

# WAF-A-MoLE: An adversarial tool for assessing ML-based WAFs

Andrea Valenza, Luca Demetrio, Gabriele Costa, Giovanni Lagorio

*University of Genova*

*IMT Lucca*

---

## Abstract

Web Application Firewalls (WAFs) are plug-and-play security gateways. They promise to enhance the security of a (potentially vulnerable) system with minimal cost and configuration. In recent years, machine learning-based WAFs are catching up with traditional, signature-based ones. They are competitive because they do not require predefined rules. Instead, they infer their rules through a learning process.

In this paper we present WAF-A-MoLE, a WAF breaching tool. It uses guided mutational-based fuzzing to generate adversarial examples. The main applications include WAF (i) penetration testing, (ii) benchmarking and (iii) hardening.

*Keywords:* web application firewall, SQL injection, penetration testing, adversarial machine learning

---

## 1. Motivation and significance

Many modern systems expose some web services over the Internet. When they are vulnerable, the security of the entire system is compromised. A widespread mitigation technique is to deploy a *Web Application Firewall* (WAF). A WAF attempts to detect malicious incoming payloads and drop them before they reach their target. Clearly, the ability to craft payloads that pass undetected gives a tremendous advantage to attackers.

WAFs are traditionally signature-based, with a predefined set of rules for attack identification. However, this approach lacks generality and requires a significant effort to maintain the rule set. For this reason, researchers have recently considered the adoption of *machine learning* (ML). ML-based WAFs overcome some of the limitations of traditional WAFs. Their detection rules are inferred from a set of payloads through a training process.

14 On the flip side, an aware attacker can take advantage of biases in the  
15 training set. For instance, the training set might miss some relevant pay-  
16 loads, so causing blind spots in the classification space. *Adversarial machine*  
17 *learning* [1] studies this phenomenon, i.e., how to craft misleading payloads,  
18 a.k.a. *adversarial examples*.

19 In this paper we present WAF-A-MoLE, a tool to generate adversarial  
20 examples for ML-based WAFs. In particular, the current version of the  
21 tool focuses on SQL injection (SQLi) attacks. WAF-A-MoLE starts from a  
22 payload and mutates it to bypass a target WAF. The tool relies on a set of  
23 semantics preserving mutation operators. The mutation process is guided by  
24 the classification confidence of the target WAF.

25 This paper is structured as follows. Section 2 describes WAF-A-MoLE  
26 and its main functionalities. Section 3 shows an example of how WAF-A-  
27 MoLE bypasses a target toy WAF. Section 4 highlights the impact of WAF-  
28 A-MoLE on the security community. Finally, Section 6 concludes the paper.

## 29 2. Software description

30 WAF-A-MoLE uses a *mutation-based fuzz testing* [2] methodology to cre-  
31 ate attacks that bypass a target WAF. More precisely, WAF-A-MoLE uses  
32 the classification score of the WAF to guide the fuzzing process by prioritiz-  
33 ing the most promising payloads. We describe the overall architecture and  
34 main functionalities of WAF-A-MoLE in the next section.

### 35 2.1. Software Architecture

36 WAF-A-MoLE is both a library and a command line tool (obtained by  
37 means of *Click*<sup>1</sup> decorators on the main exported functions) implemented  
38 in Python 3. Figure 1 shows the main workflow of WAF-A-MoLE. Briefly,  
39 the orchestrator (not shown in the figures) takes an initial payload  $p_0$ , that  
40 the target WAF detects as malicious with a confidence score  $\sigma_0 \in [0, 1]$ , and  
41 inserts it in the initially empty payload `Pool`. The `Pool`, in turn, manages a  
42 priority queue, storing payloads in decreasing ordered of their scores.

43 During each iteration, the head of the queue  $p_n$  is picked from the `Pool`,  
44 and passed to the `Fuzzer`, which randomly mutates  $p_n$  into  $p_{n+1}$  by apply-  
45 ing some mutation operators (see Section 2.2). Then,  $p_{n+1}$  is submitted to  
46 the target WAF for classification. Since we do not expect WAFs to adhere  
47 to any specific interface, WAF-A-MoLE uses specific *adapters* that ensure

