# You Only Perturb Once:
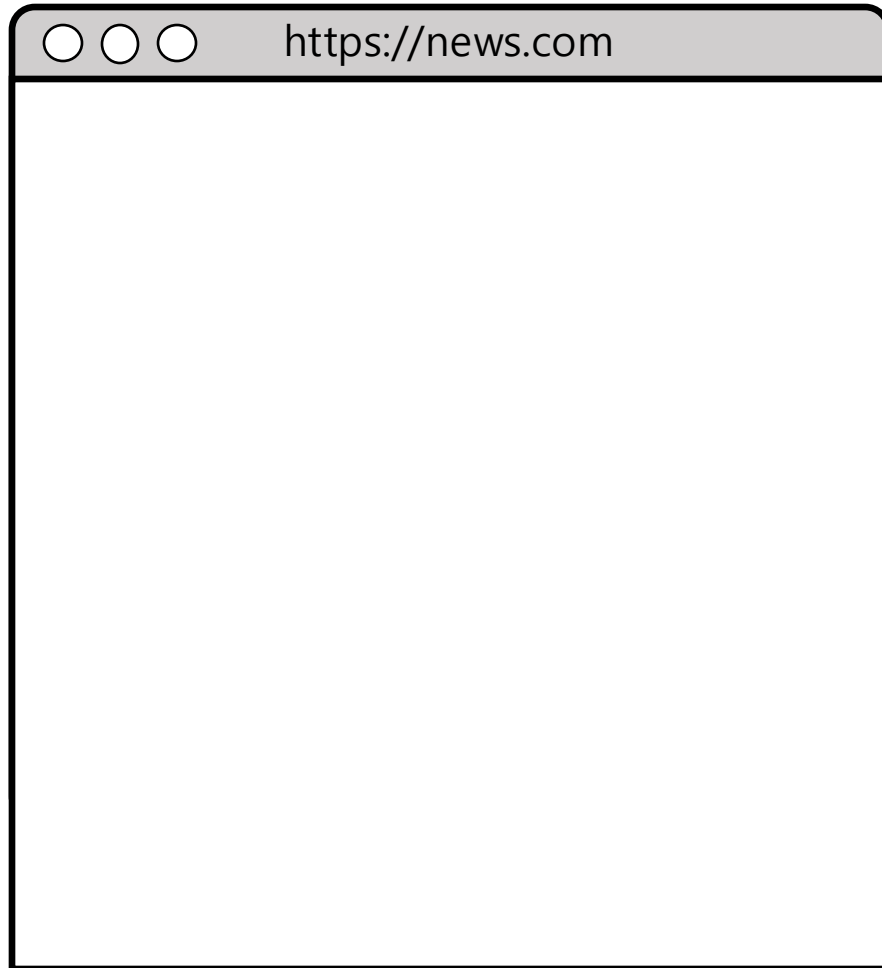# Bypassing (Robust) Ad-Blockers Using Universal Adversarial Perturbations

