

JA3cury - A new approach to TLS fingerprinting by merging fingerprinting methods

Lukáš Hejčman¹, Ing. Karel Hynek², Ing. Tomáš Čejka Ph.D.³

Abstract

TLS is the most popular encryption protocol used on the internet today. It aims to provide high levels of security and privacy for inter-device communication. However, it presents a challenge from a network monitoring and administration standpoint, as it is not possible to analyse the communication encrypted with TLS at a large scale with existing methods based on deep packet inspection. Analysing encrypted communication can help administrators to detect malicious activity on their networks, and can help them identify potential security threats. In this paper, we present a method that allows us to leverage the advantages of two TLS fingerprinting methods, JA3 and Cisco Mercury, to determine the operating system and processes of clients on multiple networks. Our method is able to achieve comparable or better results than the existing Mercury approach for our datasets whilst providing more analysis opportunities than JA3. Furthermore, by using JA3 fingerprints, we open the door to the utilisation of this approach in the wider industry, where JA3 fingerprinting is predominant.

Keywords: TLS – Fingerprint – Cisco – Mercury – JA3 – Identification – JA3cury

Supplementary Material: N/A

¹xhejcm01@stud.fit.vutbr.cz, Faculty of Information Technology, Brno University of Technology

²hynekkar@fit.cvut.cz, Faculty of Information Technology, Czech University of Technology

³cejkat@cesnet.cz, CESNET, z.s.p.o.

1. Introduction

Network traffic analysis is the process of capturing and analysing network traffic to increase the network's performance and security. Whilst many methods exist to accomplish this, it is becoming more important to utilize automatic techniques to filter the content on the network into specific categories due to the constantly increasing amount of traffic passing through networks.

However, network traffic analysis is becoming more challenging due to the rise in the use of encrypted traffic. According to a report by Google LLC[1], the amount of encrypted network traffic using the TLS protocol has been steadily increasing since at least 2014 to the current 95% of all traffic on the internet. Encryption is generally perceived as beneficial towards the privacy and security of the communication between endpoints on the internet; however, it is making the

traditional network analysis approach based on packet content inspection useless. Also, the recent report published by ENISA¹ recognizes encrypted traffic as a possible serious security threat due to hidden malicious activities[2] that current monitoring tools cannot easily detect. Therefore, it is essential to focus current research activities on encrypted traffic analysis and the retrieval of information about the connections and communicating systems.

The traditional approach for encrypted traffic analysis is its decryption (e.g., using a network proxy). However, this is very computationally expensive and is therefore not feasible for high throughput networks at a large scale. Furthermore, decrypting user communication can be seen as a security transgression and a

¹European Union Agency for Cybersecurity

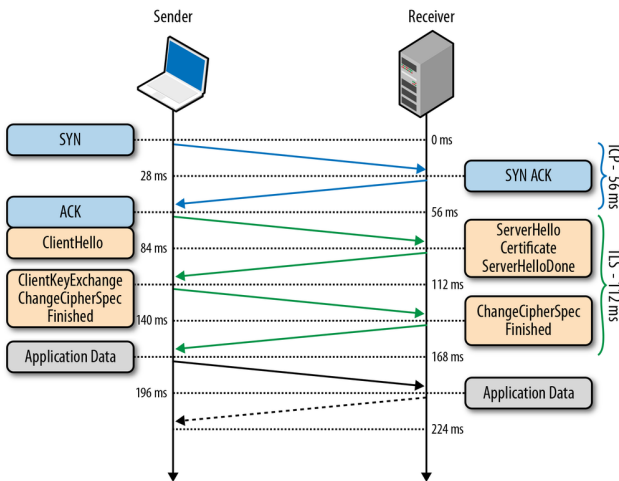


Figure 1. TLS Handshake overview

33 privacy concern. Traffic decryption is usually deployed
 34 in a highly restricted environment, such as company
 35 networks.

36 One of the methods that sufficiently preserves the
 37 privacy and security benefits of encryption is TLS fin-
 38 gerprinting. It works by gathering information about
 39 the client from the unencrypted portion of the TLS
 40 communication — the *Client Hello* packet. This packet
 41 outlines the parameters of the TLS communication sup-
 42 ported by the client. Because different combinations
 43 of applications and operating systems support particu-
 44 lar subsets of TLS communication parameters, we are
 45 able to create a database of applications and their TLS
 46 parameters. The database then allows us to evaluate
 47 the encrypted portion of the communication and infer
 48 the application name or the operating system type.

