



## E-MAIL HEADER INFORMATION

Joshua Clark

Fourth Year Project Report for the Final Honour  
School of Computer Science

May 2016

# Abstract

After extensive public education, fewer people are now clicking on links in e-mails that are disguised as phishing attacks, though the threat still remains, and considerable amounts of work has gone into exploring the demographics most likely to be targeted. As the number of technically literate people grows, this sort of attack is increasingly unlikely to be successful. Therefore, malicious entities are more likely to attempt to attack people based on the information leaked in their emails, and more specifically, the header, which most people are less likely to have some degree of control over.

The risks are not just limited to individual users, and at a corporate level, the risks posed by leaking information through e-mails could be even greater: e-mail headers can reveal the internal network structure of a company's computer systems as well as the different pieces of software that are running inside the system. Extracting the social information could be of great value for executing a phishing attack, however, there is also value in determining the specific weaknesses in a system. This can be aided through the use of vulnerability databases.

This report discusses the existing research into the information leaked by e-mail headers and presents a tool to extract such information.

# Acknowledgements

I want to thank my supervisor, Dr Jason R.C. Nurse for his assistance throughout the year; my tutor, Professor Peter Jeavons, for his unfailing help and support throughout my time at Oxford. I would like to acknowledge the support from my family and partner, Agata, for encouraging me.

# Contents

|          |                                             |           |
|----------|---------------------------------------------|-----------|
| <b>1</b> | <b>Introduction</b>                         | <b>8</b>  |
| 1.1      | Motivation . . . . .                        | 8         |
| 1.2      | Aims and Objectives . . . . .               | 8         |
| 1.3      | Typical Use . . . . .                       | 8         |
| 1.4      | Document Conventions . . . . .              | 9         |
| 1.5      | Structure . . . . .                         | 9         |
| <b>2</b> | <b>Literature Review</b>                    | <b>10</b> |
| 2.1      | General Data Leakage . . . . .              | 10        |
| 2.1.1    | Personal Data Leakage . . . . .             | 10        |
| 2.1.2    | Corporate Data Leakage . . . . .            | 10        |
| 2.2      | Data Leakage from E-Mails . . . . .         | 11        |
| 2.2.1    | E-Mail Headers . . . . .                    | 11        |
| 2.2.2    | Example Header and Pertinent Data . . . . . | 11        |
| 2.2.3    | Existing Research . . . . .                 | 12        |
| 2.3      | Existing Tools . . . . .                    | 12        |
| 2.3.1    | Google . . . . .                            | 13        |
| 2.3.2    | Microsoft . . . . .                         | 13        |
| 2.4      | Vulnerabilities . . . . .                   | 13        |
| 2.4.1    | CVE Mitre Lookup . . . . .                  | 14        |
| 2.4.2    | Norton Vulnerability Protection . . . . .   | 14        |
| 2.5      | Overview . . . . .                          | 14        |
| <b>3</b> | <b>Implementation</b>                       | <b>16</b> |
| 3.1      | Overview . . . . .                          | 16        |
| 3.2      | Definitions . . . . .                       | 16        |
| 3.2.1    | Parsing . . . . .                           | 16        |
| 3.2.2    | Database Queries . . . . .                  | 18        |
| 3.3      | Data Extraction and Parsing . . . . .       | 18        |
| 3.3.1    | Received fields . . . . .                   | 19        |
| 3.3.2    | Other fields . . . . .                      | 19        |
| 3.4      | Analysis . . . . .                          | 20        |
| 3.4.1    | Text-Based . . . . .                        | 20        |
| 3.4.2    | Database Queries . . . . .                  | 20        |
| 3.5      | Visualising the Results . . . . .           | 21        |
| <b>4</b> | <b>Evaluation</b>                           | <b>22</b> |
| 4.1      | Methodology . . . . .                       | 22        |
| 4.2      | Results . . . . .                           | 22        |
| <b>5</b> | <b>Conclusions and Future Work</b>          | <b>23</b> |
| 5.1      | Conclusions . . . . .                       | 23        |
| 5.2      | Future Work . . . . .                       | 23        |

**List of Tables**

1.1   Format of presented data found in e-mail header   . . . . .   9

## List of Figures

|     |                                                    |    |
|-----|----------------------------------------------------|----|
| 2.1 | Google Apps Toolbox E-mail header output . . . . . | 13 |
| 2.2 | Microsoft E-mail header output . . . . .           | 14 |
| 2.3 | CVE Search Results . . . . .                       | 15 |
| 2.4 | Norton Vulnerability Protection Results . . . . .  | 15 |

# List of Algorithms

|   |                                          |    |
|---|------------------------------------------|----|
| 1 | Lookup based on a known key . . . . .    | 20 |
| 2 | Lookup based on a key property . . . . . | 20 |
| 3 | Extracting CVE entries . . . . .         | 21 |

# 1 Introduction

## 1.1 Motivation

E-mail systems are now so integrated into our modern lives that we struggle to cope without them. E-mails ubiquity is also one of its largest weaknesses, a fact recognised very early on. The first spam email was sent in 1978, as documented by Templeton n.d. After spam came phishing, first described by Felix and Hauck 1987, with the first-real world use being against the customers of America Online, an ISP. However, this still relies on the targets providing their data for malicious purposes. One of the first e-mail viruses to spread was the Happy99 virus, which, other than propagating itself, had no other effect on infected systems. Later viruses would target credit-card and banking information. However, all of these techniques rely on the malicious email being received and its contents being opened. There are fewer instances recorded, however, of the information flow being sent the other way. A more subtle attack will focus on the information being sent from a legitimate user to an attacker. It is easy enough for an individual to read an e-mail header and identify interesting elements, however, on a large scale, this quickly becomes more difficult.