---

<sup>1</sup><https://click.palletsprojects.com/en/7.x/>

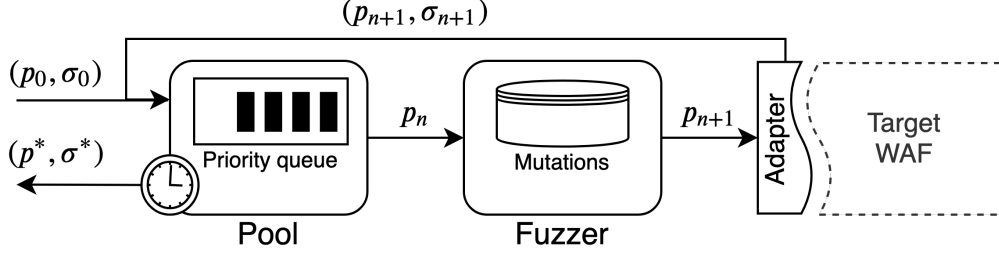


Figure 1: Main workflow of WAF-A-MoLE.

48 compatibility. The **Adapter** then returns the classification score  $\sigma_{n+1}$  of  $p_{n+1}$ ,  
49 which is fed back into the **Pool**.

50 This cycle finishes successfully whenever the best confidence score  $\sigma^*$  is  
51 less than a given threshold, or is interrupted, returning the best pair  $(p^*, \sigma^*)$   
52 found so far, because the number of iterations, queue sizes or computation  
53 time reach their maximum values.

54 In order to apply WAF-A-MoLE to different machine learning models,  
55 without incurring into a tight coupling, we designed an interface, modeled  
56 in Python as an abstract class called `Model`, which generalizes the behaviour  
57 of those models. This class provides two abstract methods, `classify` and  
58 `extract_features`, that need to be instantiated for each kind of model. That  
59 is, since, no real model matches exactly our interface, for each of them we  
60 need an adapter class that wraps the target model and exports our `Model`  
61 interface (see Figure 1).

62 We provide many wrappers out of the box, which are the ones that we  
63 used for running our experiments. They also serve as examples of how to  
64 implement new wrappers. In particular, we offer wrappers for two well-known  
65 frameworks: `SklearnModelWrapper` for *scikit-learn*<sup>2</sup>, and `KerasModelWrapper`  
66 for *keras*<sup>3</sup>.

## 67 2.2. Software Functionalities

68 As discussed in Section 2.1, the main components of WAF-A-MoLE are  
69 **Pool** and **Fuzzer**. The former handles the priority queue and termination con-  
70 ditions. Although they raise some technical issues (e.g., due to the memory  
71 usage of large data structures), these aspects belong to the generic context  
72 of program optimization. Instead, **Fuzzer** requires more attention.

<sup>2</sup><https://scikit-learn.org/stable/index.html>

<sup>3</sup><https://keras.io/>

73 Following the mutational fuzzing approach, **Fuzzer** applies a number of  
 74 mutation operators. Mutation operators act on the string representation, and  
 75 they modify the syntax of a payload without altering its semantics. Since we  
 76 focus on SQLi, the currently implemented mutation operators work on SQL.  
 77 We describe them below.

78 *CS.* The *Case Swapping* operator randomly changes the capitalization of  
 79 keywords in a query (e.g., **Select** to **sELecT**).

80 *WS.* *Whitespace Substitution* leverages the equivalence between several  
 81 alternative separators (whitespaces) between query tokens. For instance,  
 82 alternative whitespaces include `\n` (line feed), `\r` (carriage return) and `\t`  
 83 (horizontal tab). Each whitespace can be replaced by an arbitrary, non-  
 84 empty sequence of whitespaces (e.g., `1 = 1` may become `1\n\t=\r 1`).

85 *CI.* The *Comment Injection* operator randomly adds an inline comment  
 86 (`/*...*/`) between two query tokens. As whitespaces, inline comments act  
 87 as token separators (e.g., modifying `1 = 1` to `1/**/= 1`).

88 *CR.* The *Comment Rewriting* operator randomly modifies the content of  
 89 a comment. This both applies to inline and trailing (`#` and `--`) comments  
 90 (e.g., `/*abc*/` may become `/*xy*/`).

91 *IE.* The *Integer Encoding* operator modifies the representation of nu-  
 92 merical constants. This includes alternative base representations, e.g., from  
 93 decimal to hexadecimal (e.g., `0x2` for `2`), as well as statement nesting (e.g.,  
 94 `(SELECT 42)` for `42`).