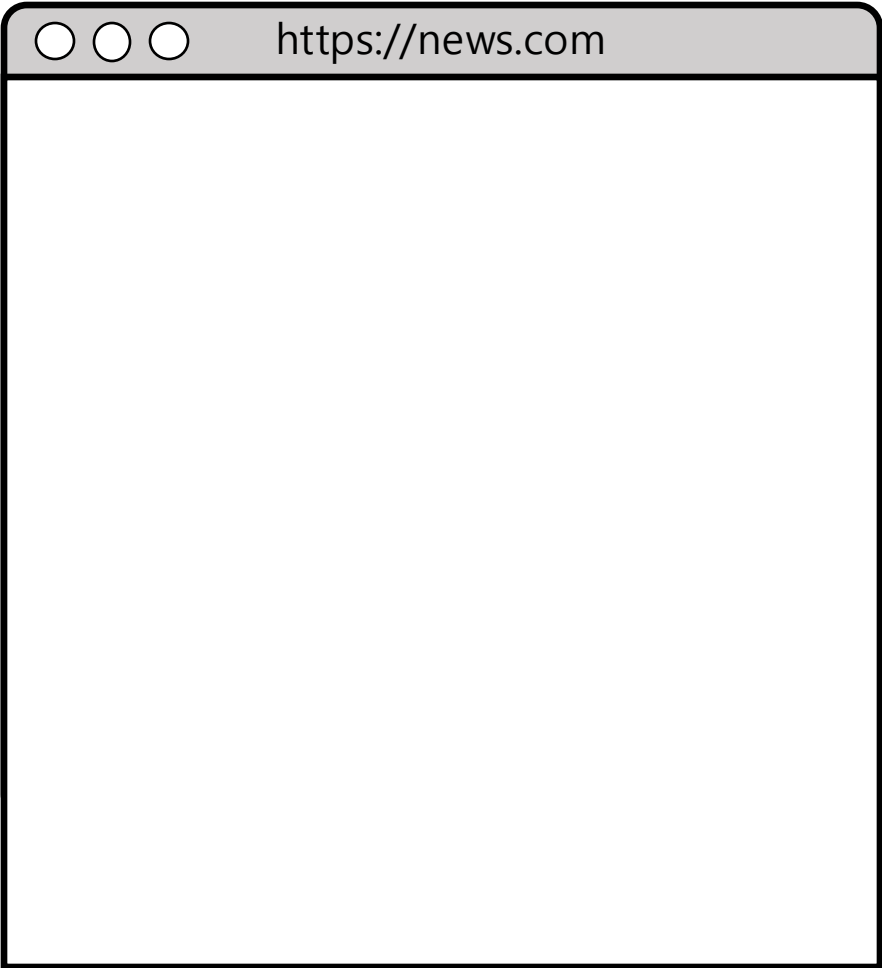**Dongwon Shin**[*], Suyoung Lee[*], Sanghyun Hong[†] and Sooel Son[*]

