to the main requests (line 27-33)
3. Identifying whether the last request in the traffic is the main request or not (line 34-37)
4. Extracting web application pages (lines 38-43)

In line 9, first, all requests in the traffic either the main requests or secondary ones are extracted and put in the *HTTPRequest* variable. In line 10, existing responses in the traffic are added to the *HTTPResponse* variable. In lines 11 and 12, the total number of requests and responses is calculated. In line 15, it is checked whether the current request is the first request, if yes, it is considered as the main request and included in the *MainRequest* variable. In line 21, it is checked whether the *Referer* field of the current request is different from the *Referer* field of the previous request, if yes, the previous request is considered as the main request.

In line 29, it is checked whether the content-type field of the current response is text/html. If the given condition is met, the current request is considered as the main request.

In line 35, it is checked whether the last HTTP response is the main response, if yes, the last request is considered as the main request as well.

In line 41, all main requests and responses are the web application pages which are shown as pairs (request, response).

3) Clustering the web application pages

In this step, BLProM clusters the web application pages. Clustering aims to put similar pages in the same cluster. This is helpful to prevent the infinite growth of the user navigation graph. The pages in the same cluster are similar to each other.

At this stage, all pages in the user's stored traffic have been extracted as pairs (main request, corresponding response to the main request). Each pair (main request, corresponding response to the main request) shows one page of the application. To extract the optimal user navigation graph, the similar extracted pages must be identified and clustered. In the user navigation graph, nodes indicate the application's unique pages and edges represent the link between the pages. To identify similar pages, a criterion should be considered by the purpose of the clustering. The type of operations that the user can perform on the page is considered as a measure for the separation of pages. In other words, two pages are similar if the user can perform the same operations on them. For example, consider two pages such that both contain only a button but the title of the buttons is different where in the first page the title is "continue" and for the second one, it is "save". These two pages are different because the user performs different operations on them. Thus, according to the criteria specified for the similar pages, the following pages are considered similar in web applications:

- In the online shop application, if the pages related to the shopping cart of goods even if they contain different items, they are considered similar pages.
- In the online shop application, the profile pages of each product are similar. because they have the same HTML structure in terms of the important HTML elements.
- Pages that display search results with different keywords are similar.
- If the structure of pages is a subset of another page in terms of the important HTML elements, these pages are considered similar.
- Two pages that both contain user comments are considered similar even if the contents are about different products.



Algorithm 1 The Pseudocode for Extracting Web Application Pages From the Traffic**INPUT:** HITRtraffic: {*HttpMessage*₁, *HttpMessage*₂, *HttpMessages*₃, ..., *HithMessage*_n}**OUTPUT:** WebPages as a set of (*Request*_i, *Response*_i)

```

1: Begin
2: let JWebPages = {}; // set of web pages
3: let HttpRequest= {}; //set of HTTP Requests
4: let HttpResponse = {}; // set of HTTP Responses
5: let MainRequest= {}; // set of main HTTP Requests
6: let MainResponse = {}; // set of main HTTP Responses
7: let i, k= 1; // counter for current HttpMessage
8: let LastReferer =  $\emptyset$  , NewRefere =  $\emptyset$ 
9: HttpRequest=ExtractReg(HITRtraffic); //extract HTTP Requests from HTTP traffic
10: HttpResponse - ExtractResp(HITRtraffic); //extract HTTP Response from HTTP traffic
11: n = extractNumber(HttpRequest) //extract total number of HTTP Requests
12: m = extractNumber(HttpResponse) //extract total number of HTTP Responses
13: // extract Main HTTP Request from HTTPRequests
14: for i: 1 ... n do
15:   if (i = 1) then
16:     add MainRequest  $\leftarrow$  HttpRequesti // First Request is a Main Request
17:     LastReferer  $\leftarrow$  Referer of HttpRequesti;
18:   else
19:     NewRefere  $\leftarrow$  Referer of HttpRequest1;
20:   end if
21:   if (NewReferer  $\neq$  LastReferer) then
22:     add MainRequest  $\leftarrow$  HttpRequesti-1;
23:     LastReferer  $\leftarrow$  NewRefere;
24:   end if
25: end for
26: //extract Main HTTP Response from HTTPResponse
27: for k : 1 ... m do
28:   if (content-type of HTTPResponsek = text/html) then
29:     addMainResponse  $\leftarrow$  HttpResponsek
30:   end if
31: end for
32: // checking last request is main request or not
33: if (HTTPResponsem is in MainResponse) then
34:   addMainRequest  $\leftarrow$  HTTPRequestm
35: end if
36: //extract set of web pages
37: size = extractNumber(MainRequest) //extract total number of Main HTTP Requests
38: for j: 1 ... size do
39:   WebPages  $\leftarrow$  (MainRequestj)MainResponsej;
40: end for
41: return WebPages;
42: end