## 1.2 Aims and Objectives

This project aims to develop a tool that can be used as described above to automatically extract information from e-mail headers and analyse its results to display the personal information contained within an e-mail's header, as well as information about the software configurations that may be found on a user's computer, or the servers used to send their e-mail.

The program would be expected to satisfy the following minimal requirements in order for it to be considered successful:

**Accuracy** — any information produced by the parser should be reflective of the input e-mail

**Representation** — the produced visualisation should be intuitive to read: each element should be presented separately from the others, and clearly labelled.

**Portability** — the visual output produced by the program should be available to the user in a variety of formats.

**Interactivity** — the program should produce sensible warnings when an e-mail that is not possible to parse has been entered.

## 1.3 Typical Use

On starting the application, the user will provide an e-mail that they wish to have analysed. This will then be parsed, and some relevant information presented in a table.

Lastly, an option is available to view the information about security vulnerabilities in a separate webpage, forming the main output of the program.



| Email Header Information                                                     |                 |                                                                    |
|------------------------------------------------------------------------------|-----------------|--------------------------------------------------------------------|
| Sender Information<br>(Name, originating domain)                             | Sender Software | Sender Usernames<br>(Presented as a list with likely organisation) |
| Graphical representation of devices used to deliver the e-mail               |                 |                                                                    |
| List of derived information including found software and similar information |                 |                                                                    |
| Histograms for vulnerability scores, separated by product                    |                 |                                                                    |
| Searchable and filterable table of discovered CVEs                           |                 |                                                                    |

Table 1.1: Format of presented data found in e-mail header

The resultant webpage will be structured as in table 1.1. It will then be possible for the user to click on the representations of the devices to find out more information. It will also be possible to search within the vulnerability list to find more information, as well as filter by impact and availability details.

## 1.4 Document Conventions

By convention, when class diagrams are used, ovals will represent traits/abstract classes, with rectangles representing concrete implementations. Objects are indicated using bold lines in the diagram, and are similar in behaviour to statically declared objects in many languages.

Whenever **this font face** is used, the text is referring to either some implementation level object or class, or text entered into the application or used for testing.

## 1.5 Structure

Chapter 2 begins by discussing the existing research on the subject as well as existing publicly-available tools to analyse headers. I then use these as a basis to discuss features that would be expected to appear in a header analyser looking for leaked information and vulnerabilities.

The implementation's high-level structure and details will be discussed in Chapter 3, and algorithms presented in pseudo-code where necessary. A full listing will be presented at the end of this document in an appendix. The results of the analysis of the headers will be discussed in Chapter 4, beginning with the methodology used, and presenting a number of results. Finally, Chapter 5 will discuss my conclusions and areas of further improvement.

## 2 Literature Review

In this chapter, we will discuss the nature of existing threats to data, and the ongoing research in this area. We will then consider the specific threats posed by e-mail.

### 2.1 General Data Leakage

The importance of data leakage is gaining more importance as the amount of information stored about entities increases, and the risks are being considered more carefully. From the obvious ramifications for businesses discussed in Papadimitriou and Garcia-Molina 2011: the loss of trust and legal action resulting from the discovery of leaking data, to the more personal issues discussed in Irani et al. 2011: the possibility of using the discovered data to discover passwords or to physically identify them. 53% of Americans can be uniquely identified by their birthdate, gender and location (city/town), with the number jumping 87% when using birthdate, gender and zip code.

#### 2.1.1 Personal Data Leakage

From a persona perspective, there are a number of risks. There are a significant number of social networks available, with an estimated 1.65 billion monthly active users, with a significantly higher proportion used in developed countries. Irani et al. 2011 showed the rate at which the information gathered from social networks can be used to uniquely identify an individual, with only 9 sites required before there's approximately a 70% chance that both a person's hometown and name could be recovered. A similar number of sites can give an aggregate normalised attribute leakage of 1, where the leakage is defined as

$$\Psi(F_a, P) = \frac{\sum_{f_a \in F_a} \phi(f_a, P)}{|F_a|} \text{ where } \phi(f_a, P) = [f_a \in P]$$

#### 2.1.2 Corporate Data Leakage

When companies receive user data, they often have a legal obligation to ensure that the data is protected and treated as confidential and sensitive. When this trust is broken, there are often severe consequences, both from regulators and consumers moving their business to competitors. Squicciarini, Sundareswaran, and Lin 2010 considers one way that data may be leaked, despite care being taken to ensure that it is properly encrypted and stored. By failing to protect against the indices for databases being stored insecurely, customer data may be leaked. Order-preserving encryption schemes, as described in Agrawal et al. 2004, is one way of solving this problem, to an extent.

## 2.2 Data Leakage from E-Mails

### 2.2.1 E-Mail Headers

All e-mails include additional information about the sender and receiver, some of which is used by an e-mail client in order to display more information about the message that is currently being viewed, such as its original sender, reply-to addresses and the time it was sent. Additional fields allow senders to authenticate themselves using public-key methods.

