

WAF-A-MoLE

Evading Web Application Firewalls through Adversarial Machine Learning

Luca Demetrio, Andrea Valenza, Giovanni Lagorio

Università di Genova

{luca.demetrio, andrea.valenza}@dibris.unige.it

giovanni.lagorio@unige.it

Gabriele Costa

IMT School for Advanced Studies

gabriele.costa@imtlucca.it



@zangobot
@AvalZ_



Contributions

- WAF-A-MoLE: a tool to produce adversarial examples against WAFs
- State-of-the-art ML-based WAFs benchmark and bypass
- Dataset of benign and malicious SQL queries

SQL Injections (SQLi)

- Injects arbitrary behaviour
- OWASP Top Ten #1 since 2010
- Should be sanitized

```
SELECT * FROM users  
WHERE username='admin' OR 1=1#'  
and password = 'zenhack'
```

Web Application Firewalls

- Inspects incoming payloads to detect (and block) attacks
- Traditionally signature-based (e.g., white/black lists, regular expressions)
- WAFs treat the symptoms, not the disease

WAFs powered by Machine Learning

- Learns the *model* from the *syntax* of standard attacks (e.g., ZAP, sqlmap)
 - Uses *proximity* instead of exact rule matching
 - Assigns a *confidence* score to input payloads
-
- WAF-Brain¹, End-to-end recurrent neural network
 - Token-based detectors²
 - SQLiGoT³, graph-based algorithm

¹ BBVA/waf-brain: WAF-Brain - the clever and efficient Firewall for the Web (<https://github.com/BBVA/waf-brain>)

² “Classification of malicious web code by machine learning”, *Komiya, Ryohei and Paik, Incheon and Hisada, Masayuki*, iCAST 2011

“SQL Injection detection using machine learning”, *Joshi, Anamika and Geetha, V*, ICCICCT 2014

³ “SQLiGoT: Detecting SQL injection attacks using graph of tokens and SVM” *Kar, Debabrata and Panigrahi, Suvasini and Sundararajan, Srikanth*, Computers & Security, 2016

Training and Benchmark

		A	R	P
ModSecurity CSR	Paranoia 1/2	86.10%	86.10%	100%
	Paranoia 3/4	91.85%	91.85%	100%
	Paranoia 5	96.46%	96.46%	100%
WAF-Brain	RNN	98.27%	96.73	99.8%
Token-based	Naive Bayes	50.16%	98.71%	50.08%
	Random forest	98.33%	98.33%	100%
	Linear SVM	98.75%	98.76%	100%
	Gaussian SVM	97.82%	97.82%	100%
SQLiGoT	Dir. Prop.	90.61%	97.30%	85.82%
	Undir. Prop.	96.38%	97.31%	95.54%
	Dir. Unprop.	90.52%	97.12%	85.80%
	Undir. Unprop.	96.25%	97.05%	95.53%

Are ML WAFs effective?

Trained on our dataset

Performances compared to
signature based WAF
(ModSecurity CSR)

Mangling SQLi Syntax

`admin' OR 1=1#` \equiv `admin' OR 0X1=1 or 0x726!=0x726 OR
0x1Dd not IN/*(seleCt 0X0)>c^Bj
>N]*/ ((SeLeCT 476) ,(SELECT(SElEct
477)) ,0X1de) oR 8308 noT lIkE
8308\x0c AnD truE OR 'FZ6/q' LiKE
'fz6/qI' and TRUE and
'>U' != '>uz' #t '%' 03;Nd`

Same attack, different syntax

Adversarial Machine Learning

- Introduces the presence of an **adversary**
- Looks for blind spots in **corner-cases**
- **Evasion attack:** malicious payloads seen as legit
- **Constraint:** preserves attack semantics