```



Table 1. Attribute Vector of a Page in osCommerce Web Application.

inputs	null
Buttons	html.body.button.Reviews# html.body.div.div.div.div.form.div.div.span.s pan.button.Add to Cart
anchors	null
image	html.body.div.div.div.a.img

```

<html>
  <title> Hello! </title>
  <body>
    <p>
      <button value= "continue"> Continue</button>
    </p>
  </body>
</html>

```

Figure 4. An Example of HTML Code.

Definition of Document Object Model (DOM) path for an HTML element: DOM path of an element is the position of the element in the HTML code.

For example, in Figure 4 the DOM path of the button (DOM_{button}) is Document.Html.Body.P.Button.

Definition 1. [Similar Pages] the similar pages are those the user can perform the same operations on them and are identical in terms of the position of the important HTML elements in the page. The important HTML elements in the page include buttons, images, inputs, and anchors.

The clustering process includes three steps:

1. Extracting attributes vectors of the page
2. Identifying the subset pages
3. Clustering pages

In the following, these steps are discussed in detail.

1. Extracting attributes vectors of the page

The BLProM shows each page of the application as a pair (main request, response). In this step, BLProM extracts the corresponding attributes vectors of each page by applying a data mining operation on the above pair. The BLProM models each page using the following attribute vector:

WebPages = the total pages in an application
 $\forall w \in \text{WebPages}$

$$w = (DOM_{inputs}, DOM_{buttons}, DOM_{anchors}, DOM_{imgs})$$

$$DOM_{inputs} = \prod_i^n DOM(input_i)$$

$$DOM_{buttons} = \prod_i^n DOM(buttons_i)$$

$$DOM_{anchors} = \prod_i^n DOM(anchor_i)$$

$$DOM_{imgs} = \prod_i^n DOM(img_i)$$

- $DOM(input)$: DOM path of $\langle input \rangle$ tag in the page + the value of type attribute in the $\langle input \rangle$ tag + the value of name attribute in the $\langle input \rangle$ tag (in the absence of name attribute, the value attribute is considered).
- $DOM(button)$: DOM path of the button in the page + the title of button.
- $DOM(anchor)$: DOM path of $\langle a \rangle$ tag in the page.
- $DOM(img)$: DOM path of the existing image in the page.

Suppose the web application page contains several buttons; in this case, the second element of the page attribute vector is a set of the DOM paths of buttons in the page that are separated by "#". Figure 5 shows one of the osCommerce¹ web application pages. Table 1 shows the attribute vector of the page in Figure 5. As shown, the input element and the anchor elements are null, it means the page does not contain the above tags.

2. Identifying similar pages

After extracting the attribute vector of each page, it is necessary to identify similar pages. Those pages that their attribute vectors are a subset of another page or have fully similar attribute vectors are considered as similar pages.

According to Definition 1, the attribute vector of each web application page has four elements. The attribute vector of page 1 is considered the same as the attribute vector of page 2 if:

- All vector elements of page 1 equal to corresponding elements in the vector of page 2.
- All vector elements of page 1 are a subset of corresponding elements in the vector of page 2.
- All vector elements of page 2 are a subset of corresponding elements in the vector of page 1.

¹ <https://www.oscommerce.com/>



Microsoft IntelliMouse Explorer

[MSIMEXP]

Microsoft introduces its most advanced mouse, the IntelliMouse Explorer! IntelliMouse Explorer features a sleek design, an industrial-silver finish, a glowing red underside and taillight, creating a style and look unlike any other mouse. IntelliMouse Explorer combines the accuracy and reliability of Microsoft IntelliEye optical tracking technology, the convenience of two new customizable function buttons, the efficiency of the scrolling wheel and the comfort of expert ergonomic design. All these great features make this the best mouse for the PC!