95 *OS.* Some operators can be replaced by others that behave in the same  
 96 way. For instance, `1 = 1` (equality check) is simulated by `1 LIKE 1` (pattern  
 97 matching).<sup>4</sup> We call this mutation *Operator Swapping*.

98 *LI.* A *Logical Invariant* operator modifies a boolean expression by adding  
 99 terms that preserve its semantics (e.g., `1 = 1` is equivalent to `1 = 1 AND True`).

100 In defining our mutation operators, we took inspiration from some mali-  
 101 cious payload samples such as those listed by Awesome WAF<sup>5</sup> and Payloads  
 102 All The Things<sup>6</sup>. All in all, our operators generalize the techniques for pro-  
 103 ducing payloads similar to those mentioned above.

### 104 3. Illustrative Example

105 In this section, we provide a demonstration of WAF-A-MoLE applied to a  
 106 toy WAF. The toy WAF assigns a score to a (non-empty) payload  $p$  through

<sup>4</sup>Notice that, in general, `LIKE` is not equivalent to `=`. However, the equivalence holds when restricting to specific domains, e.g., comparison between integer constants.

<sup>5</sup><https://github.com/0xInfection/Awesome-WAF>

<sup>6</sup><https://github.com/swisskyrepo/PayloadsAllTheThings>

Op.	Mutant	$\sigma$	Op.	Mutant	$\sigma$
WS	admin'\t OR\n 1=1#	0.75	CI	admin' OR 1=1/**/#	0.67
CR	admin' OR 1=1#abcde	0.63	IE	admin' OR 0x1=1#	0.75
OS	admin' OR 1 LIKE 1#	0.79	LI	admin' OR 1=1 AND 0<1#	0.68

Table 1: Mutants and classification scores for the toy WAF.

the following function:

$$\sigma(p) = \min \left\{ 1, \frac{3 \cdot S(p)}{T(p)} \right\}$$

where  $S(p)$  is the number of special characters ' , = and \_ (a single white space) and  $T(p)$  is the total number of characters in  $p$ . For instance

$$\sigma(\text{admin' OR 1=1\#}) = \min \left\{ 1, \frac{3 \cdot 4}{14} \right\} \approx 0.86 \quad \sigma(\text{admin}) = \min \left\{ 1, \frac{3 \cdot 0}{5} \right\} = 0$$

Assuming the acceptance threshold of the toy WAF to be  $1/2$  (that is  $p$  is rejected when  $\sigma(p) > 1/2$ ) the first payload above is rejected. Table 1 reports the  $\sigma$  values of some mutants obtained through the application of the operators of Section 2.2.

Let assume that WAF-A-MoLE generated the payloads of Table 1. They are ordered by their  $\sigma$  values and inserted in the payload pool. Then, the next mutation round starts from the payload with the lowest  $\sigma$ , i.e., the one generated by CR.

#### 4. Impact

The impact of our tool on penetration testing activities is straightforward. Penetration testers can use WAF-A-MoLE as an off-the-shelf utility to craft attacks. For instance, we used WAF-A-MoLE against 9 instances of WAFs taken from literature to assess its effectiveness. The results are promising, and we provide an excerpt in Figure 2. Briefly, WAF-A-MoLE rapidly decreases the confidence score of the considered WAFs. The plot on the left shows how confidence decreases with the number of applied mutations. The plot on the right shows how it decreases over time, with a logarithmic scale. Since our approach is inherently stochastic, we ran our tool multiple times and chose the best run (i.e., the one that reached the threshold in the fewest mutation rounds) for each classifier.

For our analysis, we considered different classifiers:

