

# BugBox : A Vulnerability Corpus for PHP Web Applications

Gary Nilson

*Computer Science Department  
University of Maryland*

Jeff Stuckman

*Computer Science Department  
University of Maryland*

Kent Wills

*Computer Science Department  
University of Maryland*

Jim Purtilo

*Computer Science Department  
University of Maryland*

## Abstract

Web Applications provide a robust amount of code vulnerabilities that are exploited on a routine basis, but a large, diverse, automated, ease of use corpus for analyzing these vulnerabilities is not publicly available. It is recognized that obtaining an adequate vulnerability dataset is a significant hurdle in empirical vulnerability research. Many studies suffer from using vulnerability data-sets that are too small, and they are rarely made publicly available [ref]. The purpose of BugBox is to fill this void, and in doing so encourage the quality and repeatability of results. We propose a scalable framework for web application corpus design and further explain the motivation and architecture behind BugBox, an open-source collection and management framework for PHP web application vulnerabilities. Through our framework implementation and corpus we can enhance testing vulnerability indicators and metrics, developing vulnerability definition representations, testing intrusion detection systems, or creating demonstrations for training purposes.

## 1 Introduction

**Overview** With the surge of internet blogging sites <sup>1</sup>, it has never before been so critical to release and maintain secure web applications. Many blog administrators rely on these web applications, which tend to be open source, to be secure frameworks. While blogging engines do represent a majority of web applications used on the internet, there are many other applications available to the community. Due to the community of users in this space, many web applications are subject to a rich variety of exploit types, such as: cross-site scripting (XSS), cross-site request forgery (CSRF), SQL injection, buffer overflow, and more. Of the web applications offered to the community, many popular applications, such as wordpress,..., are designed in PHP. Currently, PHP has been reported

by Top Cyber Security Risks to be a target for security threats [ref].

**Database** Popular vulnerability/exploit databases such as the National Vulnerability Database (NVD), Open Source Vulnerability Database (OSVDB), Common Vulnerability and Exposures Database (CVE), or the Exploit Database (EDB) log PHP application vulnerabilities/exploits as a service to the community.

**Use** While many of these databases exist and disclose the details of the vulnerability/exploit, many do not provide vulnerable code that can be used for analysis. While sufficient information is not included on these sites for code analysis, we use these databases as a template to write exploits from which we are able to gather vulnerable code samples with a systematic approach. Code samples that can, in-turn, be further analyzed.

**Corpus** Understanding the labor involved in developing a corpus that suits the needs of a researcher, we have built an elegant framework and corpus that contains over "x" exploits and "y" vulnerability samples that is made publicly available. We hope to gain community support in order to further expand our corpus size and diversity of exploits. In order to have a publicly available corpus that can be sustained in the academic community, we provide a framework that is easy to maintain, quick to use, automated, fully compartmentalized, and scalable. Through the use of abstraction and encapsulation, the management of application, environment, and exploitation has been streamlined in BugBox. The design is intended to be make it practical to manage a large database of exploits, along with their target environments. Our corpus is more than an exploit repository, it contains the software and its dependencies in which a vulnerability exists, the configuration of the software in it's vulnerable state, exploit code that will trigger the security breach, and data that describes any distinguishing attributes of the bug. [define attributes further]

**Benefiting Projects** Our corpus would have immediate benefit for testing static analysis tools,<sup>2</sup> penetration

testing and training security teams with vulnerability injection,<sup>3</sup> and computing attack surface metrics.<sup>4</sup>

## 2 Empirical Vulnerability Research

A key question in empirical vulnerability research is how to determine which functionality, or piece of code creates a vulnerable condition in a program. In studying this problem, sometimes referred to as *vulnerability localization*, it is necessary to have a large quantity of structured data in which one can formulate and test hypotheses. Potential approaches to vulnerability research include static and dynamic analysis of the software, therefore it is important that the corpus include both the vulnerable program's code, and some way to perform analysis as the system is compromised.

We use a trace based collection approach in order to capture code vulnerabilities in software. Benefits of a trace based approach include: a global approach to collection.