The format of e-mail headers was first defined in RFC 822, and further refined in subsequent RFCs. The standard for e-mails was then formalised precisely in RFC 5322.

### 2.2.2 Example Header and Pertinent Data

In the example below, and text highlighted with red, like so is information about the receiver. Information about the sender, their hardware or software is highlighted in green, like so; and information gathered about intervening devices is highlighted in blue, like so.

```
Delivered-To: joshuaclark94@gmail.com
Received: by 10.25.150.146 with SMTP id y140csp543137lfd;
  Sat, 6 Feb 2016 08:49:56 -0800 (PST)
X-Received: by 10.112.12.2 with SMTP id u2mr8302831lbb.145.1454777396580;
  Sat, 06 Feb 2016 08:49:56 -0800 (PST)
Return-Path: <agatabor@poczta.onet.pl>
Received: from smtpo75.poczta.onet.pl (smtpo75.poczta.onet.pl. [141.105.16.25])
  by mx.google.com with ESMTPS id o199si122556361fb.94.2016.02.06.08.49.56
  for <joshuaclark94@gmail.com>
  (version=TLS1_2 cipher=ECDHE-RSA-AES128-GCM-SHA256 bits=128/128);
  Sat, 06 Feb 2016 08:49:56 -0800 (PST)
Received-SPF: pass (google.com: domain of agatabor@poczta.onet.pl designates
  141.105.16.25 as permitted sender) client-ip=141.105.16.25;
Authentication-Results: mx.google.com;
spf=pass (google.com: domain of agatabor@poczta.onet.pl designates 141.105.16.25
  as permitted sender) smtp.mailfrom=agatabor@poczta.onet.pl
Received: from [10.26.196.156] (client-8-32.eduroam.oxuni.org.uk [192.76.8.32])
  (Authenticated sender: agatabor@poczta.onet.pl)
  by smtp.poczta.onet.pl (Onet) with ESMTPA id 3pyKNH4ffyzT6tkv8
  for <joshuaclark94@gmail.com>; Sat, 6 Feb 2016 17:49:50 +0100 (CET)
Date: Sat, 06 Feb 2016 16:49:07 +0000
Subject: Test e-mail
Message-ID: <j66i9tkyhy3l4v77erlw1gne.1454777347191@email.android.com>
Importance: normal
From: Agata <agatabor@poczta.onet.pl>
To: Joshua Clark <joshuaclark94@gmail.com>
MIME-Version: 1.0
Content-Type: multipart/alternative;
  boundary="--_com.android.email_1892258509098440"

-----_com.android.email_1892258509098440
```

```
Content-Type: text/plain; charset=utf-8
Content-Transfer-Encoding: base64
```

```
CgoKC1N1bnQgZnJvbSBteSBTYW1zdW5nIEdhbGF4eSBzbWVydHBob251Lg==
```

```
-----com.android.email_1892258509098440
Content-Type: text/html; charset=utf-8
Content-Transfer-Encoding: base64
```

```
PGh0bWw+PGhlYWQ+PG1ldGEgaHR0cC1lcXVpdj0iQ29udGVudC1UeXB1IiBjb250ZW50PSJ0ZXh0
L2h0bWw7IGNoYXJzZXQ9VVRGLTgiPjwvaGVhZD48Ym9keT48ZG12IHNOeWxlPSJ3b3JkLWJyZWFr
O2t1LXAuYXN0eWw+PGJyPjxicj48YnI+PGJyP1N1bnQgZnJvbSBteSBTYW1zdW5nIEdhbGF4eSBz
bWVydHBob251Ljxicj48L2Rpdj48L2JvZG1sPg==
```

```
-----com.android.email_1892258509098440--
```

The particularly interesting portions of the e-mail header include the IP addresses of the various servers the message has travelled through, allowing their approximate location to be determined. Additionally, the information on the protocol being used and the software being run allows for anyone with access to mail headers to find more information about the attacks a device and its software may be vulnerable to.

### 2.2.3 Existing Research

In Nurse et al. 2015, the idea of using the information available in an email header was mooted, turning the previously standard threat of malware and phishing contained in received e-mails on its head, and instead presenting the threat in outgoing emails, and the personally identifying information (PII) contained therein. Many emails leaked information about employers, e-mail services and applications used, and IP address. Initial examination of a variety of e-mail headers found within my own inbox also revealed a plethora of information, including phone carriers, preferred languages, and system usernames. It is conceivable therefore, that it is possible to automate at least part of this, and present the information that can be extracted, in a white-hat tool to allow people to audit the information that they are revealing. The obvious malicious use-case involves using such information as part of a spear-phishing exercise.

An alternative vulnerability presents itself in the information about systems that may be revealed. Many email clients embed identifying information, and there are multiple databases available to allow specific threats to be identified. This could allow a malicious entity to compromise the security of a target machine, and gain access to the data stored on that machine and available on any connected network devices. Work started in Joshi, Lal, and Finin 2013 discusses the need to aggregate data about vulnerabilities from multiple sources to present a more complete and coherent picture, which is also likely to then contain more accurate data.

Al-zarouni 2004 presents an alternative set of results, describing how an individual can seek to protect themselves against malicious e-mails, using the contents of e-mail headers. Various discrepancies between forged e-mail addresses and legitimate messages are described.

