



An empirical test of the accuracy of an attack graph analysis tool

| | |
|------------------|--|
| Journal: | <i>Information and Computer Security</i> |
| Manuscript ID: | IMCS-06-2014-0036.R2 |
| Manuscript Type: | Original Article |
| Keywords: | Attack graphs, Computer Networks, Computer Security, Security, Assessments |
| | |

SCHOLARONE™
Manuscripts

Review

1 Introduction

Securing computer networks is a complicated and difficult task. Computer networks in today's enterprises often consist of many hosts. These hosts run many different operating systems and software applications of different versions and configurations. Experience suggests that, in a typical enterprise, a considerable portion of these hosts will contain publicly known vulnerabilities that may be exploited to obtain user and host privileges in the organization. However, not all vulnerabilities are equally important to remediate. Some can be used to provide the attacker with all of the privileges of a host (e.g., root access on a Linux machine), and some can only be used to provide or impact a subset of them (e.g., reading certain parts of the memory in a machine). Furthermore, vulnerabilities may require certain preconditions to be met to be exploitable. For example, it might be required that the attacker be able to interact with the machine physically or that the attacker already possess the credentials of some user in the network. As a result of factors such as these, some vulnerabilities may be difficult for an outsider to exploit on public networks, and some may be easily exploited. Additionally, the successful exploitation of some vulnerabilities may provide access that makes it possible to exploit other vulnerabilities that are not directly exploitable from public networks. Because of these contingencies, analyzing how to prioritize remediation options or determine the present risks can become overwhelmingly complex.

Attack graphs have been designed to assist decision makers in this analysis. According to Heberlein et al. (2004), "*one of the primary focuses of the attack graph efforts is to identify how an adversary can chain together vulnerability exploitation to increase his capability.*" An attack graph aims to answer questions such as the following: To which hosts can an attacker on the Internet gain access? In which ways can attackers gain root access to host X, Y, or Z? Which attacks will become impossible if the firewall is set to block port 80? Which attacks will become impossible if all instances of CVE-2014-0497 are removed?

Clearly, a tool that is able to answer questions of this sort will support a decision maker who is considering different remediation options or performing a risk analysis. In addition,

1
2
3 this type of tool can also be used in conjunction with other
4 techniques, such as intrusion detection systems (Roschke et al.,
5 2010) or forensics tools (Liu et al., 2012), to improve their
6 analysis capabilities. Many research papers have been produced
7 on attack graphs, and several software tools have been
8 developed. However, the practical utility and validity of these
9 tools in a realistic setting is unclear, and there are reasons to
10 doubt that attack graphs yield accurate results. For example, the
11 utility of the attack graph approach is dependent on the
12 availability of information about the analyzed computer
13 network and the vulnerabilities associated with it. The
14 designers of tools aiming to support decision makers propose
15 that such information should be gathered with the help of
16 vulnerability and network scanners (see (Ou et al., 2005)
17 (Jajodia, 2007) (Artz, 2002)). However, it is known that
18 vulnerability scanners are limited in their accuracy and only
19 detect approximately half of the vulnerabilities in a network
20 (Holm et al., 2011). Because of such potential issues, this paper
21 attempts to answer the research question: *How well do attack
22 graphs predict the success or failure of attacks under realistic
23 conditions?*

24
25
26 This paper presents a test of MulVAL (Ou et al., 2005), one of
27 the more commonly cited and more mature attack graph tools.
28 In this test, MulVAL is fed network configuration data and
29 vulnerability data collected using the Nexpose vulnerability
30 scanner (from Rapid7) and is used to analyze eight computer
31 networks of different size and complexity. The output of this
32 analysis is compared to observations of successful attacks
33 performed during an offensive cyber security exercise. The
34 accuracy of MulVAL's predictions is reported in terms of how
35 well the predictions correspond to the observations made by
36 attackers on successful and unsuccessful attacks. To our
37 knowledge, this is the first empirical test of an attack graph tool
38 against observed attack attempts. Based on these empirical
39 results, suggestions are made as to how predictions can be
40 improved.

41
42
43 The outline of this paper is as follows. Section 2 gives an
44 overview of attack graphs, with a particular emphasis on the
45 variant used in MulVAL. Section 3 describes the methods used
46 in the test. Section 4 presents the results, which are
47 subsequently discussed in section 5.
48
49
50
51
52
53
54
55
56
57
58
59
60

2 Attack Graphs

Multiple articles provide reviews and overviews of different attack graph approaches. Lippmann and Ingols (2005) reviewed existing approaches in 2005 and classified them based on how they viewed the goals of the attackers, how they were generated, and how well they scaled. Problems were identified with respect to scalability, obtaining attack (or vulnerability) details, and identifying what attackers could connect to from different locations in the network. Multiple methods of representing attack graphs have evolved (Alhomidi and Reed, 2012). Heberlein et al. (2012) offers a canonical representation of attack graphs, illustrated in Figure 1.

In this representation, the nodes of the graph represent the state of the entire computer network at a given time. Here, a state characterizes both the state of the computer network itself (including hosts, software, configurations, authorized users and vulnerabilities) and the identities that the attacker owns (i.e., user accounts under the attacker's control). With this state space as a starting point, an attack graph shows how the attacker can move from one state to another using the vulnerabilities and controlled identities present in the network. An attacker can thus either change the state of the network itself (e.g., by reconfiguring a firewall) or obtain a new identify (e.g., by stealing a user's password or session).