Available Options:

Model:

PS/2

Reviews

\$64.95



Add to Cart

Figure 5. A Page of osCommerce Application.

Algorithm 2 The Pseudocode for Identifying the Similar Pages

INPUT: $w_1 = (DoM_{w_1}(input), DoM_{w_1}(button), DoM_{w_1}(anchor), DoM_{w_1}(img))$,

$w_2 = (DoM_{w_2}(input), DoM_{w_2}(button), DoM_{w_2}(anchor), DoM_{w_2}(img))$

OUTPUT: Boolean flag // true means two pages are the same

```

1: Begin
2: Let flag, input, button, anchor, img=false;
3: if  $DoM_{w_1}(input) \subseteq DoM_{w_2}(input)$  or  $DoM_{w_2}(input) \subseteq DoM_{w_1}(input)$  then
4:   input=true;
5: end if
6: if  $DoM_{w_1}(button) \subseteq DoM_{w_2}(button)$  or  $DoM_{w_2}(button) \subseteq DoM_{w_1}(button)$  then
7:   button=true;
8: end if
9: if  $DoM_{w_1}(anchor) \subseteq DoM_{w_2}(anchor)$  or  $DoM_{w_2}(anchor) \subseteq DoM_{w_1}(anchor)$  then
10:  anchor=true;
11: end if
12: if  $DoM_{w_1}(img) \subseteq DoM_{w_2}(img)$  or  $DoM_{w_2}(img) \subseteq DoM_{w_1}(img)$  then
13:  img=true;
14: end if
15: if input and button and anchor and img then
16:  flag=true;
17: end if
18: return flag;
19: end

```

- If one or more vector elements of page 1 are a subset of their corresponding elements in the vector of page 2, the rest of the vector elements of page 1 must be the same with their corresponding elements in the vector of page 2.
- The null element is a subset of every element.

Similar pages are identified according to the above-mentioned attributes. Algorithm 2 illustrates the pseudocode for identifying similar pages.

3. Clustering web application pages

After identifying the similar pages, they are

put in the same cluster. The pages in a cluster are similar to each other and refer to a unique page of the application. Algorithm 3 shows the pseudocode for clustering web pages. In line 7, it is checked whether two pages w_i and w_j are the same, if yes, they are put in the same cluster.

4) Extracting user navigation graph

In this step, BLProM connects the obtained clusters that each one represents a unique web application page. Each cluster has a set of similar pages, each of these pages has URI and Referer field. Thus, each cluster contains a set of URIs and a set of Referers for the pages in that cluster. It should be noted that the Referer field is the URI



Algorithm 3 The Pseudocode for Clustering the Web Application Pages

INPUT: WebPages = $\{w_1, w_2, w_3, \dots, w_n\}$
OUTPUT: the web application model M as a set of web page clusters C;
the web page clusters C as a set of web pages;

- 1: Begin
- 2: Let $M = \emptyset$; // set of page clusters
- 3: Let $k=1$; // number of web page clusters
- 4: **for** $i : 1 \dots n$ **do**
- 5: $C_k \leftarrow w_i$
- 6: **for** $j = i + 1 \dots n$ **do**
- 7: **if** *SimilarWebPages*(w_i, w_j) **then**
- 8: WebPages \leftarrow WebPages $- w_i$;
- 9: $n = \text{length of WebPages}$;
- 10: $C_k \leftarrow w_j$
- 11: **end if**
- 12: **end for**
- 13: $k++$;
- 14: **end for**
- 15: return C;
- 16: end

of the previous web application page that the user visited. For extracting the edges of user navigation graph, the produced clusters are checked to find which Referer set of clusters has the intersection with the URI set of the cluster. When the cluster is found, these two clusters are connected. Suppose that the URI set of cluster C1 has an intersection with the Referer set of cluster C2, then the path from cluster C1 to cluster C2 ($C1 \rightarrow C2$) is created. In other words, the edge C1C2 is one of the edges in the user navigation graph. Algorithm 4 shows the pseudocode for extracting the edges from the user navigation graph. The URI set and the Referer set of each cluster are respectively obtained according to lines 5 and 6 of the pseudocode in Algorithm 3. In line 10, if the intersection of the URI set of each cluster with other clusters' Referer set is not null, the CiCj edge is added to the edge set of the graph.