## 2.3 Existing Tools

Several tools already exist online to display the information that is found in e-mail headers. Tools from Microsoft and Google exist to analyse the contents of e-mail headers. These tools

|                    |                                                                     |  |  |  |  |  |
|--------------------|---------------------------------------------------------------------|--|--|--|--|--|
| <b>MessageId:</b>  | <201107031502.p63F2i2m001182@nyork.iii.com>                         |  |  |  |  |  |
| <b>Created at:</b> | Sun Jul 03 2011 08:02:18 GMT-0700 (PDT) ( Delivered after 19 mins ) |  |  |  |  |  |
| <b>From:</b>       | New York Public Library                                             |  |  |  |  |  |
| <b>To:</b>         | some_random_user@gmail.com                                          |  |  |  |  |  |
| <b>Subject:</b>    | New York Library items due soon.212-555-2329                        |  |  |  |  |  |

| # | Delay   | From                  |   | To                     | Protocol | Time received                           |
|---|---------|-----------------------|---|------------------------|----------|-----------------------------------------|
| 0 |         | localhost.localdomain | → | nyork.iii.com          | ESMTP    | Sun Jul 03 2011 08:02:18 GMT-0700 (PDT) |
| 1 | 19 mins | nyork.iii.com         | → | sienna.pobox.com       | ESMTP    | Sun Jul 03 2011 08:21:05 GMT-0700 (PDT) |
| 2 | 1 sec   | localhost             | → | sienna.pobox.com       | ESMTP    | Sun Jul 03 2011 08:21:06 GMT-0700 (PDT) |
| 3 |         | sienna.pobox.com      | → | [google] mx.google.com | ESMTP    | Sun Jul 03 2011 08:21:06 GMT-0700 (PDT) |
| 4 | 1 sec   |                       | → | [google] 10.52.65.169  | SMTP     | Sun Jul 03 2011 08:21:07 GMT-0700 (PDT) |
| 5 | 3 sec   |                       | → | [google] 10.229.234.71 | SMTP     | Sun Jul 03 2011 08:21:10 GMT-0700 (PDT) |

Show Raw header

Figure 2.1: Google Apps Toolbox E-mail header output

clearly display the information displayed in the header, showing the key-value pairs, and the set of servers the message transferred through and the protocols used.

### 2.3.1 Google

The Google Apps Toolbox features an e-mail header analyser<sup>1</sup>. An example of the output of the utility is found in figure 2.3.1.

One of the most useful features from the Google Apps Toolbox is the information provided about the servers the message travelled through. This tool shows the details of the time taken for each hop, and the protocol used.

### 2.3.2 Microsoft

The Microsoft Message Header Analyzer<sup>2</sup> and showing sample results in figure 2.3.2

## 2.4 Vulnerabilities

Beware of bugs in the above code; I have only proved it correct, not tried it.

Donald Knuth

Problems in software are nothing new, and seeking to exploit these issues is almost as old. As the security implications behind flawed software became more widely recognised, reducing their impact wherever possible became the next most important step. The MITRE Corporation operates the National Cybersecurity Federally Funded Research and Development Centre,

<sup>1</sup>Found at <https://toolbox.googleapps.com/apps/messageheader/>

<sup>2</sup>Found at <https://testconnectivity.microsoft.com/MHA/Pages/mha.aspx>

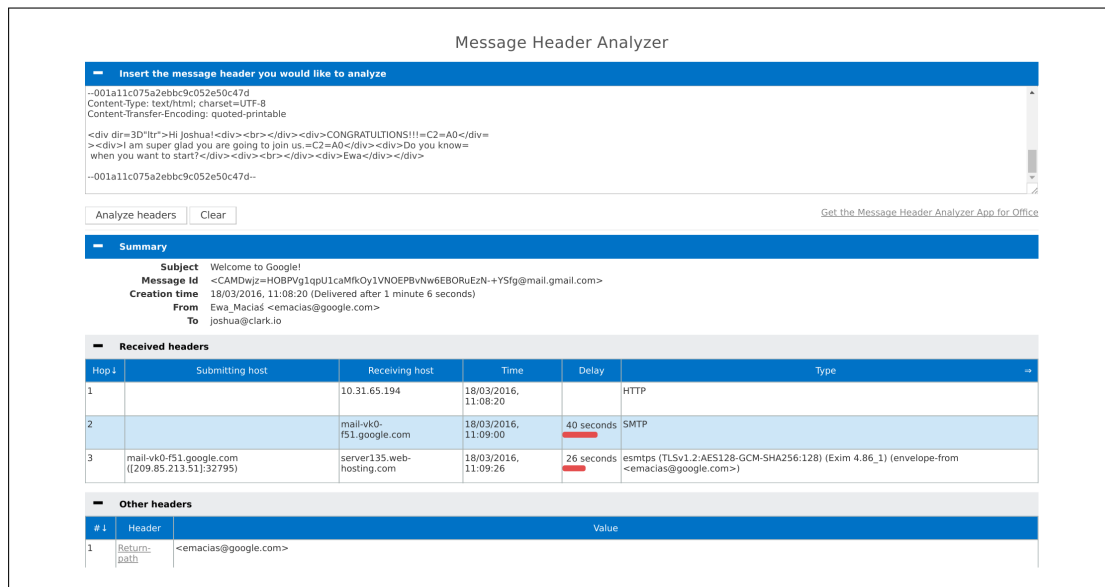


Figure 2.2: Microsoft E-mail header output

which exists to maintain a database of these vulnerabilities, which are referred to as Common Vulnerabilities and Exposures (CVE).