Heberlein et al. (2004) also describe some features of the analyses used in attack graph approaches. These features can be summarized as follows:

- **Monotonic or non-monotonic exploits.** To simplify the analysis and offer scalability, it is commonly assumed that exploits are monotonic, i.e., that once the attacker has obtained some identity, that identity is never lost. MulVAL assumes that exploits are monotonic.
- **Single path or all paths.** For simplicity, the analysis is sometimes limited to generating a single attack path, i.e., the analysis stops as soon as one possible path to the specified state is found. MulVAL finds all paths.
- **Forward or backward chaining.** In a forward chaining technique, the analysis starts with an assumption of the current state and assesses which states are reachable. In a backward chaining technique, an end-state of interest

1
2
3 is defined, and the analysis assesses which states may
4 lead to that end-state. MulVAL uses backward chaining.
5

6 The mature attack graph tools from academia, namely MulVAL
7 (Ou et al., 2005) (Homer and Ou, 2009), the TVA tool (Jajodia,
8 2007) (Noel et al., 2009) (Jajodia and Noel, 2010), and NetSPA
9 (Artz, 2002) (Lippmann, 2002) (Chu et al., 2010) (Ingols et al.,
10 2009), share the same features, except that the TVA tool uses
11 forward chaining. There are, however, other differences
12 between these tools. For example, MulVAL uses Datalog rules
13 to specify its input to Prolog, whereas the TVA tool uses more
14 loosely defined input formats and operates on matrices
15 representing the attack steps. On a conceptual level, it can be
16 argued that the tools share the same problems and weaknesses,
17 and the accuracy of one of these tools ought to reflect the
18 accuracy of the other tools. Thus, while MulVAL is used in this
19 test, the results ought to be generalizable to similar attack graph
20 approaches.
21
22
23
24
25

26 Theoretically, there is little reason to question the validity of
27 the inferences produced by any of the attack graph tools. If the
28 tools' algorithms are provided correct input data, they will most
29 likely produce correct output data and yield accurate results.
30 Unfortunately, fully correct input data are rarely available to
31 the decision makers in enterprises, leading to the question of
32 how accurate results will be under realistic conditions. No
33 reports detailing tests of the accuracy under the conditions
34 proposed in research papers, in which the analyses describe
35 attack paths based on input from vulnerability scanners, can be
36 found for any of these tools. The paper that comes the closest to
37 providing a test of accuracy is the test performed by Zhang et
38 al. (2011). In this test, seven servers were assessed using
39 MulVAL on several occasions. The different results produced
40 by the tool on these different occasions were correlated with,
41 and explained by, changes in the state of the computers (e.g.,
42 new vulnerabilities). Consequently, the test assessed the
43 internal validity of the tool, but did not assess its validity in
44 practice or how well it predicted the success of real attacks.
45
46
47
48
49
50

51 **3 Methods and Materials**

52 This section describes the methods and materials used in the
53 test. Section 3.1 describes the criteria used to assess the
54 accuracy and utility of the predictions and the analysis method.
55 Section 3.2 describes the computer networks on which the
56
57
58
59
60

1
2
3 predictions were made. Section 3.3 describes the attackers in
4 the test, the scenario by which they were guided and the actual
5 attacks they performed. Section 3.4 describes the tool
6 configuration and how the output of the tool was codified.
7

8 9 *3.1 Assessment criteria*

10 This test aimed to evaluate the accuracy of the results produced
11 by an attack graph tool. There are different ways to view the
12 output of such tools, with implications for the assessment
13 criteria they ought to be evaluated against. These issues are
14 discussed below.
15

16
17 First, the inclusion or exclusion of the paths in a graph may be
18 performed differently. It might be expected that the attack
19 graph will 1) show all paths that can definitely be used, 2) only
20 exclude paths that definitely cannot be used, or 3) aim to assign
21 all possible attack attempts to the class that fits best. In this test,
22 we assume that the attack graphs adhere to the third option and
23 identify both attacks that are possible and attacks that are
24 impossible. This criterion is well in line with the common view
25 of attack graphs. For example, Alhomidi and Reed (2012) state
26 that attack graphs “show all ways of how an attacker violates a
27 security policy.”
28
29

30
31
32 Second, the output can be interpreted in a possibilistic or
33 probabilistic fashion. When the output is interpreted
34 possibilistically, all attacks with a likelihood of success above
35 zero are included in the attack graph, even if the likelihood that
36 they can be exploited in practice is infinitely small. In a
37 probabilistic interpretation, it is expected that attempts to
38 perform the attack steps in the graph are likely to succeed, e.g.,
39 because they are more likely than some threshold value. In any
40 of these interpretations, it should be expected that the inclusion
41 of an attack path and the possibility of compromising a
42 machine correspond to a higher success rate than the excluded
43 paths and those machines for which there is no possibility of
44 compromise and that all successful attacks are included in the
45 attack graph. Put differently, a possibilistic method should at
46 least be able to predict when an attacker will fail any attempted
47 attack and include all successful attacks. There should therefore
48 be an agreement between the attack graph’s predictions and the
49 observed success rates.
50
51

52
53
54
55
56 Third, one may require correctness at different levels of
57 abstraction. These levels of abstraction may vary from the
58
59
60

1
2
3 lowest level, of individual exploits and ports, to higher levels,
4 such as when the output (e.g., the number of paths) is
5 considered to provide only indications or examples of how
6 vulnerable a network is. This assessment focuses on the ability
7 to predict which machines attackers can execute code on and
8 the abstraction level of hosts and their attack paths. This is the
9 attack type and the abstraction level used in the vast majority of
10 all proposals related to attack graphs. The tool is assessed by
11 comparing the attack paths produced by the tool with the
12 attackers' observations and producing a confusion matrix for
13 the predictions (i.e., the possible and impossible attack paths)
14 and the (successful and failed) attacks against designated target
15 machines.
16
17
18
19

20 3.2 Attacked computer networks

21 In this test, over a thousand virtual machines were deployed,
22 together forming computer networks of various sizes and
23 complexity to represent a synthetic threat environment. Of
24 these machines, 199 machines in eight fictitious organizations
25 of different types were monitored and assigned as targets. The
26 199 machines used 88 different virtual machine templates and
27 were all configured differently in some regard, e.g., with
28 respect to users. Some organizations' networks consisted of
29 only a few computers, representing a small organization or
30 personal network; other organizations' networks consisted of
31 several VLANs with firewalls limiting the access possibilities
32 between them. Figure 2 illustrates one of the monitored
33 computer networks.
34
35
36
37
38

39 Multiple operating systems and applications were instantiated
40 in these networks. Different patch levels and versions of
41 Windows (including Windows 2000, XP, and 7 and Windows
42 Server 2003 and 2008) and several versions of Linux-based
43 distributions (including Gentoo, Debian, and Ubuntu) were
44 utilized. These ran multiple desktop and server applications.
45 Among others, the client-side applications included different
46 versions of Adobe Reader, software development tools such as
47 Visual Studio, Internet Explorer, and Firefox, and applications
48 of the Microsoft Office suite or Open Office. Server-side
49 applications included different versions of Wordpress,
50 phpMyAdmin, IIS, domain controllers, network infrastructure
51 services (e.g., DNS and DHCP), and FTP servers. The aim was
52 that the deployed applications should be representative of the
53 standard software found in most enterprises' computer
54
55
56
57
58
59
60