$$x^* = \underset{x,}{\operatorname{argmin}} D(f(x), c_t)$$

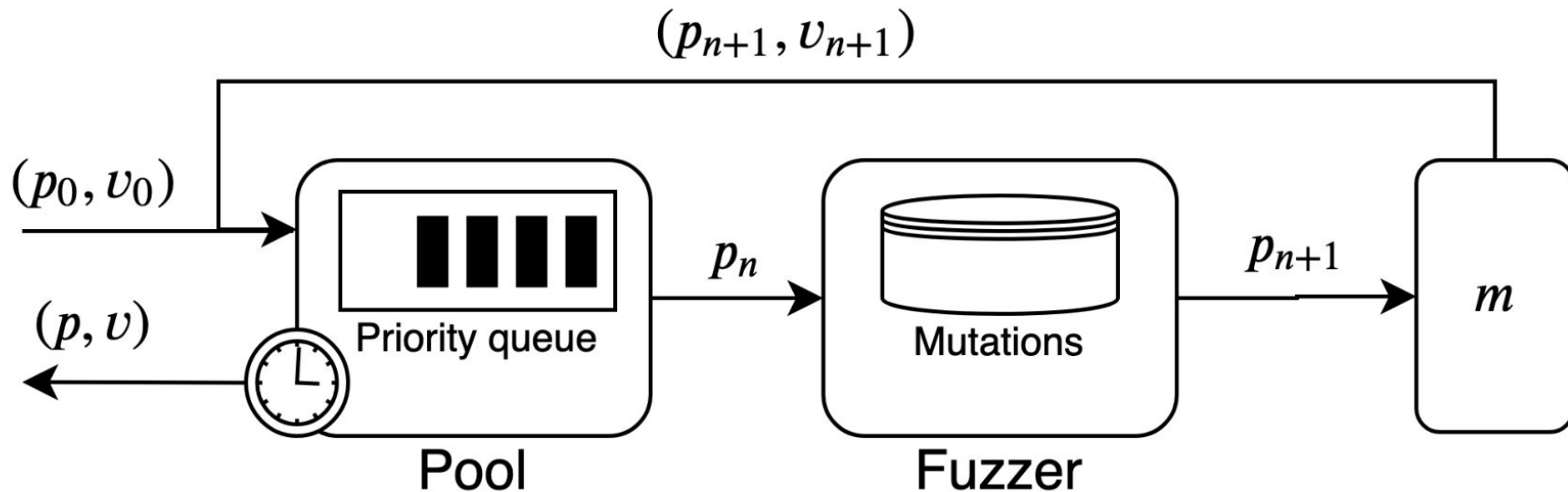
subject to $C(x)$

WAF-A-MoLE

Gray-box: only uses the confidence score assigned by the WAF

Efficient: prioritizes most promising (i.e., fittest) candidates

Plug-and-play: can be used against every SQLi detector (using the proper driver)



Mutational Fuzzing

input: Model m , Payload p_0 , Threshold t

output: $\text{head}(Q)$

$Q := \text{create_priority_queue}()$

$v := \text{classify}(m, p_0)$

enqueue(Q, p_0, v)

while $v > t$

$p := \text{mutate}(\text{head}(Q))$

$v := \text{classify}(m, p)$

enqueue(Q, p, v)