[*]KAIST, [†]Oregon State University

**ACSAC 2024**

# Online Advertising & Tracking Service (ATS)

# Online Advertising & Tracking Service (ATS)

https://news.com

**Request to https://news.com**

https://news.com

**ATS publisher**

# Online Advertising & Tracking Service (ATS)

https://news.com

Request to https://news.com

https://news.com

**ATS publisher**

# Online Advertising & Tracking Service (ATS)

https://news.com

## HTML

```
<body>
  <iframe src = 'http://ad.com/show_ad'>

  ...

  <script src = 'http://ad.com/track_user.js'>

  ...
</body>
```
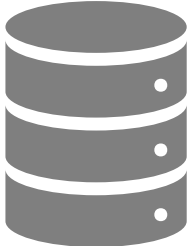
# Online Advertising & Tracking Service (ATS)
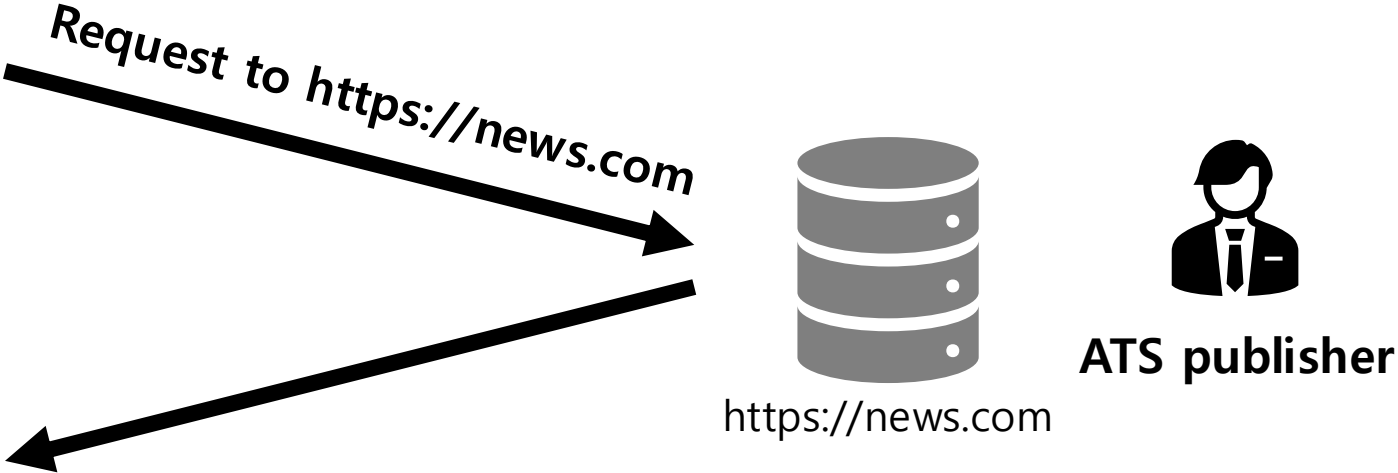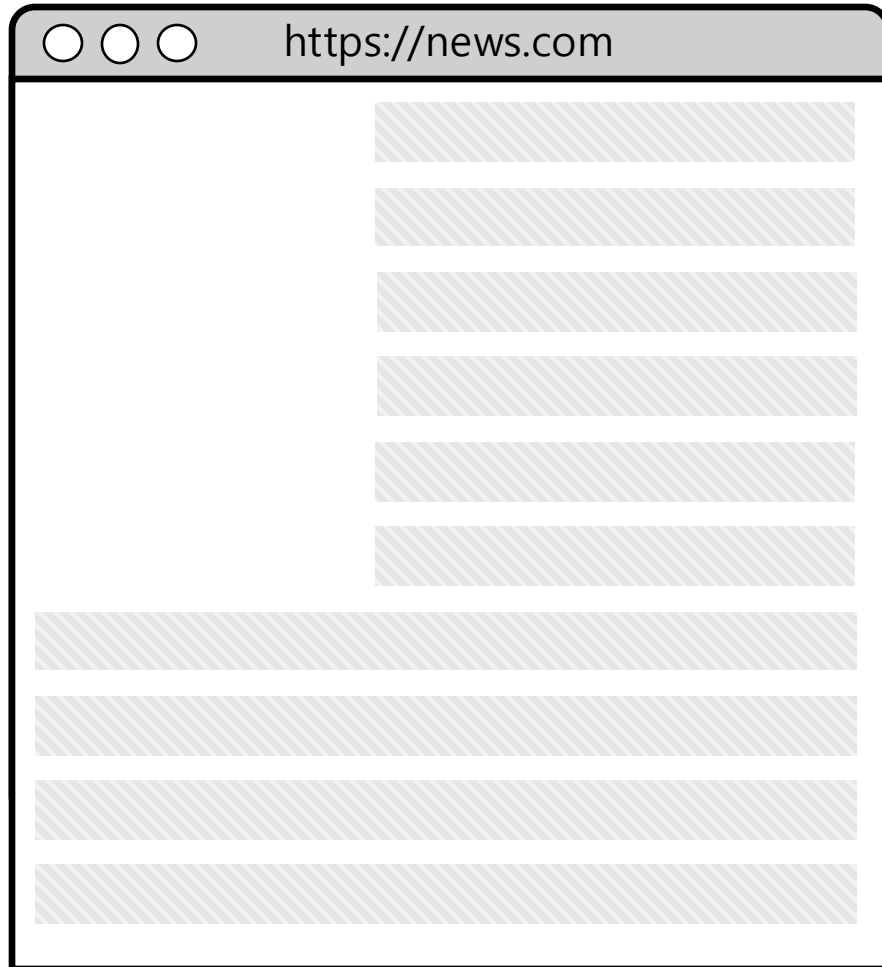
# Online Advertising & Tracking Service (ATS)