1
2
3 networks. However, to enable a meaningful exercise for the
4 attacking teams and to produce enough data for the test, these
5 applications were more vulnerable than the ones found in the
6 typical enterprise (i.e., they had not been updated and patched
7 recently). Another difference between this environment and the
8 computer networks typical of organizations is the lack of
9 custom-built applications (e.g., interconnected spreadsheet
10 applications) and larger enterprise systems (e.g., ERP systems).
11
12

13
14 To allow interaction with the users on the machines, scripts
15 implemented in Auto IT (Jonathan Bennett AutoIt Consulting
16 Ltd, 2013) emulated user actions. These scripts used the
17 applications to send emails to each other, surf websites in the
18 environment, open emails and attachments and access files on
19 local machines. To produce a realistic behavioral pattern, the
20 emulated behaviors were performed according to a predefined
21 instruction list created based on the actions of real users in an
22 office environment. The instruction list was produced by
23 collecting the historic usage of web browsers and desktop
24 applications and the outgoing emails of 17 users in three
25 organizations (two research organizations and one game
26 developer). To generate a variety of user behavior, the historic
27 records (covering months to years of computer usage) were
28 split into one-week instruction sets, which allowed thousands
29 of instruction lists (i.e., “users”) to be created. The activities
30 performed by the scripted user agents thus followed the same
31 sequence and had the same intensity as real users during a
32 typical workweek. However, the agents were dumb in the sense
33 that they did not engage in dialogs. Their responses to the
34 emails sent to them were to click all links and open all
35 attachments, i.e., they were always cooperative as seen from
36 the attackers’ point of view. Furthermore, they only performed
37 standard user behaviors; more sporadic tasks, such as the
38 installation of custom software applications or administrative
39 tasks, were not performed in this test.
40
41
42
43
44
45
46
47

48 3.3 Attackers, scenario and attacks

49 Two independent teams attacked the computer networks to find
50 secret keys hidden within them. One team consisted of security
51 researchers from the Swedish Defence Research Agency, and
52 the other team consisted of security specialists from the
53 Swedish Armed Forces Network and Telecommunications
54 Unit. Both teams restricted themselves to using only publicly
55
56
57
58
59
60

1
2
3 available tools and publicly known exploits, e.g., the tools and
4 exploits packed with the Backtrack 5 operating system.
5

6
7 The attacks started at noon on a Tuesday in November of 2012
8 and continued until noon on Thursday of the same week, with
9 no scheduled interruptions or pauses. The attackers had no prior
10 knowledge of which machines had secret keys hidden in them
11 but were told that they were in some of the eight computer
12 networks. Their mission was to compromise machines and
13 extracts the keys. As a result, a mix of reconnaissance and
14 penetration activities was conducted. The two teams logged
15 134 penetration attempts and 31 network scans aimed at the
16 networks monitored in this test. Additionally, over 100 other
17 types of reconnaissance activities were undertaken, e.g., ARP
18 scans. These penetration attempts led to 40 successfully
19 compromised machines. In all of these, privileges were
20 obtained to execute code as root, system administrator, or
21 similar.
22
23
24
25

26 The reconnaissance activities and attacks, together with their
27 successes, were recorded by each attacker in a log. The
28 accuracy of these logs was assessed by comparing them to
29 recordings made by screen capture software installed on the
30 machines used by some of the attackers. These comparisons
31 showed very high agreement between what was done and what
32 was logged. All of the deviations identified concerned
33 reconnaissance activities, e.g., network scans.
34
35
36

37 *3.4 Tool configuration and input data*

38 The MulVAL project is open source and available for
39 download (see Ou et al. (2013)). MulVAL requires information
40 on vulnerabilities, subnets, users and host access control lists as
41 input. For each vulnerability, the following fields are specified:
42 the hostname, Common Vulnerabilities Enumeration ID,
43 program (target service or application), range list (if it is
44 remotely or locally exploitable), type of loss (textual), severity
45 (high/medium/low), access control requirement
46 (high/medium/low), service and port. Subnets are described by
47 the hosts they include, and the host access control lists describe
48 by which protocol and on which port entities (i.e., hosts or
49 subnets) they are allowed to communicate. MulVAL also
50 makes it possible to include users and their accounts in the
51 analysis. These users can be labeled as incompetent, meaning
52 they are susceptible to social engineering attacks. In this test,
53 machines with active user agents were labeled as incompetent
54
55
56
57
58
59
60

1
2
3 to reflect that the users opened all email attachments and
4 clicked links indiscriminately.
5

6 Descriptions of the systems used in this test were provided to
7 MulVAL by extracting data from a) the system configuration
8 files used to instantiate the computer networks in the cyber
9 range and b) vulnerability scans of the targeted hosts, produced
10 by Nexpose. The system configuration files were validated to
11 provide a reliable source of information and were used for all
12 information except the vulnerabilities and host access control
13 lists. Vulnerability scans were used to collect information about
14 the vulnerabilities in the hosts. These scans were performed
15 repeatedly on the individual vulnerabilities on each of the
16 deployed hosts (88 different variants on 199 machines) in a
17 controlled environment to ensure a reliable output. The host
18 access control lists were created based on the IP table files in
19 the deployed firewalls and the settings in the host firewalls.
20
21
22
23
24

25 To perform these steps and support the analysis, a web-based
26 tool was created that imported the system configuration files
27 and vulnerability scans and performed the analysis of interest.
28 MulVAL was set to analyze all the ways an attacker located on
29 the internet (i.e., outside the computer networks) could execute
30 code on the hosts in the targeted computer networks using
31 different privilege levels. MulVAL operates on the *root*, *user*
32 and less precise *someUser* privileges. Backward chaining was
33 then used to determine all of the ways the machines could be
34 reached. Individual attack paths and attack steps in the resulting
35 attack graph were identified through a breadth-first search on a
36 digraph produced based on the XML file generated by
37 MulVAL. The networks used in this test generated a
38 considerable number of possible paths. With a search depth of
39 60 steps, 495,360 existing attack paths were found for the eight
40 networks and 199 machines. These steps were compared to the
41 steps attackers used when evaluating the prediction accuracy.
42
43
44
45
46
47

48 **4 Results**

49 The prediction accuracy of MulVAL is presented in section 4.1.
50 As will be shown, the accuracy of MulVAL's prediction poorly
51 fit the attackers' successes and failures. Section 4.2 goes a step
52 further and investigates the reasons for MulVAL's incorrect
53 predictions.
54
55
56
57
58
59
60

4.1 Prediction accuracy

All attacks performed by the attackers in this test led to the ability to execute code at the highest privilege level (viz. root, system administrator or similar) on the targeted machines. Table 1 summarizes the results and the relationships between attackers' successes in compromising machines and the tool's predictions of the ability to execute code on machines at the highest privilege level (called *root* in MulVAL). In other words, Table 1 describes the relationship between the attackers' ability to obtain root privileges on machines and the tools predictions for obtaining root privileges on the machines.