### 2.4.1 CVE Mitre Lookup

There are a number of tools to look up CVEs<sup>3</sup> and showing sample results in figure 2.4.1. There are a number of limitations to the results returned by the CVE Mitre tool. Firstly, little context is returned: information about scores, the impact and access information are omitted, for example. Additionally, the process of finding relevant vulnerabilities is further slowed down by the necessity to search for specific terms one at a time. Additionally, automated tools exist at a consumer and enterprise level that will automatically scan a computer or network to detect installed software configurations and show the results.

### 2.4.2 Norton Vulnerability Protection

For example, the now deprecated Norton Vulnerability Protection tool, as shown in Figure 2.4.2<sup>4</sup> lists the programs and the total number of vulnerabilities found, providing more information on each program. This method has the advantage of indicating the specific programs that have vulnerabilities, with the aim of allowing a user to update their vulnerable applications, however it does not allow for more fine-grained information.

## 2.5 Overview

<sup>3</sup>Fount at <https://www.cve.mitre.org/find/index.html>

<sup>4</sup>Available at [community.norton.com](https://community.norton.com)

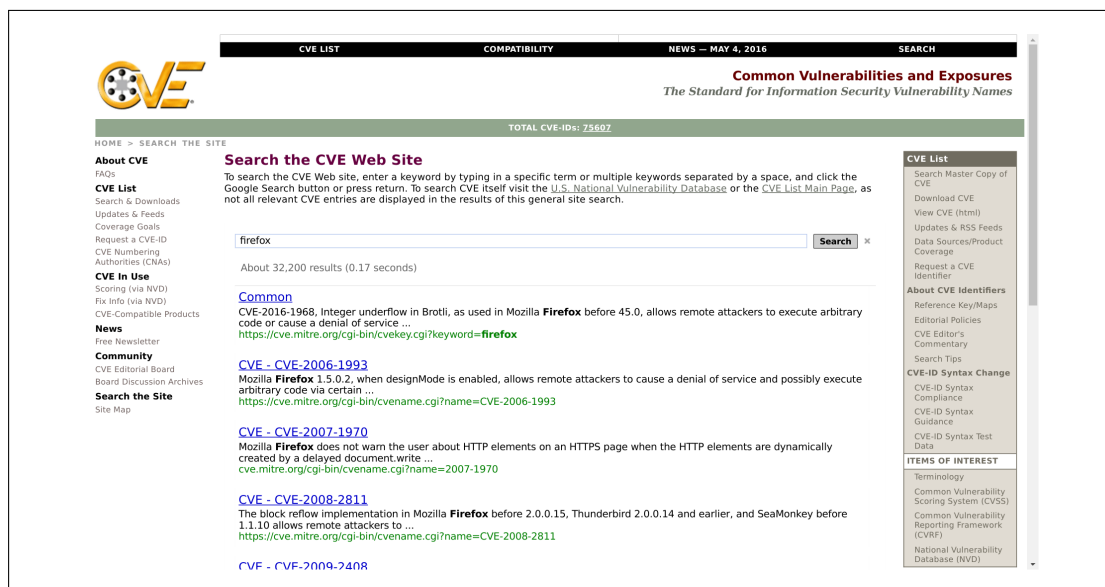


Figure 2.3: CVE Search Results

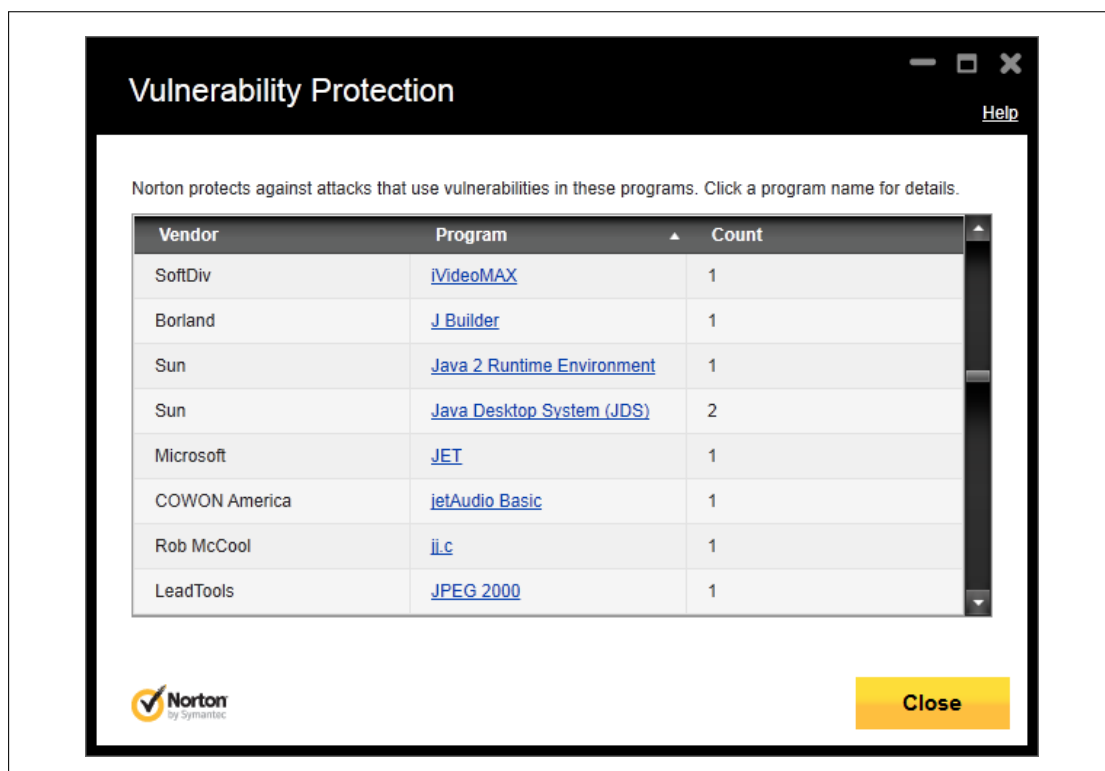


Figure 2.4: Norton Vulnerability Protection Results

## 3 Implementation

### 3.1 Overview

The analysis is implemented as a series of stages, firstly, the e-mail header is parsed, to extract important information to a predefined set of Java objects. This is followed by the analysis phase, where the resultant data is passed to a set of analyser modules, each running separately. Finally, this information is presented to the user. This chapter presents each of these stages in detail.

### 3.2 Definitions

The following covers the essential definitions required for the notation and concepts that will be discussed in this document.

#### 3.2.1 Parsing

In order to aid the parsing of the e-mail header, a combination of regular expressions and context-free grammars are needed, and defined as follows.

**Alphabets and Languages** A set of symbols, usually denoted as  $\Sigma$ . A language is a subset of  $\mathcal{P}(\Sigma)$ .

The following special classes are provided as part of the Perl-Compatible Regular Expression library, and are subsets of the alphabet of Unicode characters, defined in Group et al. n.d.