https://news.com

**Ad**

**Request to https://news.com**

https://news.com

**ATS publisher**

**Request to http://ad.com**

http://ad.com

**ATS provider**

**They track users' browsing history!**

# ATS blockers



https://news.com

Request to https://news.com

https://news.com

ATS publisher

Request to http://ad.com

http://ad.com

ATS provider

Ad

KAIST Web Security & Privacy Lab Oregon State University SAIL

# ATS blockers



ATS publisher

ATS provider

**ATS blockers block resources fetched from ATS providers!**

# ML-based ATS Blockers

**Webpage**

Ad

https://news.com

**Graph representation**

**Extracted features**

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 97 | 21 | ... |

**Random forest classifier**

**Non-ATS**

**ATS**

## All these steps take place in the client-side

[1] Iqbal et al. AdGraph: A graph-based approach to ad and tracker blocking. S&P '20
[2] Sjosten et al. Filter list generation for underserved regions. WWW '20
[3] Siby et al. WebGraph: Capturing advertising and tracking information flows for robust blocking. Security '22
[4] Lee et al. AdFlush: A real-world deployable machine learning solution for effective advertisement and web tracker prevention. WWW '24

KAIST  Web Security & Privacy Lab  Oregon State University  SAIL

# Adversarial Attacks against ML-based ATS Blockers



**Abusive ATS publishers and providers** may seek to
**bypass ATS blockers** to maximize their profit!

# Adversarial Attacks against ML-based ATS Blockers

https://news.com

Request to https://news.com

**Can ATS publishers/providers bypass these ATS blockers?**

ATS provider

http://ad.com

**Abusive ATS publishers and providers** may seek to
**bypass ATS blockers** to maximize their profit!

KAIST  W Web Security & Privacy Lab  Oregon State University  SAIL

# Adversarial Attacks against ML-based ATS Blockers

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 97 | 21 | ... |

**Extracted features ($x$)**

Random forest classifier ($f$)

Non-ATS

ATS

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 98 | 81 | ... |

**Perturbed features ($x + \delta$)**

Random forest classifier ($f$)

Non-ATS

ATS

Optimize perturbation ($\delta$) **on this request node** to bypass the ATS blocker!

[5] Zhu et al. Eluding ML-based adblockers with actionable adversarial examples. ACSAC '21

KAIST · Web Security & Privacy Lab · Oregon State University · SAIL

# Adversarial Attacks against ML-based ATS Blockers

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/track_user.js | 97 | 27 | ... |

**Extracted features ($x'$)**

However, what if the adversary wants to attack **another network request**?

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 98 | 81 | ... |

Non-ATS

ATS

Random forest classifier ($f$)

Non-ATS

ATS

Random forest classifier ($f$)

**Perturbed features ($x + \delta$)**

Optimize perturbation ($\delta$) **on this request node** to bypass the ATS blocker!

[5] Zhu et al. Eluding ML-based adblockers with actionable adversarial examples. ACSAC '21

# Limitation of Per-Sample Attacks

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/track_user.js | 97 | 27 | ... |

**Extracted features ($x'$)**

Random forest classifier ($f$)

Non-ATS

ATS

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/track_user.js | 101 | 72 | ... |

**Perturbed features ($x' + \delta'$)**

Random forest classifier ($f$)

Non-ATS

The adversary has to optimize perturbation ($\delta'$) **again on this network request node!**

[5] Zhu et al. Eluding ML-based adblockers with actionable adversarial examples. ACSAC '21

KAIST    Web Security & Privacy Lab    Oregon State University    SAIL

# Limitation of Per-Sample Attacks

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/track_user.js | 97 | 27 | ... |

**Can adversaries bypass a target ATS blocker at scale?**

| http://ad.com/track_user.js | 101 | 72 | ... |

Perturbed features $(x' + \delta')$

Non-ATS

Random forest classifier $(f)$

[5] Zhu et al. Eluding ML-based adblockers with actionable adversarial examples. ACSAC '21

KAIST  Web Security & Privacy Lab  Oregon State University  SAIL

# Limitation of Per-Sample Attacks

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/track_user.js | 97 | 27 | ... |

**Multiple samples** $(x_1, x_2, ..., x_n)$

Random forest classifier $(f)$

Non-ATS

ATS

Can adversaries bypass the detection of multiple network requests **using a single perturbation**?