1. WAF-Brain<sup>7</sup>, a deep neural network trained on raw characters contain-

<sup>7</sup><https://github.com/BBVA/waf-brain>

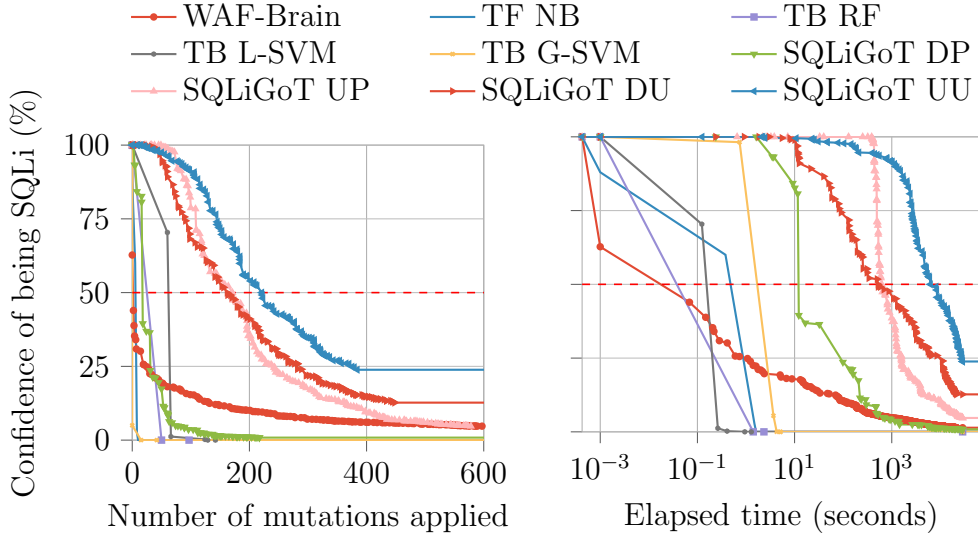


Figure 2: WAF-A-MoLE applied to the `admin' OR 1=1#` payload

- 133 ing legitimate and malicious payloads,
- 134 2. Token-based models [3, 4], implemented using different algorithms,
- 135 built on a histogram of tokens extracted from the queries, and
- 136 3. SQLiGoT [5], a Support Vector Machine (SVM) [6] classifier that rea-
- 137 sons on top of a graph structure extracted from input queries.

138 WAF-A-MoLE bypasses WAF-Brain in 7 mutation rounds. Some Token-  
 139 based approaches performed worse than WAF-Brain, Token-based Random  
 140 Forest and Gaussian SVM variants reached the threshold in respectively 2  
 141 and 3 mutation rounds. The Linear SVM and Naive Bayes variants performed  
 142 better, with 24 and 46 rounds. SQLiGoT proved to be the most resilient: one  
 143 of its variants, namely the Undirected Proportional one, reached the thresh-  
 144 old after 134 rounds, and the Undirected Unproportional version needed 290  
 145 steps.

146 The full details about our experiments can be found at <https://github.com/AvalZ/waf-a-mole>.

148 The by-product of WAF-A-MoLE are *adversarial examples* [7] for the  
 149 target WAF. Adversarial examples are a cornerstone in *adversarial training*  
 150 [7, 8, 9]. ML-based WAFs are trained on datasets that very rarely charac-  
 151 terize the entire classification domain. Developers can use the adversarial  
 152 examples to re-train their WAF, covering areas originally not included in the  
 153 training dataset [7, 8, 9]. Intuitively, training the classifier with regular data  
 154 and adversarial examples leads to a more robust model. **On the other hand,**  
 155 **an adversarially trained model loses accuracy w.r.t. its standard counterpart,**

156 as the problem to be learned is more complex. Although the methodology  
157 of [7, 8, 9] is not applied to our working domain, i.e., SQL payloads, we  
158 believe that a similar technique can be ported in our context. In this way,  
159 WAF-A-MoLE can support the WAF hardening process.

160 Using WAF-A-MoLE, we showed that ML-based WAFs are vulnerable  
161 to adversarial attacks. We believe that the main reason is the gap between  
162 the syntax level (of the WAFs classification) and the semantic level (of the  
163 vulnerable application). This observation pushes forward an open research  
164 question: to what extent (syntax-based) WAFs prevent injection attacks?  
165 WAF-A-MoLE candidates to be a valuable assessment tool to support this  
166 research line.

## 167 5. Related work

168 *Machine learning based WAFs.* Ceccato et al. [10] propose a clustering method  
169 for detecting SQL injection attacks against a victim service. The algorithm  
170 learns from the queries that are processed inside the web application un-  
171 der analysis using an unsupervised one-class learning approach, namely K-  
172 medoids [11]. New samples are compared to the closest medoid and flagged  
173 as malicious if their edit distance w.r.t. the chosen medoid is higher than  
174 the diameter of the cluster. Kar et al. [5] develop *SQLiGoT*, an SVM that  
175 express queries as graphs of tokens, whose edges represent the adjacency of  
176 SQL-tokens. Pinzon et al. [12] explore two different directions: visualiza-  
177 tion and detection, achieved by a multi-agent system called *idMAS-SQL*. To  
178 tackle the task of detecting SQL injection attacks, the authors set up two  
179 different classifiers, namely a Neural Network and an SVM. Makiou et al. [4]  
180 develop a hybrid approach that uses both machine learning techniques and  
181 pattern matching against a known dataset of attacks. The learning algorithm  
182 used for detecting injections is a Naive Bayes [13]. They look for different 45  
183 tokens inside the input query, chosen by domain experts. Similarly, Joshi et  
184 al. [3] use a Naive Bayes classifier that, given a SQL query as input, extracts  
185 syntactic tokens using spaces as separators. The algorithm produces a fea-  
186 ture vector that counts how many instances of a particular word occurs in  
187 the input query. The vocabulary of all the possible observable tokens is set  
188 a priori. Komiya et al. [14] propose a survey of different machine learning  
189 algorithms for SQL injection attack detection.