A global approach to collection can be more beneficial than using static analysis tools to find vulnerabilities.

[many things with references go here]

Demonstrate that BugBox ....

Predicting which metrics

Code inspection

Trace-based vulnerability definitions

System taint analysis

“Vulnerability localization” Characteristics:

1. Collecting Vulnerable Application source code
2. Collecting vulnerability details (line-based, run-based, trace-based)

## 3 Requirements

BugBox is designed to work with the Debian GNU/Linux distribution and compatible distributions. It can be distributed as a self-contained virtual machine, or as a package that can be installed on an existing system. The machine must have sufficient storage (roughly 4 GB per OS environment, and 2 GB for the application, engine, and exploit sources [Confirm these #s]). Dependencies for BugBox include MySQL, Selenium Server, debootstrap, and the Advanced Packaging Tool (APT.)

## 4 Environment

[Why we need different OS environments, how the corpus manages them]

The corpus contains a main engine, written in python, to load a web application, gather traces, run exploits, and cleanup. We provide an engine to abstract away the details of the setup process so that contributors can focus on creating new exploit scripts. This separation accomplished two tasks, it makes developing for the corpus more approachable for programmers that want to contribute and provides a robust framework that can provide more of a challenge to seasoned programmers.

We use a virtual jail, the linux chroot environment. Staging the application in the chroot environment provides locality, reproducibility, and stability.

**The web application remains local.** Many applications, such as wordpress, have a MYSQL backend. If we were to host the application on a different server or VM we would have to ensure that the MYSQL database is properly setup each time. With the chroot environment, the process is simplified, we keep a MYSQL database outside of the CHROOT Jail in order to facilitate MYSQL access.

**Current tests are independent from future tests.** We load a clean application into the chroot jail every time we wish to run a new test. This ensures that there is no corruption of the original Web Application and provides reproducible results when testing.

**The web application cannot contaminate our testing environment.** If a web application crashes due to the malicious script, we can ensure that it does not crash our corpus environment, worst case the chroot is corrupted, and we can ensure that the crash cannot have un-intended side effects in our testing environment.

### 4.1 File Structure

The file structure mimics the modularity of our corpus. Currently we store backups of the Web Applications in the /backups directory and the actual applications in the /packages directory. We store a separate backup in order to verify that the selenium scripts did not modify the Web Application in any way. If the Web Application is corrupt, then we simply copy the web application from the backup folder to the package folder.

Currently as simple incremental backup system is used to ensure that each chroot jail remains un-tainted. Periodically, an original instance of a chroot system will overwrite the

Future implementations will store only the md5 hash of the program on the corpus server and the backups on an external server. This will allow for a quick comparison followed by a remote copy if necessary. This sep-

aration will ensure that backups will never be corrupted through use of the corpus.

All the exploits that are available are located under the /scripts folder. While script is referenced as Selenium Script throughout the document, here scripts also include services such as: deploying the web application in the chroot jail, starting the trace collection, and running the script. All mounted chroot environments reside in /live\_systems.

insert file system diagram

## 5 Framework

The BugBox framework has five main components: web application store, testing environment, testing scripts, and host system. Figure 1 illustrates the general structure of the corpus environment, with arrows showing lines of control or communication. BugBox is a modular framework design where individual components can be interchanged without effecting the overall functionality of the framework.

At it's core, the organization of the vulnerability engine is driven by the Python module and package system. The system breaks down into the following four python modules:

1. corpus.**Engine**
2. corpus.**Targets**
3. corpus.**Exploit**
4. corpus.**SeleniumDriver**

The **Engine** drives the environment setup, tear-down, and exploitation process. It contains most of the logic for doing work with exploits. The **Targets** module has as submodules each application and application plugin that are associated with an exploit. **Exploit** is the superclass for each exploit in the corpus, defining interfaces and attributes that the engine uses to manage the environment and exploitation. The **SeleniumDriver** is a wrapper class for the Selenium's Firefox web driver.