In this step, BLProM connects the obtained clusters that each one represents a unique web application page. Each cluster actually has a set of similar pages, each of these pages has URI and Referer field. Thus, each cluster contains a set of URIs and a set of Referers for pages in the cluster. It should be noted that the Referer field is actually the URI of the previous web application page that the user visited. For extracting the edges of user navigation graph, the produced clusters are checked to find which Referer sets of clusters has the intersection with the URI set of the cluster, when the cluster is found, these two clusters are connected. Suppose that the URI set of cluster C1 has intersection with the Referer set of cluster C2, then the path from cluster C1 to the cluster C2 ($C1$

$\rightarrow C2$) is created. In other words, the edge C1C2 is one of the edges in the user navigation graph. Algorithm 4 shows the pseudocode for extracting the edges from the user navigation graph. The URI set and the Referer set of each cluster are respectively obtained according to the lines 5 and 6 of the pseudocode in Algorithm 3. In line 10, if the intersection of URI set of each cluster with other clusters' Referer set is not null, CiCj edge is added to the edge set of the graph. The user navigation graph is created according to the algorithm in Algorithm 5, where nodes are the clusters set and the identified edges are the path between clusters. In the created graph, each node indicates the unique page and the edges show the path between pages.

Definition 2. [The User Navigation Graph] This graph is shown with tuple $\langle C0, C, E \rangle$ where C is the set of nodes in the graph, $C0 \in C$ is the first (initial) node in the graph and $E \subseteq C \times C$ is the set of edges in the graph.

Algorithm 5 shows the pseudocode for extracting the user navigation graph. Line 7 indicates a cluster that contains the initial page of the application. It is considered as the initial node of the graph.

3.2 Identifying Business Processes in the Application

To identify the web application business processes in the user navigation graph, it is necessary to define the process and final node and then define the business process.

Definition 3. [The Application Process (P)] The



Algorithm 4 The Pseudocode for Extracting Edges From the User Navigation Graph

INPUT: $C = \{C_1, C_2, C_3, \dots, C_k\}$ //web pages clusters as graph nodes
 WebPages= $\{w_1, w_2, w_3, \dots, w_n\}$ //set of web pages
OUTPUT: the web application graph edges E as a set of edges

```

1: Begin
2: Let  $E = \emptyset$ ; // set of web application graph Edges
3: for  $j : 1 \dots k$  do
4:   Let  $URI_{cj}, Referer_{cj} = \emptyset$ ;
5:    $URI_{cj} = URI_{cj} \cup \text{ExtractURI}(w)$  for any  $w \in C_j$ 
6:    $Referer_{cj} = Referer_{cj} \cup \text{ExtractReferer}(w)$  for any  $w \in C_j$ 
7: end for
8: for  $i : \dots k$  do
9:   for  $j : i + 1 \dots k$  do
10:    if  $(URI_{ci} \cap Referer_{cj} \neq null)$  then
11:       $E \leftarrow E + C_i C_j$ 
12:    end if
13:  end for
14: end for
15: return E;
16: end

```

Algorithm 5 The Pseudocode for Extracting the User Navigation Graph

INPUT: $C = \{C_1, C_2, C_3, \dots, C_k\}$ //web pages clusters as graph nodes
 w_1 // First web page
OUTPUT: the web application navigation graph $\langle C_0, C, E \rangle$

```

1: Begin
2: Let  $C_0 = \emptyset$ 
3:  $E = \text{ExtractGraphEdges}$ ; //set of web application graph Edges
4: for  $j : 1 \dots k$  do
5:   if  $C_k$  contains  $w_1$  then
6:      $C_0 = C_0 + C_k$ 
7:   end if
8: end for
9: return  $C_0, C, E$ ;
10: end

```

process P in the application is a sequence of nodes and edges in the user navigation graph like E_1, E_2, \dots, E_k , where $E_i \in E$, and $E_i = C_{i-1} C_i$.

Algorithm 6 shows the pseudocode for extracting processes.