49 This information about the user application and
 50 operating systems is beneficial for network adminis-
 51 trators and security experts. It allows the detection of
 52 network policy violations, malware presence, or just
 53 insecure and outdated versions of operating systems
 54 or applications.

55 One of the most common fingerprinting approaches
 56 is the JA3 fingerprint introduced by Salesforce[3]. JA3
 57 is the de-facto industry standard, supported in high-
 58 performance monitoring tools (such as Flowmon ADS[4])
 59 and IDSs (such as Suricata[5]). Even though JA3 is
 60 very popular, it lacks a well-maintained and curated
 61 public fingerprint database. The existing open-source
 62 databases usually lack information, are unmaintained,
 63 and sometimes even do not follow the standard JA3
 64 fingerprint format.

65 Anderson et al.[6] from Cisco research presented
 66 a more accurate fingerprinting approach called Mer-
 67 cury [7]. The mercury fingerprint approach allows a
 68 much more accurate client taxonomy due to the vastly
 69 larger amounts of information stored in the database.

70 Even though the Mercury database includes a well
 71 maintained public database and it is in active devel-
 72 opment, its adoption across the industry is poor. One
 73 of the reasons might be the long plain-text fingerprint,
 74 which put even more strain on network analysis tools.

75 Therefore, we propose a new approach to TLS fin-
 76 gerprinting called JA3Cury, which combines the JA3
 77 fingerprint with other context information from the
 78 Mercury database. Surprisingly, even though it uses a
 79 smaller fingerprint, the JA3Cury matching algorithm
 80 achieved larger accuracy than the original database
 81 based on Mercury fingerprints in our testing environ-
 82 ment. Furthermore, we are able to introduce this larger
 83 precision to systems that currently use the JA3 ap-
 84 proach without redesigning or reengineering these sys-
 85 tems.

86 Our paper is organized as follows: section 2 briefly
 87 summarizes the TLS protocol and the concept of TLS
 88 fingerprinting. Section 3 describes the JA3 and Mer-
 89 cury fingerprinting methods, section 4 introduces the
 90 JA3Cury approach, section 5 provide comparisons of
 91 results obtained with JA3Cury and Mercury, and finally
 92 section 6 concludes the paper.

2. TLS Fingerprinting

93 Transport Layer Security (TLS) is a cryptographic pro-
 94 tocol designed to facilitate secure and encrypted com-
 95 munication between two parties. It is based on the now
 96 deprecated Secure Socket Layer (SSL) protocol. TLS
 97 is the de-facto standard for secure communication on
 98 the internet, as it is the protocol used by HTTP Secure
 99 (HTTPS).
 100

101 Before two clients can communicate through a
 102 TLS secured connection, they must first agree on the
 103 parameters of the connection, such as the ciphers sup-
 104 ported by both sides. This negotiation happens during
 105 the “handshake” phase of the communication. An
 106 overview of the handshake can be seen in figure 1.

107 The *Client Hello* and *Server Hello* messages are
 108 always sent without encryption, because the param-
 109 eters of the communication haven’t yet been agreed
 110 upon. This gives us the opportunity to intercept these
 111 messages and analyse their contents.

112 Both fingerprinting methods work by selecting a
 113 subset of data from the Client Hello packet of the TLS
 114 communication, compiling them into some format,
 115 and comparing them with a database of collected and
 116 annotated fingerprints.

117 3. Existing solutions

118 JA3

119 The JA3 fingerprinting method works by only extract-
120 ing information from the following five fields of the
121 Client Hello packet:

- 122 • TLS version
- 123 • Supported cipher suites
- 124 • TLS extension headers
- 125 • Elliptic curves
- 126 • Elliptic curves point formats

127 The fingerprint is then generated by concatenating
128 these fields in their decimal representation into a string
129 separated by comas, to generate the following:

130 Version,Ciphers,Extensions,EC,ECPF

131 The TLS standard does not require all these field
132 to be present in the Client Hello packet[8]. If any
133 field is missing, it is replaced with an empty string in
134 the fingerprint representation. The fingerprint is then
135 hashed with MD5.