As Table 1 shows, the attackers successfully compromised 40 hosts in the computer networks as root. Only 6 (15%) of these hosts were compromised using an attack path included in MulVAL's output. The attackers failed to compromise 159 of the hosts they probed, with 85 (53%) of these predicted to be impossible by MulVAL. The attacker could only execute code as root on 6 (8%) of the 80 machines that MulVAL predicted could be compromised to give root-level code execution privileges. Thus, 74 (93%) of the machines that MulVAL predicted were possible to obtain root access on were not compromised. Of the 119 machines MulVAL predicted as impossible to compromise, 34 (29%) were successfully compromised as root. In three of these 34 cases, MulVAL reported the machine as reachable but failed to describe a path that included the exploited vulnerability.

4.2 Reasons for incorrect predictions

The successful attacks missed by MulVAL can be explained by a combination of inaccurate vulnerability scans and inaccurate interpretations of vulnerability information. This section provides further details on this.

MulVAL requires information about the vulnerabilities in the computer network. This information is typically, as in this case, collected using a network vulnerability scanner. As could be expected, the imperfect information provided by the vulnerability scanner influences the results negatively. All 34 hosts reported as false negatives in Table 1 were compromised using a vulnerability Nexpose did not report. However, this is only a part of the explanation. If these exploited vulnerabilities that Nexpose missed are manually added to the scan results, only seven of the 34 become true positives. The remaining 27 machines are still assessed as impossible to execute code on

1
2
3 with root privileges, but all are possible to execute code on with
4 lower privileges. The reason for this is the way MulVAL
5 interprets and processes vulnerability information.
6

7
8 All the vulnerabilities missed by the scanner (CVE-1999-0504,
9 CVE2003-0352, CVE-2006-3439, CVE-2007-3039, CVE-
10 2007-1748, CVE-2008-4250, and CVE-2010-0478) are known
11 to be able to yield root/administrator/system privileges without
12 subsequent attacks that yield escalate privileges. An indication
13 of this is that all but one (CVE-1999-0504) are marked as
14 having full impact on confidentiality, integrity and availability
15 in the US National Vulnerability Database. However, MulVAL
16 labels these, and all other remote exploits, as yielding the
17 access level of *someUser*. Local exploits are labeled as giving
18 root privileges. In fact, further inspection of the data shows that
19 none of the thirteen true positives obtained with the corrected
20 vulnerability scans were predicted entirely correct. All these
21 attacks involved a privilege escalation attack that increased the
22 *someUser* privileges to *root* privileges or reused *root*-level
23 accounts reached elsewhere in the attack graph.
24
25

26
27
28 If vulnerabilities are added manually and MulVAL's
29 interpretation is overridden so that they yield *root*-level
30 privileges, all 40 successful attacks are predicted. All 40 are
31 also predicted as possible to compromise if it is judged as
32 sufficient if MulVAL identifies that code can be executed with
33 some privileges (i.e., *root*-level privileges are not required);
34 however, under such conditions, the vast majority of machines
35 not compromised are also predicted as possible to compromise.
36 Thus, in summary, a combination of missing vulnerability
37 information and poor processing of vulnerability information
38 explains the false negatives.
39
40
41
42

43
44 The 74 false positives are more problematic to find the reason
45 for in this test. When attackers fail to compromise a machine,
46 this provides a clear indication that the machine is more
47 difficult to compromise than the other machines that were
48 compromised. However, their failure does not mean that attacks
49 are impossible to accomplish. Because of this, it is unclear how
50 many of the 74 false positives for root-level access would have
51 been true positives in another test. The 74 false positives
52 include predictions of attacks that were attempted and failed,
53 but it also includes cases where the attackers scanned the host
54 but failed to find an attack worth trying. Furthermore, typical
55 attacks performed against the 74 machines are password-
56
57
58
59
60

1
2
3 guessing attacks against SSH logins, malware attached to email
4 and exploitation of web application vulnerabilities. Thus, while
5 attackers tried to compromise the machines, they sometimes
6 did not find a path at all, and they rarely tried a stable exploit
7 included in MulVAL's attack paths.
8
9

10 11 **5 Discussion and Future Work**

12 Attack graph tools such as MulVAL are easy to use. Provided
13 that host access control lists can be produced (e.g., from
14 firewall rules) and vulnerability information can be collected
15 (using a vulnerability scanner), the tools can perform their
16 analysis. The output of the analysis contains both a list of the
17 privileges each attack leads to (e.g., on which machines the
18 attack can execute code) and a full graph of the attacks
19 providing those privileges (e.g., the software vulnerabilities
20 exploited). A decision maker may try to work with this
21 information directly or add an analysis technique on top of the
22 attack graph tool's output. For example, critical paths can be
23 identified to produce a prioritized list of mitigation options
24 using ranking algorithms, as described in Sawilla and Ou
25 (2008). Given the sheer number of paths that are produced
26 (almost 500,000 in this test), such post-analysis methods are
27 likely to be needed to make the analysis results
28 comprehensible.
29
30
31
32
33
34

35 This test did not investigate techniques that could build on
36 attack graphs or how the techniques could influence the
37 analysis results. Instead, the test directly investigated the
38 accuracy of MulVAL when used in combination with a
39 vulnerability scanner. The accuracy of such an analysis would
40 serve as the foundation for analysis methods based on attack
41 graphs. Unfortunately, the results show that MulVAL failed to
42 predict which attacks red teams could accomplish during a two-
43 day exercise:
44
45
46