Definition 4. [the final nodes in the user navigation graph (F)] The final nodes refer to the completion of a business process that occurs when the application reaches there.

The final nodes can be detected by examining the HTTP responses. For example, in the process of buying a product, a phrase like "Thank you for your purchase" is displayed after completion of the purchase. By specifying a set of these phrases and searching them in the responses, the final nodes can be identi-

fied. Some keywords used for identifying the final node are Thank, Congratulations, Successfully, Log Off and Search Results. Additionally, some buttons in the web application page are good signs that help the identification of the final node in the graph. The examples of final nodes include the page after clicking the "save" button, the page after clicking the "creation" button and the page after clicking the "submit" button.

Definition 5. [the application business process] The business process BP in the application is a process that has at least one of the following conditions:

1. The first node of the process is the initial node in the user navigation graph (C_0) and the process end node is the final node in the user navigation graph (F).
2. If the process passes its first node again, it means



Algorithm 6 The Pseudocode for Extracting Processes**INPUT:** the web application first node C_0 the web application Graph edges E **OUTPUT:** the web application graph process P as a set of web application process

```

1: Begin
2: Let  $P_0 = \emptyset$ ; // set of web application processes
3: Let StartEdge=  $\emptyset$ ; // set of graph edge that begin from  $C_0$ 
4: StartEdge= ExtractGraphEdges( $C_0$ ) //Extract all edges from  $C_0$ 
5: EndPoint = ExtractEndPoint(StartEdge) //Extract the end point of edges
6: if ( $ExtractProcess(EndPoint, E) \neq null$ ) then
7:   return  $P=E+ExteractProcess(Endpoint, E)$ ; for any edge  $E \in StartEdge$ 
8: else
9:   return  $E$ ;
10: end if
11: end

```

the first node and the end node of the process are the same and the process length is greater than two. (If there is a return to the passed node in the process and the created loop length is more than two, this is a business process).

All processes in the graph from the initial node (the application initial page) to the identified final nodes, as well as the processes of their initial node and the final node are the same; and all of them are the application business processes. Algorithm 7 shows the pseudocode for identifying the application business process. In line 4, the application business processes are extracted. In line 5, the final nodes of the graph are extracted. In line 7, the processes that start with the initial node and end with the final node are identified as the business process are stored in the variable BP. In line 9, the processes with repeated nodes are detected and in line 11, among the detected processes, if their initial node and their final node are the same and their length is greater than two, they are added to the variable BP as the business process.

4 EXPERIMENTAL RESULTS

The testbed used in this section is a network consisting of a web server (test target) and two clients (BLProM system and legal user). The web server and the clients are loaded on a virtual machine. The web server and the clients' profiles are shown in Table 2.

The web applications listed in Table 3 are installed on the web server (test target) and then we plan to identify the business layer of the web applications.

The legal user first starts using the selected web applications. The user crawls all permitted parts of the web application. HTTP traffic of the legal user is given to BLProM as its input.

5 EVALUATION

BLProM's goal is to identify the business layer of the web application. We can identify business logic vulnerability by identifying the business layer of the web application. BLProM detects the business processes of the web application. Identifying business processes is the main step in dynamic security testing of the web application in the business layer.

We compare BLProM with OWASP ZAP. The OWASP Zed Attack Proxy (ZAP) is a free web scanner. It scans web applications and automatically finds some security vulnerabilities. ZAP is the only free web scanner that has API for extracting web application graph. The web application pages are graph nodes and the relations among pages are shown as graph edges. The main difference between BLProM and ZAP is in detecting similar pages. ZAP cannot detect similar pages in the web application but BLProM can. ZAP's graph only shows the relation among scanned pages but BLProM generates the optimal graph. About the accuracy of the generated graph, both BLProM and ZAP are the same.

To evaluate the proposed approach, we first show the accuracy of clustering by the following criteria:

- True Positive: Samples that fit well into their correct clusters.
- False Positive: Samples that fit in a cluster that do not belong to that cluster.
- False Negative: Samples that do not fit in a cluster but they belong to that cluster.
- Recall: It is calculated by the following formula:

$$recall = \frac{TruePositive}{TruePositive+FalseNegative}$$

- Precision: It is calculated by the following formula:



Table 2. Testbed Profiles.