136 For example, the following string is a valid JA3
137 fingerprint in decimal format:

```
769,4-5-47-51-50-10-22-19-9-  
21-18-3-8-20-17-255,,,
```

138 This string would then be hashed with MD5 to
139 generate the string

```
b677934e592ece9e09805bf36cd68d8a
```

140 which would then be used as the primary key in
141 the fingerprint database.

142 The JA3 databases usually contain 3 fields: the
143 JA3 string, the JA3 hash, and the description. However,
144 the description of the fingerprint which is then used
145 for identification is not standardized; this results in
146 many fingerprints containing information about the
147 application in a format which is completely unsuitable
148 for further classification and analysis. For example,
149 the string

```
BurpSuite Free  
(Tested: 1.7.03 on Windows 10),  
eclipse,JavaApplicationStub,idea
```

150 contains information about the application, the oper-
151 ating system, and some further processes without
152 conforming to a specified format, and thus makes it
153 impossible to parse in large quantities. This results in
154 databases where the fingerprint classification must be
155 taken at face value and no further analysis is possible.

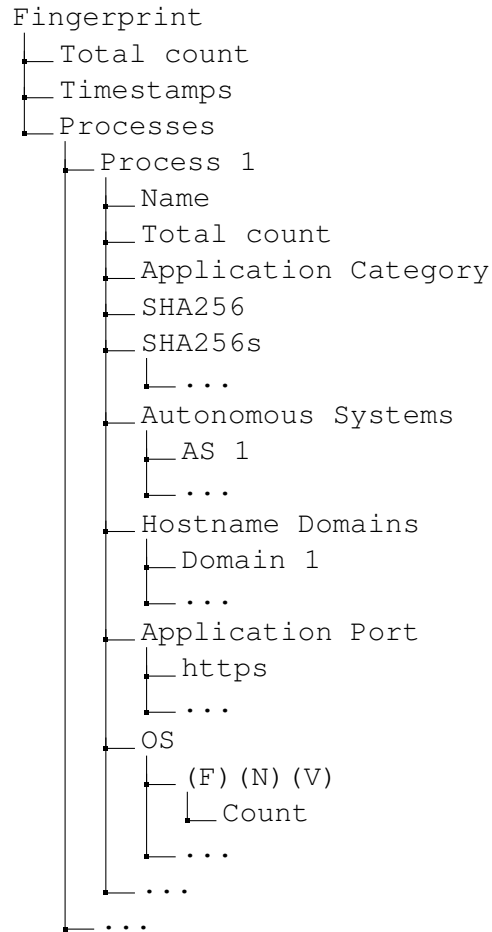


Figure 2. The structure of an entry in the Mercury database.

Mercury 156

The Mercury fingerprint format was developed by
Cisco by David McGrew and Blake Anderson. The
fingerprint itself contains much more information, as it
is basically a string representation of important fields
in the Client Hello packet in hexadecimal format.

The overall format of the fingerprint is the follow-
ing:

```
(version) (cipher suites) 164  
((extensions)...) 165
```

Where (version) is the hex representation of
the advertised TLS version, (cipher suites) is a
list of hex values of cipher suited offered by the client,
and ((extensions)...) contains the hex repre-
sentation of the extensions and their values (where
applicable).

Furthermore, compared to JA3, the Cisco Mer-
cury database contains much more information, and
is formatted so that it encourages further analysis of
the results. This format is visualized in figure 2. The
database was created by a novel approach of fusing net-
work endpoint data and captured fingerprints[9], so it

178 contains detailed process and contextual information.

179 The approach of Mercury prefers generating a
180 knowledge base about a network based on a few clients,
181 which is then used for detection[10]. However, the
182 Mercury GitHub repository also includes a well main-
183 tained open source database which can be used for
184 identification without the need to generate custom
185 knowledge bases.

186 The database fields in the Mercury database are
187 much more complex than in a JA3 databases. The
188 fields don't contain a simple 1 : 1 mapping of finger-
189 print to process, but instead contains many possible
190 processes for each fingerprint, including the number
191 of times they were encountered when building the
192 knowledge base. This number can then be used in
193 further classification. In the latest version of the Mer-
194 cury database which we used for our experiments, the
195 largest number of distinct processes mapped to a single
196 fingerprint was 9.