Randomly *mutates* the payload

Selects the best variants

Favors the survival of fittest
candidates with a priority queue

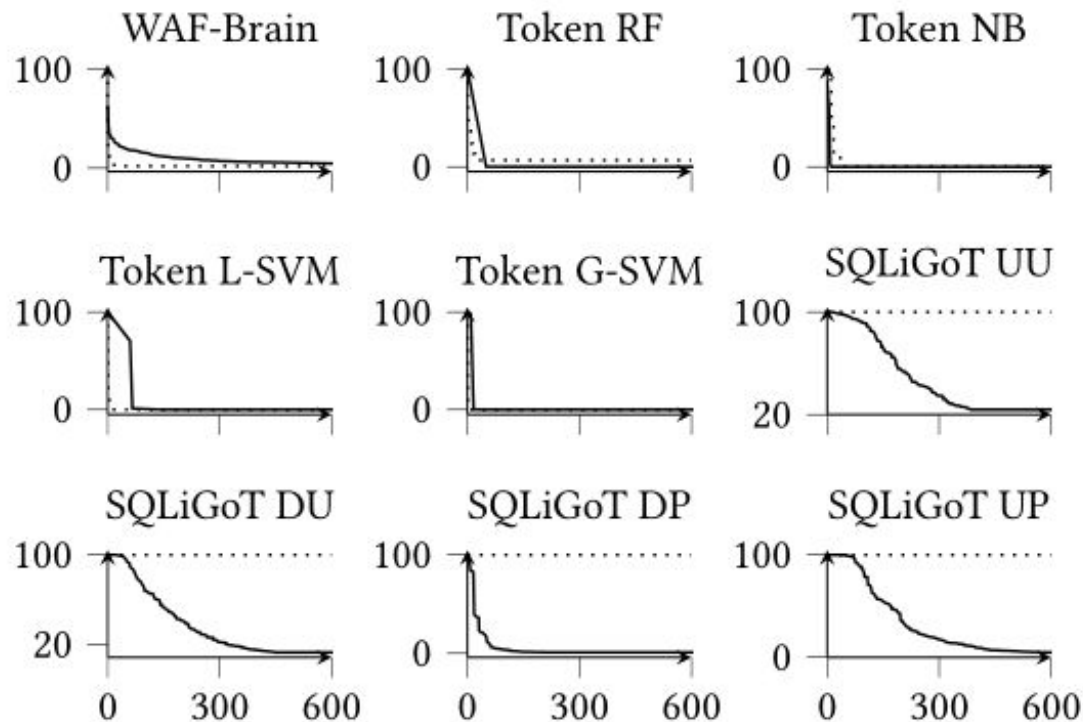
Mutation Operators and Example

Operator	Short definition	Example
Case Swapping	$CS(\dots a \dots B \dots) \rightarrow \dots A \dots b \dots$	$CS(\text{admin}' \text{ OR } 1=1\#) \rightarrow \text{ADmIn}' \text{ oR } 1=1\#$
Whitespace Substitution	$WS(\dots k_1 k_2 \dots) \rightarrow \dots k_1 _ k_2 \dots$	$WS(\text{admin}' \text{ OR } 1=1\#) \rightarrow \text{admin}' \backslash \text{n OR } \backslash \text{t } 1=1\#$
Comment Injection	$CI(\dots k_1 k_2 \dots) \rightarrow \dots k_1 /**/ k_2 \dots$	$CI(\text{admin}' \text{ OR } 1=1\#) \rightarrow \text{admin}' /**/ \text{OR } 1=1\#$
Comment Rewriting	$CR(\dots /s_0*/\dots \#s_1) \rightarrow \dots /s'_0*/\dots \#s'_1$	$CR(\text{admin}' /**/ \text{OR } 1=1\#) \rightarrow \text{admin}' /*abc*/ \text{OR } 1=1\# \text{xyz}$
Integer Encoding	$IE(\dots n \dots) \rightarrow \dots 0x[n]_{16} \dots$	$IE(\text{admin}' \text{ OR } 1=1\#) \rightarrow \text{admin}' \text{ OR } 0x1=1\#$
Operator Swapping	$OS(\dots \oplus \dots) \rightarrow \dots \boxplus \dots$ (with $\oplus \equiv \boxplus$)	$OS(\text{admin}' \text{ OR } 1=1\#) \rightarrow \text{admin}' \text{ OR } 1 \text{ LIKE } 1\#$
Logical Invariant	$LI(\dots e \dots) \rightarrow \dots e \text{ AND } \top \dots$	$LI(\text{admin}' \text{ OR } 1=1\#) \rightarrow \text{admin}' \text{ OR } 1=1 \text{ AND } 2 < > 3\#$



Confidence over Mutation Rounds

WAF-A-MoLE (line) vs. random fuzzer (dots)



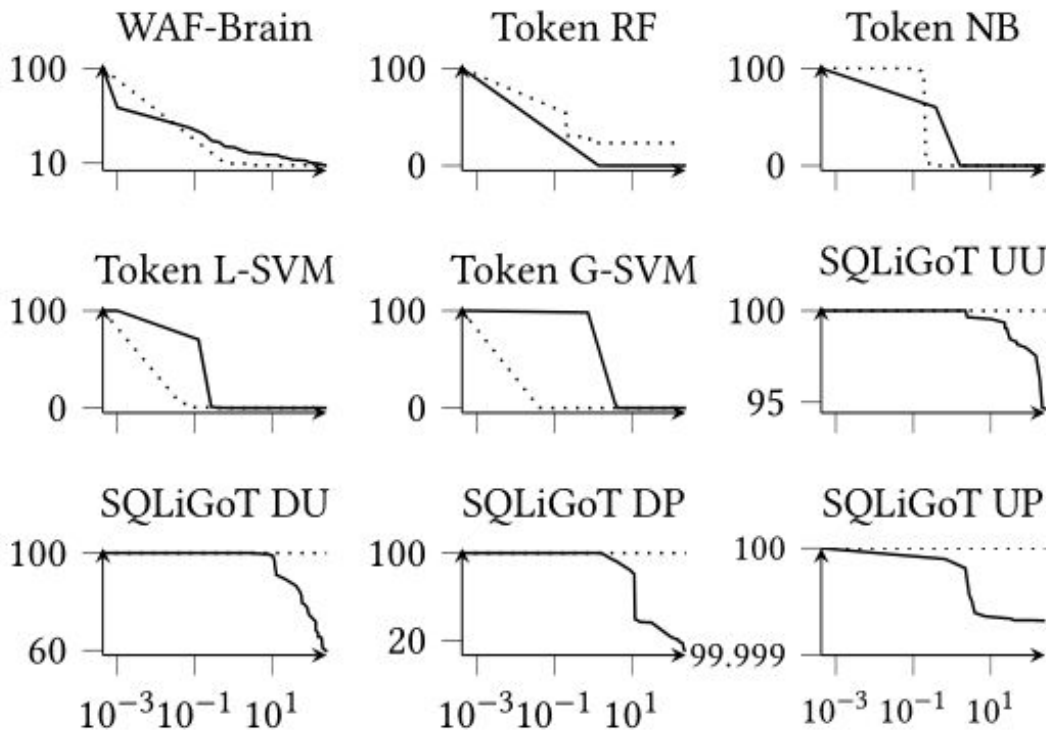
X mutation rounds (1,...,600)
Y WAF confidence (%)

WAF-Brain and Token-based evaded even by **random** (!)

SQLiGoT variants are robust to random mutations, but **vulnerable** to our approach

Confidence over Time

WAF-A-MoLE (line) vs. random fuzzer (dots)



X elapsed time (0,...,100 sec)

Y WAF confidence (%)

WAF-Brain and Token-based evaded in **few seconds**

SQLiGoT variants need more time, especially UP

Conclusions

- The attacker harnesses **infinite** representations of the same payload
- Complexity of features **does not imply** adversarial robustness
- Statistical approaches **may be ineffective** against an aware adversary
- Unexpectedly, **random fuzzing** is sometimes effective
- Guided approach successfully **evades all models**

Thank you for your attention!

- Tool <https://github.com/AvalZ/waf-a-mole>
- Dataset https://github.com/zangobot/wafamole_dataset



{luca.demetrio, andrea.valenza}@dibris.unige.it



@zangobot
@AvalZ_