### 5.1 Vulnerability Engine

### 5.2 Target Module

The anatomy of a "Target" module

### 5.3 Exploit

Management actions, recording commands attributes,

```
__init__()
exploit()
```

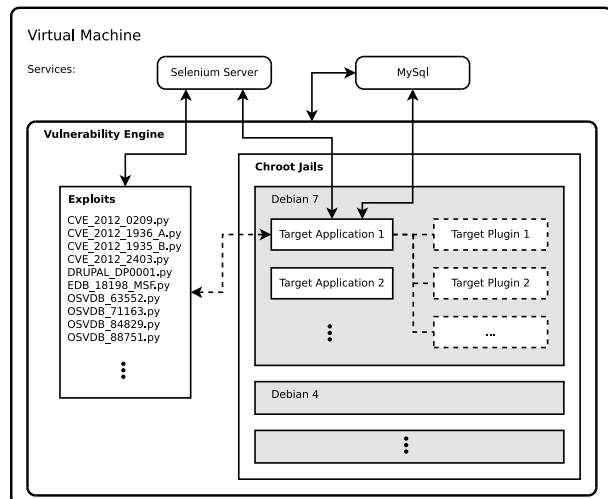


Figure 1: System Diagram

```
check() [not yet]
cleanup()
```

#### 5.3.1 As a Standalone Application

Each exploit is defined in it's own python file as a module. In Python, modules can either be imported from other modules, or executed directly from the command-line. The exploits written for BugBox take advantage of this property to give the researcher the option of writing scripts to do work on a set of exploits, or to invoke one exploit at a time. When run from the command-line, the supplied arguments are passed directly to the vulnerability engine, which supports the options shown below:

Usage: python OSVDB\_89960.py [options]

Options:

start:	Start exploit instance
stop:	Stop exploit instance
exploit:	Run the exploit
check:	Check if the corresponding environment is running
xdebug_on:	Turn on xdebug autotrace
xdebug_off:	Turn off and collect xdebug autotrace

#### 5.3.2 As a Python Module

Scripting of experiments. corpus.Query [Work in progress] interface for using queries of exploit metadata to manage sets of exploits/environments.

```

import corpus

class Exploit (corpus.Exploit):

    def __init__(self):
        corpus.Exploit.__init__(self, {
            'Name' : "CVE_2012_2403",
            'Description' : "Creates a post containing a XSS payload.",
            'References' : [['CVE', '2012-2403'],
                           ['OSVDB', '81463']],
            'Target' : "Wordpress 3.3.1",
            'Type' : "XSS",
            'VulWikiPage' : "http://seamster.cs.umd.edu/CVE-2012-2403"
        })

    return

    def exploit(self):

        payload = "<a href=\"#\" title=\"XSS http://example.com/onmouseover\"
            \"=eval (unescape (/ %61%6c%65%72%74%28%31%29%3b%61%6c%65%72\"
            \"%74%28%32%29%3b%61%6c%65%72%74%28%33%29%3b/.source)) //\"
            \">XSS</a>\"

        driver = self.create_selenium_driver()
        driver.get ("http://localhost/wordpress/?p=1")
        driver.find_element_by_id("author").clear()
        driver.find_element_by_id("author").send_keys("selenium script")
        driver.find_element_by_id("email").clear()
        driver.find_element_by_id("email").send_keys("selenium@python.org")
        driver.find_element_by_id("url").clear()
        driver.find_element_by_id("url").send_keys("www.python.org")
        driver.find_element_by_id("comment").clear()
        driver.find_element_by_id("comment").send_keys(payload)
        driver.find_element_by_id("submit").click()
        driver.cleanup()

    return

if __name__ == "__main__": # Exploit invoked from command-line
    engine = corpus.Engine(Exploit())
    engine.parse_args(sys.argv)

```

### 5.3.3 Selenium Scripting

We aim to gather code that is vulnerable. We drastically can reduce the amount of computation needed by reducing the code size that needs to be analyzed through the use of past, found exploits.

In order to better isolate vulnerable code, we create a selenium script in Python, replicating actions a user would take in performing an exploit on a web application. While the script is running, we can have XDebug monitor the execution and output the final execution trace

for the selenium script execution. We can see future additions to this by turning XDebug on and off through cookie manipulation while the script is running, furthermore reducing the size of the vulnerable code.