197 4. JA3cure

198 To leverage both the widespread usage of JA3 through-
199 out the industry and the better classification and anal-
200 ysis opportunities presented by the Mercury finger-
201 printing approach and database, we have devised an
202 approach we call JA3cure. With this approach, we
203 are able to search for JA3 fingerprints in the Mercury
204 database by converting existing Mercury fingerprints
205 to their corresponding JA3 representation.

206 This conversion works by extracting the relevant
207 information from the Mercury fingerprint, formatting it
208 as a JA3 fingerprint and hashing it using the MD5 hash
209 function. The result is a fingerprint database based on
210 the Cisco Mercury database that can be indexed using
211 JA3 fingerprints.

212 An example of this conversion can be seen in fig-
213 ure 3. As is shown in the figure, the conversion from
214 Mercury to JA3 is destructive; some information is lost
215 during the conversion. This means that the Mercury
216 database that previously contained only unique finger-
217 print entries now contains duplicate fingerprint entries
218 for some fingerprints. Out of the 9,060 entries in the
219 database we used, 1,914 unique fingerprints were lost;
220 this is a 21.1% decrease in fingerprint count.

221 Initially, we were worried this would lead to a de-
222 crease in accuracy of client and process identification
223 compared to the original database; however, as we will
224 show in Section 5, the accuracy in our experiments
225 remained the same or even increased for some certain
226 scenarios.

227 This decrease in the number of unique fingerprints
228 also influenced the approach of our classification algo-
229 rithms to fingerprint collisions. Kotzias et al. found
230 that around 7.3% of JA3 fingerprints cause collisions
231 with each other[11]. However, since our approach of
232 converting fingerprints introduces collisions into the
233 database, our classification algorithms were designed
234 to deal with them and they didn't cause a perceptible
235 decrease in detection accuracy.

236 Furthermore, it is important to note that the JA3cure
237 approach congregates the Client Hello classifications
238 over some time period, rather than identifying and
239 classifying a single Client Hello. This has led to
240 better classification results overall, as gathering the
241 information over many handshakes increases the op-
242 erating system detection accuracy, as well as overall
243 process detection accuracy due to the systems ability
244 to overcome statistical anomalies.

245 5. Detection

246 Thanks to the complexity of the Mercury/JA3cure
247 database, we were able to try many different finger-
248 printing approaches with varying degrees of accuracy.
249 Overall, we created 7 algorithms for traffic classifica-
250 tion. Each algorithm took into account different com-
251 binations of information from the database, as well as
252 some contextual information about the Client Hello
253 packet, such as the destination port, server domain
254 name, etc.

255 During detection, we compared three sets of re-
256 sults:

- 257 • As a baseline measurement, we used an unmod-
258 ified Mercury classification created with the of-
259 ficial `pmercury` utility.
- 260 • Our classification algorithms performed using
261 the unmodified version of the database using
262 Mercury fingerprints.

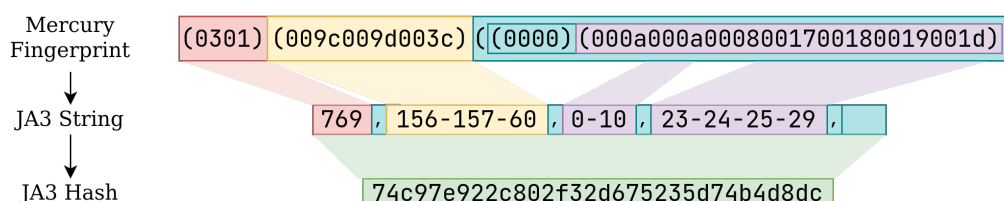


Figure 3. Converting a Mercury fingerprint to JA3

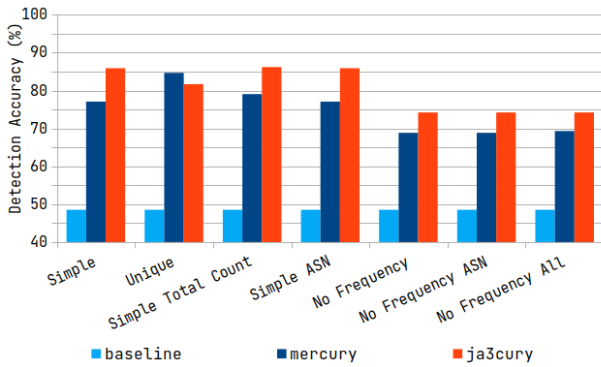


Figure 4. Process classification results.