190 *Adversarial ML attacks.* Among all the techniques proposed in the state-of-  
191 the-art that leverage on white-box gradient techniques [7, 15, 16], we focus on  
192 black-box attacks, as they are similar to our method. Ilyas et al. [17] use the  
193 Natural Evolution Strategy (NES) [18] to guide the creation of adversarial

194 examples against well-known image classifiers. Xu et al.[19] propose a genetic  
195 algorithm that automatically learns which mutations should be applied to  
196 PDF malware to bypass a target classifier. Anderson et al. [20] train an  
197 agent to learn the best sequences of mutations applied to Windows malware  
198 to fool a target classifier. Chen et al. [21] estimate the target function’s  
199 boundary locally around a particular input. Then, they guide the generation  
200 of adversarial examples by computing an approximated gradient using the  
201 values obtained from the target classifier in that local region.

## 202 6. Conclusions

203 In this paper we presented WAF-A-MoLE, a guided mutational fuzzing  
204 tool to generate adversarial examples for ML-based WAFs. The tool has  
205 several possible applications, the main one being the security assessment of  
206 the WAFs.

207 There are numerous future directions for this research line. In particular,  
208 there are three that we consider of primary importance. In the first place,  
209 we plan to apply WAF-A-MoLE to commercial WAFs. The main difficulty  
210 is that vendors usually do not share details about the internals of their prod-  
211 ucts. Hence, this direction requires establishing an agreement with vendors.  
212 Secondly, we aim at extending our approach to deal with hybrid WAFs that  
213 also consider payload signatures. For both commercial and hybrid WAFs,  
214 we would also like to explore the possibility of integrating our tool in a re-  
215 training process for ML classifiers. Finally, we are interested in finding new  
216 mutation operators as well as investigating their effectiveness when applied  
217 alone or combined with others.

## 218 Acknowledgements

219 This work was partially funded by the Horizon 2020 project “Strategic  
220 Programs for Advanced Research and Technology in Europe” (SPARTA).

- 221 [1] L. Huang, A. D. Joseph, B. Nelson, B. I. Rubinstein, J. D. Tygar, Ad-  
222 versarial machine learning, in: Proceedings of the 4th ACM workshop  
223 on Security and artificial intelligence, ACM, 2011, pp. 43–58.
- 224 [2] A. Zeller, R. Gopinath, M. Böhme, G. Fraser, C. Holler, Mutation-  
225 based fuzzing, in: Generating Software Tests, Saarland University, 2019,  
226 retrieved 2019-05-21 19:57:59+02:00.  
227 URL <https://www.fuzzingbook.org/html/MutationFuzzer.html>