**alnum** — letters and digits

**alpha** — letters

**ascii** — the set of ASCII characters (character codes 0 — 127)

**blank** — tabs or blank spaces

**cntrl** — control characters

**digit** — decimal digits

**graph** — printing characters (excluding spaces)

**lower** — lower-case letters

**print** — printing characters (including spaces)

**punct** — punctuation marks (printing characters excluding letters and spaces)

**space** — white space

**upper** — upper case letters



**word** — “word” characters (same

**xdigit** — hexadecimal digits

## Regular Languages

Regular languages are defined as follows:

- $\emptyset$  and  $\{\epsilon\}$  are regular languages
- for each  $a \in \Sigma$ ,  $\{a\}$  is a regular language
- if  $A$  and  $B$  are both regular,  $A \cup B$ ,  $A \cdot B$  and  $A^*$  are regular languages.
  - $A \cup B$  is the union of two languages.  $A \cup B = \{s : s \in A \vee s \in B\}$
  - $A \cdot B$  is the concatenation of two languages.  $A \cdot B = \{ab : a \in A, b \in B\}$
  - $A^*$  is the Kleene star of a language.

$$\begin{aligned} A_0 &= \{\epsilon\} \\ A_1 &= A \\ A_{i+1} &= \{aa' : a \in A_i, a' \in A\} \\ A^* &= \bigcup_{i \in \mathbb{N}} A_i \end{aligned}$$

## Context-Free Grammars

A context-free grammar  $G$  is defined as  $G = (V, \Sigma, R, S)$  where:

- $V$  is a variable.
- $\Sigma$  is the alphabet of symbols.
- $R$  is a relation defined over  $V \rightarrow (V \cup \Sigma)^*$
- $S$  is the start symbol

For example,  $\langle S \rangle$  is the field name with the associated productions  $\langle T \rangle \langle U \rangle$ , where  $T$  and  $U$  are productions.

$$\langle S \rangle \models \langle T \rangle \langle U \rangle$$

For example,  $\langle S \rangle$  is the field name with the associated productions  $a \langle U \rangle$ , where  $a$  is a terminal symbol.

$$\langle S \rangle \models a \langle U \rangle$$

This is then extended in the following ways used in the RFC syntax.

The square brackets are used to indicate an optional element.

$$\langle \text{field} \rangle \models \langle \text{field-name} \rangle : [\langle \text{field-body} \rangle] \text{ CRLF}$$

The asterisk is used to indicate an element that appears 0 or more times.  $n^*$  is used to indicate a component that repeats  $n$  or more times.

$$\langle \text{fields} \rangle \models \langle \text{dates} \rangle \langle \text{source} \rangle 1^* \langle \text{destination} \rangle * \langle \text{optional-fields} \rangle$$

The hash-symbol is used to indicate an element that appears a certain number of times.  $m*n$  is used to indicate a component that repeats at least  $m$  times and at most  $n$  times.

$$\langle \text{fields} \rangle \models \langle \text{dates} \rangle \langle \text{source} \rangle 1\#\langle \text{destination} \rangle * \langle \text{optional-fields} \rangle$$

The  $|$  is used to indicate a selection between a pair of elements.

$$\langle \text{fields} \rangle \models a \mid b$$

### 3.2.2 Database Queries

The following notations will be used for the CVE database queries.