Web server (test target)	CPU: Pentium dual core-2.20 GHZ OS: windows 8.1 VMware cpu: 1GHZ VMware RAM: 1G VMware OS: windows 7
Client (BLProM machine)	CPU: Intel corei7 2.20 GHZ OS: windows 8.1 VMware cpu: 1GHZ VMware RAM: 1G VMware OS: windows 7
Client (legal user)	CPU: Pentium dual core i3-3210 GHZ OS: windows 7 VMware cpu: 1GHZ VMware RAM: 1G VMware OS: windows 7

Table 3. Selected Web Applications for Evaluation.

Web application	Description
TomatoCart-1.1.8.6.1	e-commerce
osCommerce-2.3.4	e-commerce
WackoPicko	Web application for Sharing picture

Algorithm 7 The Pseudocode for Identifying the Application Business Processes**INPUT:** the web application navigation graph $\langle C_0, C, E \rangle$

the web application Graph edges E

OUTPUT: the web application graph business process PB as a set of web application process

```

1: Begin
2: Let  $P = \emptyset$ ; // set of web application processes
3: Let  $F = \emptyset$ ; //set of web application final nodes
4: P= ExtractProcess; //Extract processes in the web application
5: F= ExtractFinalNodes; //Extract Final nodes in the web application
6: for  $i : 1 \dots k$  do
7:   if  $P_k$  start by  $C_0$  and ends by F then
8:      $BP = BP + P_k$ 
9:   end if
10: end for
11: R=ExtractProcessWithRepeatedNodes(P); //Extract Process with Repeated Nodes
12: for  $j : 1 \dots m$  do
13:   if  $R_j$  has same initial and end node and  $\text{length}(R_j) \geq 2$  then
14:      $BP = BP + R_j$ 
15:   end if
16: end for
17: return BP;
18: end

```



Table 4. The Clusters of Selected Web Application Pages Evaluation.

Criteria	Web application	WackoPicko	Tomatocart	osCommerce
	#samples	89	150	210
	#clusters	29	66	40
	true positive	65	146	205
	false positive	24	4	5
	false negative	23	3	4
	recall	0.74	0.98	0.98
	precision	0.73	0.97	0.98
	f-measure	0.73	0.98	0.98

Table 5. Comparing the Proposed Approach With OWASP ZAP in Scanning WackoPicko.

Criteria	WackoPicko			
	Approaches	Proposed approach	OWASP ZAP	Percentage of improvement compared to OWASP ZAP
	# HTTP Message	89	89	–
	# Graph Nodes	22	89	75.2
	#Graph Edges	48	270	82.2
	# process (P)	12	48	75
	Average edge in each process (\bar{E})	4	46	91.3
	Average edges in all processes ($P*\bar{E}$)	48	2208	97.8
	#business processes	10	NA	–

$$precision = \frac{TruePositive}{TruePositive+FalsePositive}$$

- F-Measure: It is calculated by the following formula:

$$F - Measure = \frac{2*recall*precision}{TruePositive+FalsePositive}$$

The value of these criteria for the ClusteringWebPages algorithm (Algorithm 3) is shown in Table 4. In the first row, #samples shows the total number of web pages in HTTP traffic extracted from the ExtractWebPages algorithm in Algorithm 1. In the second row, #clusters shows the total number of clusters extracted from the ClusteringWebPages algorithm in Algorithm 3. In the next rows, the criteria listed above are calculated for each web application.

To evaluate the proposed approach, we compare the output of BLProM with the scanning output of OWASP ZAP for selected web applications. BLProM output for WackoPicko is shown in Table 5. #graph nodes shows the total number of clusters extracted from the ClusteringWebPages algorithm in Algorithm 3. #graph edges is the total number of edges extracted from the ExtractGraphEdges algorithm in

Algorithm 4. #process is the total number of paths from the first node extracted from the ExtractProcess algorithm in Algorithm 6. Average edge in each process (\bar{E}) is calculated by the sum of edges of each process divided by the total number of processes. Average edge in all process is calculated by the product of the total number of processes (P) in average edges of each process (\bar{E}). #business processes is the total number of business processes in the web application.