- 47 • The red teams were three times as likely to accomplish
48 what MulVAL predicted as impossible (29% success
49 rate) as to accomplish attacks along paths MulVAL
50 predicted as possible (8%).
- 51 • Only 6 (8%) of the 80 machines MulVAL predicted as
52 possible to obtain root level privileges on were
53 comprised during the 48 hours.
54
55
56
57
58
59
60

1
2
3 This test was performed to answer the research question: *how*
4 *well do attack graphs predict the success or failure of attacks*
5 *under realistic conditions?* Given the results of this test, the
6 answer to that question is *very poorly*. The primary reasons for
7 the false negatives are reliance on vulnerability scanners and
8 inaccurate interpretation of vulnerabilities provided by them.
9 Further discussions of the results and conditions that may have
10 skewed the results are discussed in section 5.1 below. In section
11 5.2, recommendations to researchers are given.
12
13
14

15 5.1 Possible explanations for poor accuracy

16 There are several possible explanations for the poor
17 performance of the attack graph tool in the test in addition to
18 the issues with vulnerability information and privilege levels
19 described above. Some of these can be seen as excuses for why
20 this particular test produced an unsatisfying result; others are
21 related to the test criteria and the general properties of attack
22 graphs. In the paragraphs below, some presumed objections to
23 the results are stated (in italics) together with a response to the
24 objection.
25
26
27
28

29 ***The attackers in this test were unrepresentative of the threat***
30 ***scenario in which attack graphs are supposed to be used.*** It is
31 true that the attackers are of a certain type, and they are not
32 representative of the attackers threatening enterprises in
33 general. For example, there were no malware writers, botnet
34 herds, or security-illiterate disgruntled employees involved in
35 this test. However, the type of threat for which MulVAL
36 produces predictions has not been defined, requiring
37 interpretations of its scope. Based on the reasoning provided in
38 articles on MulVAL and similar tools, it is reasonable to
39 assume that attack graphs (and MulVAL) should work when
40 there is a match between the attacks and vulnerabilities
41 modeled and the attacks and vulnerabilities the threat is capable
42 of finding, i.e., when the vulnerabilities fed into MulVAL are
43 those that the attackers might exploit. In this test, in which the
44 attackers were limited to publicly available tools and MulVAL
45 was fed the output of a vulnerability scan, the definitions of the
46 vulnerabilities were clearly matched. Furthermore, even if there
47 was a mismatch and another type of attacker was imagined for
48 the attack graphs, it is reasonable to expect that the predictions
49 would also be indicative for this type of attacker. In other
50 words, even with the wrong type of attackers, it should still be
51
52
53
54
55
56
57
58
59
60

1
2
3 expected that an accurate prediction would correlate with the
4 observations made for the other attackers.
5

6 ***The attack efforts were not independent and randomly***
7 ***distributed.*** This claim is true. In this test, the same attack may
8 have been attempted multiple times, and the attacks were
9 selected by human agents who (presumably) reasoned about the
10 best course of action before their attempts. They may have
11 spent hours on some machines and dismissed others as
12 impossible to compromise in seconds. It is reasonable to expect
13 that these humans thought in a similar manner to MulVAL and
14 considered the exploitation of existing vulnerabilities within
15 their reach. If this is the case, they should have been more
16 likely to test the attacks that MulVAL predicted as possible and
17 less likely to test the attacks that MulVAL predicted as
18 impossible. Therefore, it may be the case that a randomized
19 sample of attacks would produce a higher proportion of
20 practically unlikely attacks that would be easy to predict as
21 impossible. However, this is a poor explanation for the low
22 prediction accuracy for successful attacks or for attacks
23 predicted as possible: the attackers were more likely to succeed
24 if MulVAL said that the attack was impossible than if MulVAL
25 said it was possible. Additionally, it can be argued that a test
26 with randomly selected attack paths would lack the ecological
27 validity needed to say whether MulVAL works in an
28 operational context, in which attacks are not random.
29
30
31
32
33
34
35

36 ***The attackers' logs may have been erroneous and introduced***
37 ***biased measurement errors.*** Half of the attackers had their
38 screens recorded. No issues were found with the logs' accuracy,
39 except for some cases in which the attacker did not report a
40 failed attempt, typically when the attacker tried multiple ways
41 to compromise a machine. The omission of failed attempts does
42 not influence the results in this test, in which all of the probed
43 machines were considered interesting targets for the attackers.
44 The logs produced by the attackers probably contained some
45 flaws and missed other additional information. However, based
46 on the comparisons with the screen recordings, it is safe to say
47 that no errors were of a sufficient magnitude to threaten the
48 overall conclusions.
49
50
51
52

53 ***Unrealistically vulnerable networks were used, and this***
54 ***influenced the results.*** The computer networks used in this test
55 were certainly more vulnerable than the average enterprise's
56 computer network, with some of the machines not having been
57
58
59
60

1
2
3 updated in a decade. The high number of exploitable
4 vulnerabilities led to an unusually large number of possible
5 attack paths through the computer network. Although this is
6 unrealistic, it is difficult to see why this would cause the poor
7 accuracy of the attack graph tool. On the contrary, the use of
8 well-known vulnerabilities implies that the tool had accurate
9 information and therefore should have produced accurate
10 results.
11
12