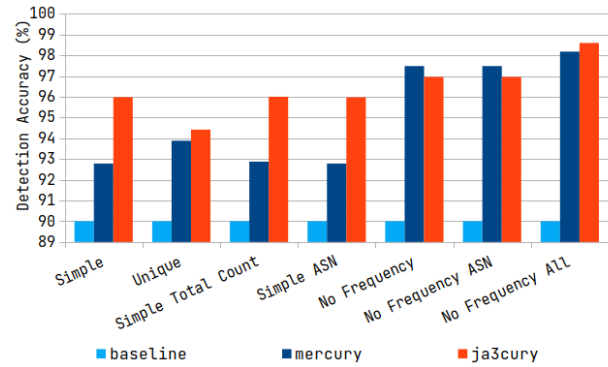


Figure 5. Category classification results.

- Our classification algorithms performed using the JA3database.

Our datasets contain over 48,000 Client Hello packets collected over 6 home networks made up of personal computers, laptops, home servers, and phones. All the major operating systems (Windows, Mac, Linux, Android, iOS) were represented in our experiments.

Some of our classification algorithms even exposed problems with the current version of the Mercury database; it was created on corporate Cisco networks, and thus skews heavily towards enterprise applications, such as Cisco Webex, and towards very specific operating systems, mainly Mac OS. However, we discovered that JA3cury is largely able to overcome this due to the meshing together of different fingerprints from the original Mercury database, which tends to average out the discrepancies.

Process and Category Detection

Each process in the database contains a classification into many categories, such as productivity, security, or gaming. This means that process and category detection are closely connected together. However, we discovered it is possible to obtain a high accuracy of category detection even with relatively low process detection accuracy. This is due to the fact that the erroneous classifications tend to get averaged out due to the vastly lower number of categories than processes, which leads to larger detection scores.

The average results for process classification of the top 5 processes for each client can be seen in figure 4.

Our JA3cury method was able to outperform the baseline results generated by the official `pmercury` for all our classification algorithms. Furthermore, our modified JA3cury database outperformed the original Mercury database in all but one experiment.

Furthermore, the category detection using our algorithms was also successful, viz figure 5. Again, our classification using JA3cury was more successful overall than either the original `pmercury` detection, or

even the detection using our algorithms and the original Mercury database.

The difference in scoring processes and categories between the default Mercury and our JA3cury approaches is well illustrated in figure 6. This graph shows the scores of different processes on the Y axis, with each process being represented by a bar on the X axis. Each process is also scored with Mercury and JA3cury. You can see that JA3cury identified more processes and categories, and it attributed a higher score to correct processes compared to Mercury.

Operating System Detection

The information about operating system classification in the Mercury database is dependent on the process classification, as the operating system information is nested inside each process (see figure 2). Furthermore, the database unfortunately does not contain information about mobile operating systems; instead, they tend to be classified as desktop operating systems with the most similar kernel architecture; MacOS for iOS devices, and Linux for Android devices.

The database contains operating system information split into three parts: the family (Linux, MacOS, Windows), the name (Windows 10 Professional, Linux 4.19, ...), and the build version (10.5.6.7, ...). For our experiment, we decided to classify the operating system using a tree structure with depth of 4, where the operating system frequency trickles down into the leaf nodes. An example of this tree can be seen in figure 8.

Furthermore, the tree is sorted such that each parent has its children ordered from the most frequent to the least frequent. This allows us to find the most probable operating system by taking the leftmost nodes. In this case, the classification would result in WinNT - Windows 10 Enterprise - 10.0.18363.