- 228 [3] A. Joshi, V. Geetha, Sql injection detection using machine learning, in:  
229 2014 International Conference on Control, Instrumentation, Communi-  
230 cation and Computational Technologies (ICCICCT), IEEE, 2014, pp.  
231 1111–1115.
- 232 [4] A. Makiou, Y. Begriche, A. Serhrouchni, Improving web application  
233 firewalls to detect advanced sql injection attacks, in: 2014 10th Interna-  
234 tional Conference on Information Assurance and Security, IEEE, 2014,  
235 pp. 35–40.
- 236 [5] D. Kar, S. Panigrahi, S. Sundararajan, Sqligot: Detecting sql injection  
237 attacks using graph of tokens and svm, *Computers & Security* 60 (2016)  
238 206–225.
- 239 [6] C. Cortes, V. Vapnik, Support-vector networks, *Machine learning* 20 (3)  
240 (1995) 273–297.
- 241 [7] I. Goodfellow, J. Shlens, C. Szegedy, Explaining and harnessing adver-  
242 sarial examples, in: *International Conference on Learning Representa-*  
243 *tions*, 2015.  
244 URL <http://arxiv.org/abs/1412.6572>
- 245 [8] K. Grosse, N. Papernot, P. Manoharan, M. Backes, P. McDaniel, Ad-  
246 versarial examples for malware detection, in: *European Symposium on*  
247 *Research in Computer Security*, Springer, 2017, pp. 62–79.
- 248 [9] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, A. Vladu, Towards  
249 deep learning models resistant to adversarial attacks, *Sixth International*  
250 *Conference on Learning Representations (ICLR)*.
- 251 [10] M. Ceccato, C. D. Nguyen, D. Appelt, L. C. Briand, Sofia: an auto-  
252 mated security oracle for black-box testing of sql-injection vulnerabili-  
253 ties, in: *Proceedings of the 31st IEEE/ACM International Conference*  
254 *on Automated Software Engineering*, ACM, 2016, pp. 167–177.
- 255 [11] L. K. P. J. RDUSSEEUN, Clustering by means of medoids.
- 256 [12] C. I. Pinzon, J. F. De Paz, A. Herrero, E. Corchado, J. Bajo, J. M.  
257 Corchado, idmas-sql: intrusion detection based on mas to detect and  
258 block sql injection through data mining, *Information Sciences* 231 (2013)  
259 15–31.
- 260 [13] M. E. Maron, Automatic indexing: an experimental inquiry, *Journal of*  
261 *the ACM (JACM)* 8 (3) (1961) 404–417.

- 262 [14] R. Komiya, I. Paik, M. Hisada, Classification of malicious web code by  
 263 machine learning, in: 2011 3rd International Conference on Awareness  
 264 Science and Technology (iCAST), IEEE, 2011, pp. 406–411.
- 265 [15] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, A. Swami,  
 266 The limitations of deep learning in adversarial settings, in: 2016 IEEE  
 267 European Symposium on Security and Privacy (EuroS&P), IEEE, 2016,  
 268 pp. 372–387.
- 269 [16] N. Carlini, D. Wagner, Towards evaluating the robustness of neural  
 270 networks, in: 2017 IEEE Symposium on Security and Privacy (SP),  
 271 IEEE, 2017, pp. 39–57.
- 272 [17] A. Ilyas, L. Engstrom, A. Athalye, J. Lin, Black-box adversarial attacks  
 273 with limited queries and information, in: Proceedings of the 35th Inter-  
 274 national Conference on Machine Learning, ICML 2018, 2018.  
 275 URL <https://arxiv.org/abs/1804.08598>
- 276 [18] D. Wierstra, T. Schaul, J. Peters, J. Schmidhuber, Natural evolution  
 277 strategies, in: 2008 IEEE Congress on Evolutionary Computation (IEEE  
 278 World Congress on Computational Intelligence), IEEE, 2008, pp. 3381–  
 279 3387.
- 280 [19] W. Xu, Y. Qi, D. Evans, Automatically evading classifiers, in: Proceed-  
 281 ings of the 2016 Network and Distributed Systems Symposium, 2016,  
 282 pp. 21–24.
- 283 [20] H. S. Anderson, A. Kharkar, B. Filar, P. Roth, Evading machine learning  
 284 malware detection, Black Hat.
- 285 [21] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, C.-J. Hsieh, Zoo: Zeroth order  
 286 optimization based black-box attacks to deep neural networks without  
 287 training substitute models, in: Proceedings of the 10th ACM Workshop  
 288 on Artificial Intelligence and Security, ACM, 2017, pp. 15–26.

289 **Current code version**

<b>Nr.</b>	<b>Code metadata description</b>	<b>Please fill in this column</b>
C1	Current code version	v1.0.0
C2	Permanent link to code/repository used for this code version	<a href="https://github.com/AvalZ/waf-a-mole">https://github.com/AvalZ/waf-a-mole</a>
C3	Legal Code License	MIT
C4	Code versioning system used	git
C5	Software code languages, tools, and services used	Python 3
C6	Compilation requirements, operating environments & dependencies	Click
C7	If available Link to developer documentation/manual	<a href="https://waf-a-mole.readthedocs.io/en/latest/">https://waf-a-mole.readthedocs.io/en/latest/</a>
C8	Support email for questions	andrea.valenza@dibris.unige.it, luca.demetrio@dibris.unige.it

Table 2: Code metadata