By setting up the application loader engine, we are able to create concise python selenium scripts for data collection. The following code shows how easy it is to run an automated exploit:

Pros: -Ease of use, great python bindings -Uses browser's javascript engine -Good for Demonstration/visualization

Cons: -Not a necessary dependency, urllib+cookielib would suffice -requires SeleniumServer to run as a service on the host system -Slower than crafting all requests directly -Many times, we still need to use auxiliary libraries to send specially crafted requests anyway

[figure of a session hijack in a web-browser?]

We are restricted to web applications because of the use of selenium scripts in our corpus. [GJN ADD NOTES FROM JEFF's PAPER ON THIS]

## 5.4 Trace Collection

**Stuff about interaction with XDebug.** XDebug is a feature-rich PHP debugger that can be used to easily collector traces enabling various run-time analyses.

## 6 Scalability

In order to create a system that can handle a growing amount of vulnerabilities that are independent from one another, we represent an exploit on an application through Selenium scripts, as noted in the Framework section. Selenium scripts allow us to perform operations on web applications without fear of working in an unintended, compromised environment.

## 7 Use Cases

By providing a corpus that explicitly logs the steps taken in accumulating the log files, we have more flexibility. This flexibility can be seen with the following example:

Jon is told by his advisor that he needs to collect more trace data in order to get a proper sample size for his research. John quickly creates a selenium script for the exploit he wants to collect and shows the advisor his results. The advisor was generally happy with the trace data that he collected, but instead wanted him to do a slight modification to the exploit. If Jon did not have the selenium script at hand, he would have to duplicate all of the work previously done. Since he does have the script on hand, he can quickly make a change to the script and re-run in seconds versus hours.

While the above process is only shown in one iteration, most students know that this is not the case. One hour of work can turn into a whole week of work without

the proper framework in place. The above situation also shows that the selenium script can be discussed with the professor to show validity of the data and provide talking points for how the exploit was applied.

## 8 Acknowledgements

Metasploit, undergraduate summer labor, etc..

Now we're going to cite somebody. Watch for the cite tag. Here it comes. The tilde character (~) in the source means a non-breaking space. This way, your reference will always be attached to the word that preceded it, instead of going to the next line.

## 9 Future Development

**Distribution** Virtual machine v.s. debian package. Pre-built chroot jails v.s. build scripts. i.e. size vs. setup process balance. Striking a balance between

**Services** Why run Selenium/Mysql in VM vs a chroot jail?

**Isolating attack event** (with xdebug manually/cookies/etc...)

**Selenium driver and aux modules** Explore the possibility of a unified communication interface. This may be necessary in order to cleanly interact with xdebug with appropriate cookies set on a per-request basis (especially when modules other than Selenium are used for communication).

**Payload standardization** For each exploit currently in the corpus, there is no standard for the payload used in the attack. Since many studies may be sensitive to the payload type and encoding, it makes sense to provide the researcher with fine-grained control over this property. The Metasploit Framework has a very robust system for managing exploits along with their payloads and encodings, and can be a model for implementing this.

## 10 Conclusion

Everything is great, we just need community involvement to mature the framework and to help write more exploits.

## 11 Availability

[Available as a 10 GB virtual machine and as a standalone debian package?]

`git://bugbox.github.com/blahblah`

<http://www.vulnerabilitywiki.com>

## References

## Notes

<sup>1</sup>Zero Comments: Blogging and Critical Internet Culture (09 August 2007) by Geert Lovink

<sup>2</sup>M. Zitser, R. Lippmann, and T. Leek, "Testing Static Analysis Tools using Exploitable Buffer Overflows from Open Source Code," ???

<sup>3</sup>J. Fonseca, M. Vieira and H. Maderia, "Traning Security Assurance Teams using Vulnerability Injection," 2008 14th IEEE Pacific Rim International Symposium on Dependable Computing

<sup>4</sup>J. Stuckman and J. Purtilo, "Comparing and Applying Attack Surface Metrics," ???