13
14 ***This was really a test of the vulnerability scanner providing***
15 ***vulnerability information to MulVAL.*** To some extent, this is a
16 valid objection to the results. The analysis was not made under
17 the premise that perfect information was available. Incomplete
18 scan results contributed to a large portion of false negatives.
19 For example, the scans did not report the CVE-2008-4250
20 vulnerability in Windows machines, a vulnerability used by
21 attackers in 38 attack paths to compromise 24 machines in this
22 test. However, the result would have been poor even if these
23 vulnerabilities had been reported by the scanner. When the scan
24 results were complemented with all vulnerabilities used by the
25 attackers, MulVAL still missed 27 of 40 successful root-level
26 code execution attacks. In fact, closer inspection shows that,
27 because of problems associated with interpreting the effect of
28 vulnerabilities, none of the predicted attacks correspond
29 perfectly to the attacks performed by the attacker. To predict all
30 successful attacks it is not sufficient to complement the
31 vulnerability scan with missed vulnerabilities; it is also required
32 that an analyst interpret the results and manually set the
33 privileges that can be obtained if the vulnerability is exploited.
34 Such manual adjustments may be possible to do before an
35 analysis, and MulVAL can be improved to guess privilege
36 levels better. However, even if the problem with the analysis
37 had been the input from the vulnerability scanner, it has been
38 repeatedly argued that attack graphs are practically useful
39 because they can be fed input from such vulnerability scans and
40 then make predictions in an automated fashion. Thus, although
41 the vulnerability scanner is important to the results, it makes
42 sense to see it as part of the solution. As a side note, the makers
43 of MulVAL advocate the use of Nessus, another vulnerability
44 scanner, by providing scripts for parsing its output files. In this
45 test, these files were adapted to use Nexpose's output instead.
46 According to the test performed in 2011 by Holm et al. (2011),
47 no significant difference should be expected in accuracy if
48 Nessus were used instead.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 *Attack graphs are possibilistic, and treating the output as a*
4 *probabilistic indicator is unfair.* As noted above, it is unclear
5 how the results of attack graphs such as MulVAL are supposed
6 to be interpreted. In this test, it was expected that MulVAL's
7 classification could be used as an indicator of the attackers'
8 capabilities. If a truly possibilistic interpretation were to be
9 used instead, an attack would be labeled as possible even if it
10 were highly unlikely that it could be accomplished and
11 impossible only if it were affirmed that it was impossible under
12 the conditions given. In the confusion matrix in Table 1, this
13 would mean that failed attacks could not be used to evaluate the
14 tool because they could be the result of, for example, bad luck
15 or an incompetent attacker. It would also imply that the
16 probability of success of an attack could be any value,
17 including extremely low values such as one percent. First, not
18 even an extreme possibilistic interpretation can explain why 34
19 of the 119 machines that were considered impossible to
20 compromise as root were actually compromised. To come to a
21 possibilistic result, the input from the vulnerability scanner is
22 corrected and MulVAL's interpretation of it is manually
23 overridden or improved. Second, the designers of MulVAL do
24 not seem to have implemented a possibilistic (worst case)
25 solution. Because it is, of course, possible that a remote exploit
26 yields root-level access, a possibilistic solution ought to have
27 indicated that this is possible rather than indicating that the
28 privileges of some user are obtained.
29
30
31
32
33
34
35
36

37 5.2 Recommendations to researchers

38 Given these results, four recommendations are provided to
39 researchers interested in attack graphs.
40

41 First, we recommend that researchers interpret the results of
42 this study with caution and a positive spirit. Although the
43 results suggest that attack graphs, when used together with a
44 normal vulnerability scanner, fail to predict what a security
45 professional can accomplish, there are several nuisance
46 variables that may have distorted the results. These nuisance
47 variables include the attackers themselves, the attack graph tool
48 used, the vulnerabilities in the computer networks and the
49 vulnerability scanner used. Although it may be hard to see how
50 any realistic configuration of these variables would result in
51 accurate predictions, further tests should be performed.
52 Furthermore, the labeling of privilege levels in MulVAL could
53 be improved to produce better predictions of privilege levels,
54
55
56
57
58
59
60

1
2
3 e.g., by guessing based on the impact vector in the Common
4 Vulnerability Scoring System (CVSS) (Mell et al., 2007), using
5 complementary vulnerability information from other sources or
6 making a possibilistic (worst-case scenario) guess.
7
8

9
10 Second, there may be techniques and algorithms that could
11 improve the accuracy of attack graphs that should be tested
12 empirically. MulVAL produced approximately 500,000 un-
13 prioritized attack paths for the eight networks included in this
14 test. A visual graph would be impossible for the human eye to
15 comprehend, and it is unclear how it would support decision-
16 making. If attack graphs were prioritized and ranked relative to
17 each other, the output of the analysis might become more
18 accurate and more useful. Several suggestions have been made
19 in this direction, including assessments of critical assets in the
20 graph (Sawilla and Ou, 2008) and probabilistic rankings of the
21 paths (Homer et al., 2010) (Singhal and Ou, 2009).
22 Alternatively, conditions not related to individual
23 vulnerabilities could be used to improve the accuracy, for
24 example, by using quantitative estimates of remote code
25 execution attacks, such as those provided in studies such as
26 (Sommetad et al., 2012). These alternatives ought to be
27 considered after techniques for interpreting the privilege levels
28 attained by exploitations of vulnerabilities are improved.
29
30
31
32
33

34 Third, the accuracy of vulnerability scanners is a serious
35 practical obstacle for attack graphs today. Provided that
36 vulnerability information is interpreted correctly, significant
37 improvements in their accuracy would lead to significant
38 improvements in the accuracy of attack graphs. Research could
39 be directed toward improving vulnerability information and
40 vulnerability scanners to accelerate their progress. Additionally,
41 improvements in the accuracy of attack graphs could be gained
42 through improvements in threat intelligence, e.g., by improving
43 the understanding of the modus operandi of presumed threat
44 agents or which vulnerabilities they are capable of exploiting.
45 With such information, it would be possible to provide the
46 attack graph tools with better data and increase their accuracy,
47 e.g., through combination with probabilistic approaches or
48 other means of ranking attack steps and paths to support
49 decision-making.
50
51
52
53
54

55 A fourth recommendation is to lower the expectations for
56 attack graphs as a practically useful vulnerability prediction or
57 analysis tool. Although no direct claims have been made about
58
59
60

1
2
3 the accuracy of attack graph tools in general, or MulVAL in
4 particular, it is easy to get the impression that a vulnerability
5 scan and attack graph analysis will serve as an effective
6 vulnerability prediction and analysis method for decision
7 makers. For example, Ou et al. (2006) start their abstract with,
8 “[a]ttack graphs are important tools for analyzing security
9 vulnerabilities in enterprise networks.” Roschke et al. (2009)
10 start their abstract with, “[a]ttack graph is used as an effective
11 method to model, analyze, and evaluate the security of
12 complicated computer systems or networks.” Williams (2008)
13 starts with the declaration, “[a]ttack graphs are valuable tools
14 in the assessment of network security, revealing potential
15 attack paths an adversary could use to gain control of network
16 assets.” Jajodia (2007) concludes that the TVA tool is “a
17 powerful approach to global network vulnerability analysis.”
18 Such statements are clearly questionable in light of the results
19 of this test, where attacks are more likely to succeed if they are
20 not predicted by the tool. But, as noted above, attack graphs
21 may become practically useful in the future. Furthermore,
22 despite the difficulty in providing accurate input data, attack
23 graphs could function as a framework on which security
24 theories could be attached, related to each other and
25 synthesized.