Existing values in the first column of Table 5 are the output obtained from the proposed approach. The values in the second column are the scanning output of ZAP. The third column shows the percentage of scanning improvement of our proposed approach compared to the ZAP scan. The results of the table indicate that the ZAP scan is a non-smart scan. As a result of this non-intelligent scan, ZAP is not able to identify business layer vulnerabilities.

BLProM output for osCommerce is shown in Table 6. The third column shows the percentage of scanning improvement compared to ZAP scanning. It is observed that web application scanning is improved by identifying the web application business layer.



Table 6. Comparing the Proposed Approach With OWASP ZAP in Scanning osCommerce.

osCommerce			
Criteria \ Approaches	Proposed approach	OWASP ZAP	Percentage of improvement compared to OWASP ZAP
# HTTP Message	170	170	–
# Graph Nodes	40	170	76.4
#Graph Edges	66	379	82.5
# process (P)	23	17	26
Average edge in each process (\bar{E})	3	113	97.3
Average edges in all processes ($P*\bar{E}$)	69	1921	96.4
#business processes	18	NA	–

Table 7. Comparing the Proposed Approach With OWASP ZAP in Scanning TomatoCart.

TomatoCart			
Criteria \ Approaches	Proposed approach	OWASP ZAP	Percentage of improvement compared to OWASP ZAP
# HTTP Message	150	150	–
# Graph Nodes	66	150	56
#Graph Edges	87	410	78.7
# process (P)	31	39	20.5
Average edge in each process (\bar{E})	4	101	96
Average edges in all processes ($P*\bar{E}$)	156	3131	95
#business processes	30	NA	–

BLProM output for TomatoCart is shown in Table 7. The third column shows the percentage of scanning improvement of our proposed approach compared to ZAP scanning.

Table 8 shows the average of the proposed approach, the average of OWASP ZAP and the average percentage of improvement in the scanning of selected web applications. For example, in the "Average edges in all processes" benchmark, our approach has been improved by about 96 percent compared to OWASP ZAP.

According to the results presented in this table, can be observed that BLProM has improved web application scanning. BLProM is aware of web application business processes. By identifying the web application business layer, web scanners can detect business layer vulnerabilities.

6 CONCLUSION

Business logic vulnerabilities are strong vulnerabilities that compromise web application security. Web

scanners cannot detect business logic vulnerabilities because they are unable to understand business logic of the web application. For detecting business logic vulnerabilities, the business logic of the web application needs to be understood. Therefore, these vulnerabilities are specific to the application and difficult to identify.

In this paper, we proposed BLProM, a black-box approach for detecting business processes of the web application. BLProM aims to ease detecting business logic vulnerabilities. BLProM output is used as input in dynamic security testing of the web applications in the business layer in order to detect business logic vulnerabilities. BLProM consists of two main steps:

- 1- extracting user navigation graph
- 2- Detecting web application business processes

At the lab, we scanned three web applications by BLProM and OWASP ZAP, an open-source web application. We showed that BLProM improved scanning about %96 compared to OWASP ZAP. BLProM improved the scanning of web applications because it clusters web pages and prevents scanning similar web



Table 8. Comparison of the Average of the Proposed Approach and the Average of ZAP Approach in the Scanning of Selected Web Applications.

Criteria	Average of selected web application		
	Proposed approach	OWASP ZAP	Percentage of improvement compared to OWASP ZAP
# Graph Nodes	42.6	136.3	69.2
#Graph Edges	67	353	81.1
# process (P)	20	36.6	40.5
Average edge in each process (\bar{E})	3.6	86.6	94.8
Average edges in all processes ($P*\bar{E}$)	91	2420	96.4

application pages.