The operating system classification was performed on all clients in each network. The comparison of result for detection of the operating system family using our classifiers can be seen in figure 7. The figure

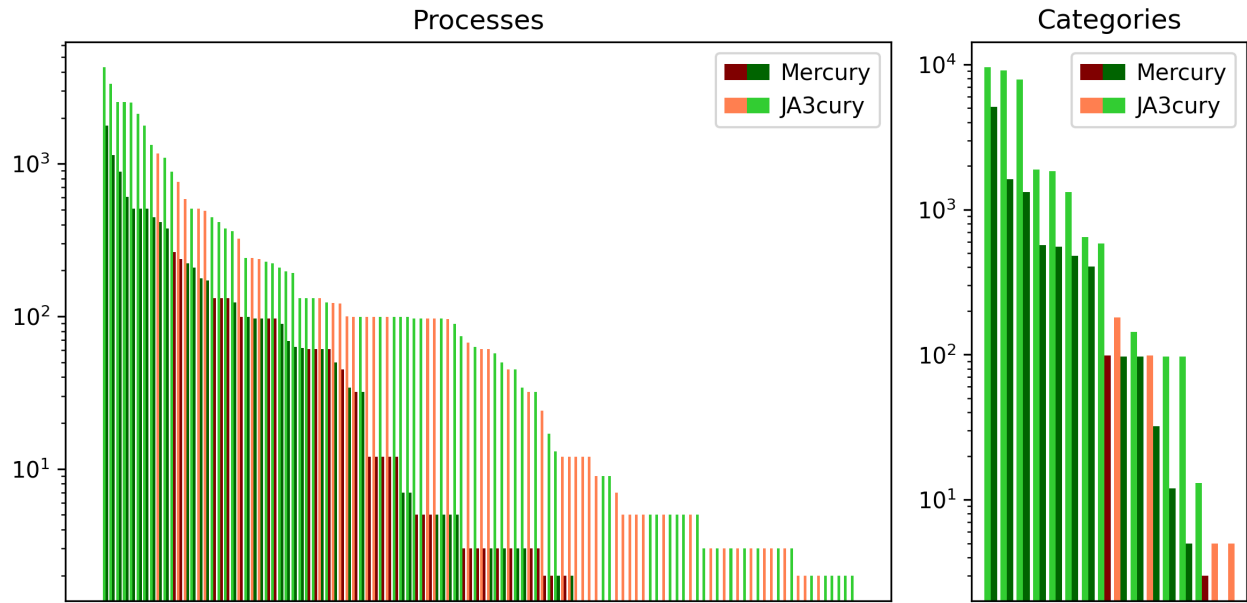


Figure 6. Detailed look at JA3curey and Mercury classification scores for one client.

341 contains only results created with our classification
 342 algorithms, because `pmercury` doesn't return infor-
 343 mation about the operating system.

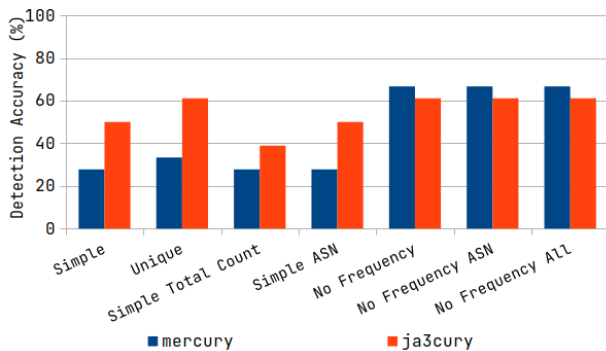


Figure 7. Operating system classification results.

344 JA3curey was able to detect the operating system of
 345 a client more accurately overall. However, the operat-
 346 ing system detection was less reliable on our datasets
 347 compared to process detection due to the prevalence
 348 of Linux machines in our datasets. The database, how-
 349 ever, contains many more entries for Windows and
 350 Mac OS, than it does for Linux. Furthermore, the
 351 fact that the database doesn't contain mobile operating
 352 systems leads to a lower accuracy as well.

353 6. Conclusion

354 In conclusion, we have developed an alternative TLS
 355 fingerprinting approach based on the strengths of JA3
 356 and Mercury fingerprinting. This approach allows us
 357 to utilize the more content rich Mercury database in a
 358 setting where we would have to rely on the results of
 359 JA3 without any further analysis. Furthermore, our ap-

proach is compatible with most existing fingerprinting
 modules thanks to the wide adoption of JA3, and can
 be used with existing modules and infrastructure with-
 out the need of re-engineering or redesigning these
 systems. Because our approach takes into account
 Client Hello packets over a time period, and doesn't
 classify each packet separately, our approach is able
 to overcome the disadvantage of missing information
 in the JA3 fingerprint compared to the full Mercury
 fingerprint.