**Set-Theoretic Operators** The operators  $F \cup G$ ,  $F \cap G$ ,  $F \setminus G$  behave as is expected for these operators, resulting in the union, intersection and difference of the sets. The only proviso being that the attribute names must match.

#### Selection

$$\sigma_{\text{product}=\text{thunderbird}} D$$

The above notation is used to indicate a search over the attribute named “product” for the string “thunderbird” in the database table  $D$ . As a single database is only being used, this may be occasionally elided. The output of this function is another object of the same type as  $D$ .

#### Projection

$$\pi_{\text{product}} D$$

The above notation is used to indicate a projection on the attribute named “product” in the database table  $D$ . The output of this function is another object of the same type as  $D$ .

**Composition** The above functions results can be coposed repeatedly to produce more specific search queries.

## 3.3 Data Extraction and Parsing

The parser’s operation completes in a number of stages, following RFC822 (Crocker 1982). The header is divided up into two disjoint sections, the routing information (**Received from...**) and the key-value map of other pertinent information.

### 3.3.1 Received fields

The received fields are the most complicated part of the e-mail header to parse, as they are described by a non-trivial grammar, presented below.

|                                              |           |                                                                                                                                              |
|----------------------------------------------|-----------|----------------------------------------------------------------------------------------------------------------------------------------------|
| $\langle \text{message} \rangle$             | $\models$ | $\langle \text{fields} \rangle * (\text{CRLF} * \text{text})$                                                                                |
| $\langle \text{fields} \rangle$              | $\models$ | $\langle \text{dates} \rangle \langle \text{source} \rangle 1 * \langle \text{destination} \rangle * \langle \text{optional-fields} \rangle$ |
| $\langle \text{field} \rangle$               | $\models$ | $\langle \text{field-name} \rangle : [ \langle \text{field-body} \rangle ] \text{CRLF}$                                                      |
| $\langle \text{field-name} \rangle$          | $\models$ | <i>any word consisting of CHAR, excluding CTLs, SPACE, and “.”</i>                                                                           |
| $\langle \text{field-body} \rangle$          | $\models$ | $\langle \text{field-body-contents} \rangle [ \text{CRLF LWSP-char} \langle \text{field-body} \rangle ]$                                     |
| $\langle \text{field-body-contents} \rangle$ | $\models$ | <i>ASCII characters</i>                                                                                                                      |
| $\langle \text{source} \rangle$              | $\models$ | $[ \langle \text{trace} \rangle ] \langle \text{originator} \rangle [ \langle \text{resent} \rangle ]$                                       |
| $\langle \text{trace} \rangle$               | $\models$ | $\langle \text{return} \rangle 1 * \langle \text{received} \rangle$                                                                          |
| $\langle \text{return} \rangle$              | $\models$ | Return-path: $\langle \text{route-addr} \rangle$                                                                                             |
| $\langle \text{recieved} \rangle$            | $\models$ | Received:                                                                                                                                    |
| $\langle \text{cont.} \rangle$               | $\models$ | $[ \text{from} \langle \text{domain} \rangle ]$                                                                                              |
| $\langle \text{cont.} \rangle$               | $\models$ | $[ \text{by} \langle \text{domain} \rangle ]$                                                                                                |
| $\langle \text{cont.} \rangle$               | $\models$ | $[ \text{via} \langle \text{atom} \rangle ]$                                                                                                 |
| $\langle \text{cont.} \rangle$               | $\models$ | $* ( \text{with} \langle \text{atom} \rangle )$                                                                                              |
| $\langle \text{cont.} \rangle$               | $\models$ | $[ \text{id} \langle \text{msg-id} \rangle ]$                                                                                                |
| $\langle \text{cont.} \rangle$               | $\models$ | $[ \text{for} \langle \text{addr-spec} \rangle ]$                                                                                            |
| $\langle \text{cont.} \rangle$               | $\models$ | $; \langle \text{date-time} \rangle$                                                                                                         |
| $\langle \text{msg-id} \rangle$              | $\models$ | $< \langle \text{addr-spec} \rangle >$                                                                                                       |
| $\langle \text{addr-spec} \rangle$           | $\models$ | $\langle \text{local-part} \rangle @ \langle \text{domain} \rangle$                                                                          |
| $\langle \text{local-part} \rangle$          | $\models$ | $\langle \text{word} \rangle * ( . \langle \text{word} \rangle )$                                                                            |
| $\langle \text{word} \rangle$                | $\models$ | $\langle \text{atom} \rangle \mid \langle \text{quoted-string} \rangle$                                                                      |
| $\langle \text{domain} \rangle$              | $\models$ | $\langle \text{sub-domain} \rangle * ( . \langle \text{sub-domain} \rangle )$                                                                |
| $\langle \text{sub-domain} \rangle$          | $\models$ | $\langle \text{domain-ref} \rangle \mid \langle \text{domain-literal} \rangle$                                                               |
| $\langle \text{domain-ref} \rangle$          | $\models$ | $\langle \text{atom} \rangle$                                                                                                                |
| $\langle \text{date-time} \rangle$           | $\models$ | $[ \text{day}, ] \text{date time}$                                                                                                           |
| $\langle \text{atom} \rangle$                | $\models$ | $1 * \text{any character excluding specials, SPACE and CTLs}$                                                                                |

An example field is as follows:

```
Received: from relay12.mail.ox.ac.uk (129.67.1.163)
  by HUB05.ad.oak.ox.ac.uk (163.1.154.231)
  with Microsoft SMTP Server id 14.3.169.1;
  Sat, 14 Nov 2015 10:55:35 +0000
```

### 3.3.2 Other fields

These are read by a Python script and output to STDOUT to be read by the Java parser in a consistent format. These are then loaded into a hashmap to allow quick lookup.