References

- [1] Common vulnerabilities and exposures. <https://cve.mitre.org/cve/cve.html>.
- [2] 2015 ITRC. Identity Theft Resource Center Breach Report Hits Near Record High in 2015. <http://www.idtheftcenter.org/ITRC-Surveys-Studies/2015databreaches.html>, Accessed Feb, 2017.
- [3] G. Pellegrino and D. Balzarotti. Toward Black-Box Detection of Logic Flaws in Web Applications. In *Network and Distributed System Security symposium 2014 (NDSS2014)*, 2014. doi:10.14722/ndss.2014.23021.
- [4] Testing for business logic, OWASP. https://www.owasp.org/index.php/Testing_for_business_logic, Accessed Feb, 2017.
- [5] D. Balzarotti, M. Cova, V. Felmetzger, and G. Vigna. Multi-module vulnerability analysis of web-based applications. In *Proceedings of the 14th ACM conference on Computer and communications security*, page 25–35. ACM, 2007. doi:10.1145/1315245.1315250.
- [6] A. Doupé, B. Boe, C. Kruegel, and G. Vigna. Fear the EAR: discovering and mitigating execution after redirect vulnerabilities. In *Proceedings of the 18th ACM conference on Computer and communications security*, page 251–262. ACM, 2011. doi:10.1145/2046707.2046736.
- [7] V. Felmetzger, L. Cavedon, C. Kruegel, and G. Vigna. Business Logic Attacks – Bots and BATs, OWASP, 2009. In *USENIX Security Symposium*, 2010.
- [8] E. Chai. Business Logic Attacks – Bots and BATs, OWASP, 2009. https://www.owasp.org/images/a/aa/OWASP_Cincinnati_Jan_2011.pdf, Accessed Feb. 2017.
- [9] X. Li and 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*, page 247–256, 2011. doi:10.1145/2076732.2076767.
- [10] M. Cova, D. Balzarotti, V. Felmetzger, and G. Vigna. Swaddler: An Approach for the Anomaly-Based Detection of State Violations in Web Applications. In *Giovanni*, pages 63–86. Springer, Berlin, Heidelberg, 2007. ISBN 978-3-540-74319-4. doi:10.1007/978-3-540-74320-0_4.
- [11] X. Li, W. Yan, and Y. Xue. SENTINEL: securing database from logic flaws in web applications. In *Proceedings of the second ACM conference on Data and Application Security and Privacy*, page 25–36, 2012. doi:10.1145/2133601.2133605.
- [12] M. Alidoosti and A. Nowroozi. BLDAST: business-layer Dynamic Application Security Tester of the web application in order to detect web application vulnerabilities against flooding DoS attacks. In *Iran Society of Cryptology Conference, shiraz, Iran*. in Persian, 2017.
- [13] M. Alidoosti, A. Nowroozi, and A. Nickabadi. Evaluating the Web-Application Resiliency to Business-Layer DoS Attacks. *ETRI Journal*, 2019. doi:10.4218/etrij.2019-0164.
- [14] M. Alidoosti and A. Nowroozi. BLTOCTTOU: business-layer dynamic application security tester of the web application in order to detect web application vulnerabilities against Race Condition attacks. In *Computer Society of Iran Conference, Tehran, Iran*. in Persian, 2018.
- [15] M. Alidoosti and A. Nowroozi. BL-ProM: Business-layer process miner of the web application. In *ISCISC*, pages 1–6. IEEE, 2018. ISBN 978-1-5386-7582-3. doi:10.1109/ISCISC.2018.8546899.
- [16] V. Crescenzi, P. Merialdo, and P. Missier. Clustering Web pages based on their structure. *Data & Knowledge Engineering*, 54(3):279–299, 2005. doi:10.1016/j.datak.2004.11.004.





Mitra Alidoosti received her B.S. and M.Sc. degrees in computer engineering from the Department of Computer Engineering, Iran University of Science and Technology, Tehran, Iran, in 2009 and 2012, respectively. Currently, she is working toward the Ph.D. degree in computer engineering at Malek-e-Ashtar University of Technology, Tehran, Iran. Her research interests are computer network security, VoIP and SIP security, and web-application security.



Alireza Nowroozi is a freelance consultant who advises government and private-sector-related industries on information technology. He has four-year experience as an academic staff member and an IT post-doctoral position with Sharif University of Technology, Tehran, Iran. He is a specialist in artificial intelligence, cognitive science, software engineering, and IT security. Besides, he is a co-founder of four IT startups.



Ahmad Nickabadi received his B.S. degree in computer engineering in 2004, and the M.Sc. and Ph.D. degrees in artificial intelligence in 2006 and 2011, respectively, from Amirkabir University of Technology, Tehran, Iran. He is currently an Assistant Professor in the Department of Computer Engineering, Amirkabir University of Technology. His research interests include statistical machine learning and soft computing.