32 33 34 **6 Conclusions**

35 This test determines that the attack graph tool MulVAL
36 predicts human attackers’ successes poorly when used together
37 with the vulnerability scanner Nexpose. Only 8 percent of the
38 machines predicted as possible to compromise were
39 compromised; 29 percent of the machines predicted as
40 impossible to compromise were compromised. This inaccuracy
41 is due to the combination of inaccurate vulnerability scans and
42 improper interpretation of the privileges that vulnerabilities
43 grant. If vulnerabilities are manually added and manually
44 corrected to provide the right privilege level, all but one
45 compromised machine are predicted as possible to compromise.

46 47 48 49 50 51 **7 References**

52 Alhomidi, M. a. and Reed, M.J. (2012), “Attack graphs
53 representations”, *2012 4th Computer Science and*
54 *Electronic Engineering Conference (CEEC)*, Ieee, pp. 83–
55 88.
56
57
58
59
60

- 1
2
3 Artz, M.L. (2002), *Netspa: A network security planning*
4 *architecture*, Massachusetts Institute of Technology,
5 available at:
6 [http://scholar.google.com/scholar?hl=en&btnG=Search&q](http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:NetSPA+:+A+Network+Security+Planning+Architecture+by#0)
7 [=intitle:NetSPA+:+A+Network+Security+Planning+Archi](http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:NetSPA+:+A+Network+Security+Planning+Architecture+by#0)
8 [tecture+by#0](http://scholar.google.com/scholar?hl=en&btnG=Search&q=intitle:NetSPA+:+A+Network+Security+Planning+Architecture+by#0).
9
- 10
11 Chu, M., Ingols, K., Lippmann, R., Webster, S. and Boyer, S.
12 (2010), “Visualizing attack graphs, reachability, and trust
13 relationships with NAVIGATOR”, *Proceedings of the*
14 *Seventh International Symposium on Visualization for*
15 *Cyber Security*, ACM, pp. 22–33.
16
- 17
18 Heberlein, T., Bishop, M., Ceesay, E., Danforth, M.,
19 Senthilkumar, C. and Stallard, T. (2004), “A Taxonomy
20 for Comparing Attack-Graph Approaches”, *netsq.com*.
21
- 22
23 Holm, H., Sommestad, T., Almroth, J. and Persson, M. (2011),
24 “A quantitative evaluation of vulnerability scanning”,
25 *Information Management & Computer Security*, Vol. 19
26 No. 4, pp. 231–247.
27
- 28
29 Homer, J., Manhattan, K., Ou, X. (Simon) and Schmidt, D.
30 (2010), *A Sound and Practical Approach to Quantifying*
31 *Security Risk in Enterprise Networks*, *people.cis.ksu.edu*,
32 Kansas, available at:
33 [http://people.cis.ksu.edu/~xou/publications/tr_homer_080](http://people.cis.ksu.edu/~xou/publications/tr_homer_0809.pdf)
34 [9.pdf](http://people.cis.ksu.edu/~xou/publications/tr_homer_0809.pdf) (accessed 24 June 2010).
35
- 36
37 Homer, J. and Ou, X. (Simon). (2009), “SAT-solving
38 approaches to context-aware enterprise network security
39 management”, *IEEE Journal on Selected Areas in*
40 *Communications*, Citeseer, Vol. 27 No. 3, pp. 315–322.
41
- 42
43 Ingols, K., Chu, M., Lippmann, R., Webster, S. and Boyer, S.
44 (2009), “Modeling Modern Network Attacks and
45 Countermeasures Using Attack Graphs”, *Annual*
46 *Computer Security Applications Conference*, IEEE, pp.
47 117–126.
48
- 49
50 Jajodia, S. (2007), “Topological analysis of network attack
51 vulnerability”, *Proceedings of the 2nd ACM symposium on*
52 *Information, computer and communications security -*
53 *ASIACCS '07*, ACM Press, New York, New York, USA,
54 p. 2.
55
- 56
57 Jajodia, S. and Noel, S. (2010), *Advanced cyber attack*
58 *modeling analysis and visualization*, Rome, NY, available
59 at:
60

1
2
3 <http://oai.dtic.mil/oai/oai?verb=getRecord&metadataPrefix=html&identifier=ADA516716> (accessed 1 April 2014).

6 Jonathan Bennett AutoIt Consulting Ltd. (2013), “AutoIt Script Editor”, available at: <http://www.autoitscript.com/site/autoit/> (accessed 13 January 2014).

11 Lippmann, R. (2002), *Netspa: A network security planning architecture*, Network Security, Massachusetts Institute of Technology.

16 Lippmann, R. and Ingols, K. (2005), *An annotated review of past papers on attack graphs*, Lexington, Massachusetts, available at: <http://www.dtic.mil/cgi-bin/GetTRDoc?AD=ADA431826&Location=U2&doc=GetTRDoc.pdf> (accessed 14 September 2010).

23 Liu, C., Singhal, A. and Wijesekera, D. (2012), “Using Attack Graphs in Forensic Examinations”, *2012 Seventh International Conference on Availability, Reliability and Security*, Ieee, pp. 596–603.

28 Mell, P., Scarfone, K. and Romanosky, S. (2007), *A Complete Guide to the Common Vulnerability Scoring System (CVSS), Version 2.0*.

33 Noel, S., Elder, M., Jajodia, S., Kalapa, P., O’Hare, S. and Prole, K. (2009), *Advances in Topological Vulnerability Analysis, 2009 Cybersecurity Applications & Technology Conference for Homeland Security*, IEEE, Washington, DC, pp. 124–129.

39 Ou, X. (Simon), Boyer, W.F. and Zhang, S. (2013), “MulVAL: A logic-based enterprise network security analyzer”, available at: <http://www.arguslab.org/mulval.html> (accessed 4 March 2014).

45 Ou, X. (Simon), Boyer, W.W.F. and McQueen, M.A. (2006), “A scalable approach to attack graph generation”, *Proceedings of the 13th ACM conference on Computer and communications security*, ACM, Alexandria, Virginia, USA, pp. 336–345.