## 3.4 Analysis

After completing the parsing of the field, it is then ready to be analysed for different features. All of the analysers implement the **HeaderAnalyser** interface, requiring information about the header to be analysed, and the currently running application. All of these then implement the **Runnable** interface, allowing the class to be run asynchronously.

### 3.4.1 Text-Based

The fields from the header are analysed in different modules, with searches being performed for specific strings. Of particular interest to Oxford Nexus users is the “X-Oxford-Username” string, containing the username of the individual that sent the message. As confirming the username is a fairly standard security procedure for an IT support technician, having access to this information could allow a phisher in a later stage of an attack to increase their credibility.

In some cases, the likely keys that are being searched for are known in advance, and can then be checked against the hash-map of entries.

An example of this approach is for the specific check for an Oxford username:

```
Input: Header  
Output: Any Oxford-based username that is found  
if X-Oxford-Username  $\in$  Header.KvMap then  
|   return Header.KvMap(X-Oxford-Username);  
end
```

**Algorithm 1:** Lookup based on a known key

Alternatively, we may be interested in properties of the keys, necessitating a search over the keys.

```
Input: Header  
Output: Any information relating to Microsoft Exchange that is found  
foreach Key k  $\in$  Header.KvMap do  
|   if k starts-with X-MS-Exchange then  
|   |   return Header.KvMap(k);  
|   end  
end
```

**Algorithm 2:** Lookup based on a key property

### 3.4.2 Database Queries

Using the results gathered from the text-based queries and analysis of the received fields, relevant software configurations are extracted and queried against results in the CVE database. These are then parsed and collated in preparation for displaying the outputs.

As more information is found, more details of products used will also become available. These are added asynchronously.

**Input:** Header product name  $p$   
**Output:** CVE Entries  
 $cve\_list \leftarrow \emptyset$ ;  
**foreach**  $s \in \sigma_{vector \neq LOCAL} \sigma_{product=p} D$  **do**  
     $cve\_builder \leftarrow \text{blank } cve$ ;  
     $cve\_builder.id \leftarrow \pi_{CVE-ID} s$ ;  
    ... – extract other features;  
     $cve\_list \leftarrow cve\_list \cup \text{make}(cve\_builder)$ ;  
**end**  
**return**  $cve\_list$ ;

**Algorithm 3:** Extracting CVE entries

### 3.5 Visualising the Results

Using a pre-existing template, the results from the e-mail analysis will be presented in a temporary webpage, which can then be saved independently. Other than the referenced JavaScript libraries, the document requires no additional information or database access, allowing it to be quickly shared.

## **4 Evaluation**

### **4.1 Methodology**

### **4.2 Results**

## **5 Conclusions and Future Work**

### **5.1 Conclusions**

### **5.2 Future Work**

# Bibliography

- [1] Rakesh Agrawal et al. “Order preserving encryption for numeric data”. In: *Proceedings of the 2004 ACM SIGMOD international conference on Management of data*. ACM. 2004, pp. 563–574.
- [2] D. Crocker. *STANDARD FOR THE FORMAT OF ARPA INTERNET TEXT MESSAGES*. STD 11. RFC Editor, Aug. 1982.
- [3] Jerry Felix and Chris Hauck. “System security: a hacker’s perspective”. In: *Interex Proceedings* 1 (1987), pp. 6–6.
- [4] PHP Group et al. *PHP: Character classes - Manual*. URL: <https://secure.php.net/manual/en/regexp.reference.character-classes.php>.
- [5] Danesh Irani et al. “Modeling unintended personal-information leakage from multiple on-line social networks”. In: *Internet Computing, IEEE* 15.3 (2011), pp. 13–19.
- [6] Akanksha Joshi, Ravendar Lal, and Tim Finin. “Extracting cybersecurity related linked data from text”. In: *Semantic Computing (ICSC), 2013 IEEE Seventh International Conference on*. IEEE. 2013, pp. 252–259.
- [7] Jason RC Nurse et al. “Investigating the leakage of sensitive personal and organisational information in email headers”. In: *Journal of Internet Services and Information Security (JISIS)* 5.1 (2015), pp. 70–84.
- [8] Panagiotis Papadimitriou and Hector Garcia-Molina. “Data leakage detection”. In: *Knowledge and Data Engineering, IEEE Transactions on* 23.1 (2011), pp. 51–63.
- [9] Anna Squicciarini, Smitha Sundareswaran, and Dan Lin. “Preventing information leakage from indexing in the cloud”. In: *Cloud Computing (CLOUD), 2010 IEEE 3rd International Conference on*. IEEE. 2010, pp. 188–195.
- [10] Brad Templeton. *Reaction to the DEC Spam of 1978*. URL: <http://www.templetons.com/brad/spamreact.html>.
- [11] Marwan Al-zarouni. *Tracing E-mail Headers*. 2004.