Our method has proven to be at least as accurate
 as default Mercury fingerprinting. Furthermore, when
 used outside of corporate networks, it tends to be more
 accurate. Furthermore, the process category was de-
 tected correctly for all major categories.

The major area of further development could in-
 clude increasing the accuracy of operating system de-
 tection and the addition of mobile operating systems
 into the Mercury database.

Acknowledgments

We would like to thank Blake Anderson and David
 McGrew for their comments and feedback.

References

- [1] Google Transparency Report. HTTPS encryption on the web – Google Transparency Report, 2021.
- [2] ENISA. Encrypted Traffic Analysis, April 2020.
- [3] John Althouse, Jeff Atkinson, and Josh Atkins. `salesforce/ja3`, June 2017. original-date: 2017-06-13T22:54:10Z.

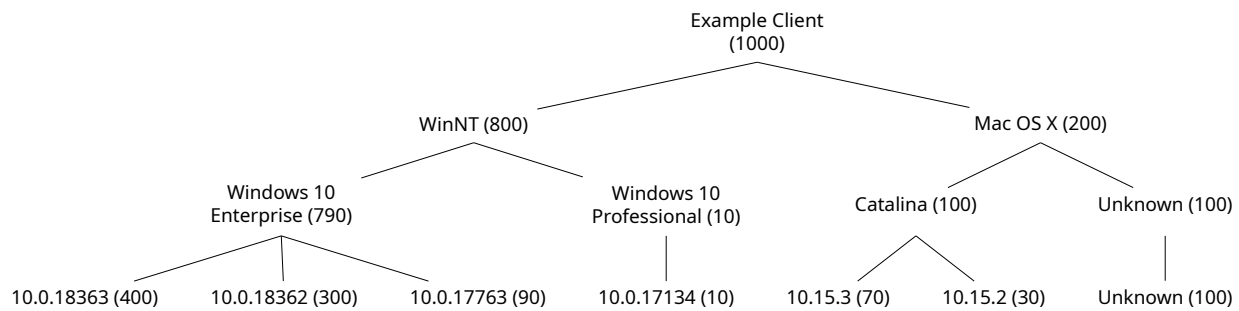


Figure 8. OS Classification Tree

- 389 [4] Flowmon. Encrypted Traffic Analysis.
- 390 [5] Suricata. 6.17. JA3 Keywords — Suricata 6.0.1
391 documentation, 2019.
- 392 [6] Blake Anderson and David McGrew. Accurate
393 TLS Fingerprinting using Destination Context
394 and Knowledge Bases. *arXiv:2009.01939 [cs]*,
395 September 2020. arXiv: 2009.01939.
- 396 [7] Blake Anderson, David McGrew, Brandon En-
397 right, Lucas Messenger, Adam Weller, An-
398 drew Chi, and Shekhar Acharya. *cisco/mer-
399 cury*, August 2021. original-date: 2019-08-
400 30T21:58:25Z.
- 401 [8] Eric Rescorla. RFC 8446 - The Transport Layer
402 Security (TLS) Protocol Version 1.3. Technical
403 report, August 2018.
- 404 [9] Blake Anderson, David McGrew, and Keith
405 Schomburg. The generation and use of tls finger-
406 prints, Jan 2019.
- 407 [10] Blake Anderson and David McGrew. Video cor-
408 respondece in regards to cisco cognitive intelli-
409 gence and cesnet collaboration., Mar 2021.
- 410 [11] Platon Kotzias, Abbas Razaghpanah, Johanna
411 Amann, Kenneth G. Paterson, Narseo Vallina-
412 Rodriguez, and Juan Caballero. Coming of age:
413 A longitudinal study of tls deployment. In *Pro-
414 ceedings of the Internet Measurement Confer-
415 ence 2018*, IMC '18, page 415–428. Association
416 for Computing Machinery, Oct 2018.