51
52 Ou, X. (Simon), Govindavajhala, S. and Appel, A.W. (2005), “MulVAL: A logic-based network security analyzer”, *Proceedings of the 14th conference on USENIX Security Symposium-Volume 14*, USENIX Association, p. 8.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

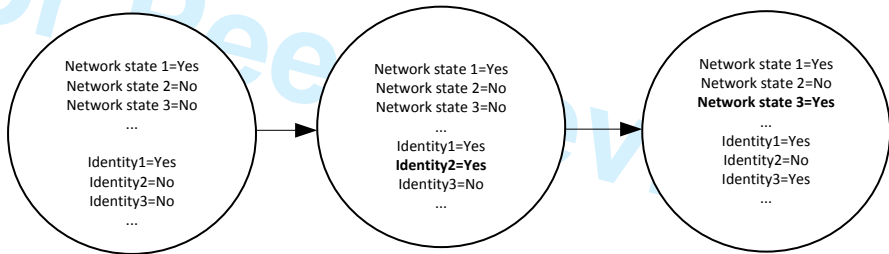
- 1
2
3 Roschke, S., Cheng, F. and Meinel, C. (2010), "Using
4 vulnerability information and attack graphs for intrusion
5 detection", *2010 Sixth International Conference on*
6 *Information Assurance and Security*, IEEE, pp. 68–73.
7
- 8
9 Roschke, S., Cheng, F., Schuppenies, R. and Meinel, C. (2009),
10 "Towards unifying vulnerability information for attack
11 graph construction", *Information Security*, pp. 218–233.
12
- 13 Sawilla, R. and Ou, X. (Simon). (2008), "Identifying critical
14 attack assets in dependency attack graphs", *13th European*
15 *Symposium on Research in Computer Security*
16 *(ESORICS)*, Springer, pp. 18–34.
17
- 18
19 Singhal, A. and Ou, X. (Simon). (2009), "Techniques for
20 enterprise network security metrics", *Proceedings of the*
21 *5th Annual Workshop on Cyber Security and Information*
22 *Intelligence Research: Cyber Security and Information*
23 *Intelligence Challenges and Strategies*, ACM, p. 25.
24
- 25
26 Sommestad, T., Holm, H. and Ekstedt, M. (2012), "Estimates
27 of success rates of remote arbitrary code execution
28 attacks", *Information Management & Computer Security*,
29 Vol. 20 No. 2, pp. 107 – 122.
30
- 31 Williams, L. (2008), *GARNET: A graphical attack graph and*
32 *reachability network evaluation tool*, (Goodall, J.R.,
33 Conti, G. and Ma, K.-L.,Eds.)*5th International Workshop,*
34 *VizSec 2008*, Springer Berlin Heidelberg, Cambridge, MA,
35 USA.
36
- 37
38 Zhang, S., Ou, X. (Simon), Singhal, A. and Homer, J. (2011),
39 *An empirical study of a vulnerability metric aggregation*
40 *method*, *csrc.nist.gov*, available at:
41 [http://csrc.nist.gov/staff/Singhal/xou-anoop-](http://csrc.nist.gov/staff/Singhal/xou-anoop-workshop2011-paper.pdf)
42 [workshop2011-paper.pdf](http://csrc.nist.gov/staff/Singhal/xou-anoop-workshop2011-paper.pdf) (accessed 27 July 2011).
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 *Figure 1. A canonical representation of attack graphs, with*
4 *transitions between states in the network and identities*
5 *controlled by the attacker, adapted from Heberlein et al.*
6 *(2012).*
7

8
9 *Figure 2. Example of an organization monitored and targeted*
10 *in the test. Circular nodes are network interfaces, and*
11 *rectangular nodes are machines. The globe is the link to the*
12 *“Internet” of the test environment.*
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

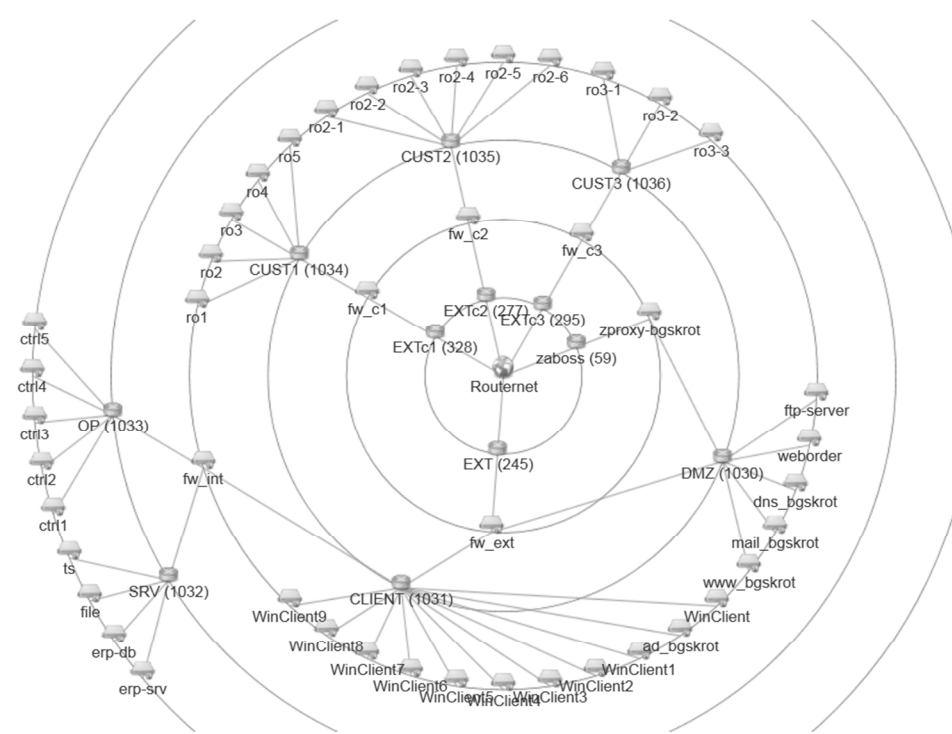
For Peer Review

1
2
3
4
5
6
7
8
9



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

For Peer



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1. Confusion matrix of MulVAL’s predictions for code execution privileges as root and attackers success at obtaining such privileges.

| | | Code execution by attackers as root | |
|---------------------------------------|------------|-------------------------------------|----|
| | | Yes | No |
| Prediction for code execution as root | Possible | 6 | 74 |
| | Impossible | 34 | 85 |

For Peer Review