Re

http://ad.c

Perturbed features $(x' + \delta')$

Non-ATS

Random forest classifier $(f)$

[5] Zhu et al. Eluding ML-based adblockers with actionable adversarial examples. ACSAC '21

KAIST  Web Security & Privacy Lab  Oregon State University  SAIL

# You Only Perturb Once

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/track_user.js | 97 | 27 | ... |

# We propose YOPO!

http://ad.com/track_user.js | 101 | 72 | ...

Perturbed features $(x' + \delta')$

**Random forest classifier** $(f)$

[5] Zhu et al. Eluding ML-based adblockers with actionable adversarial examples. ACSAC '21

KAIST   Web Security & Privacy Lab   Oregon State University   SAIL

# Our Contributions

- We show that an adversary can generate **<u>a single and cost-effective universal perturbation</u>** that bypasses recent ML-based ATS blockers.

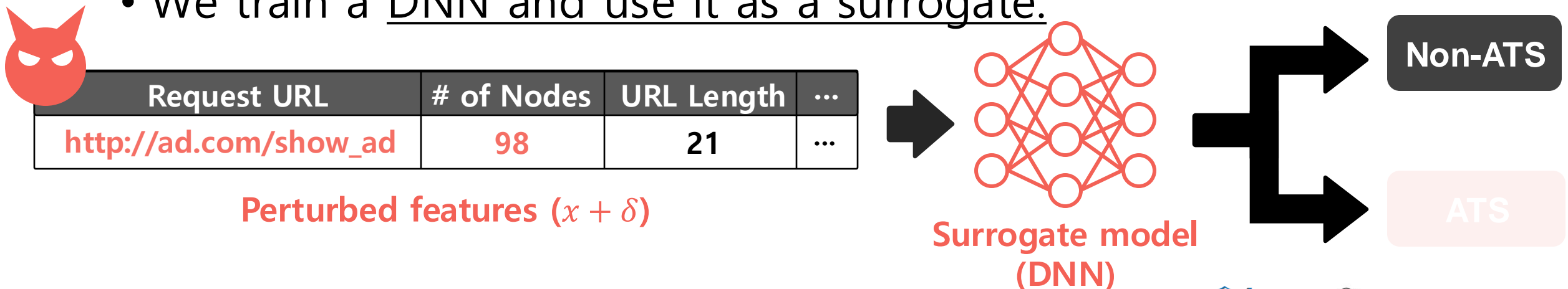- We design and implement a novel framework (YOPO) where one can **<u>generate a universal adversarial perturbation (UAP)</u>** against these ATS blockers.

- We propose two new **<u>mitigation strategies</u>** by analyzing the factors attributing to this vulnerability.

# Challenge #1: Perturbation Optimization

- Random forest classifiers are <u>not differentiable</u>.

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 97 | 21 | ... |

**Extracted features ($x$)**

**Random forest classifier ($f$)**

Non-ATS

ATS

- We train a <u>DNN and use it as a surrogate.</u>

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 98 | 21 | ... |

**Perturbed features ($x + \delta$)**

**Surrogate model (DNN)**

Non-ATS

ATS

# Challenge #2: HTML Manipulation

- Perturbed features should be **reflected in an HTML format**.

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 98 | 21 | ... |

**Perturbed features ($x + \delta$)**

**Webpage**     **Graph representation**
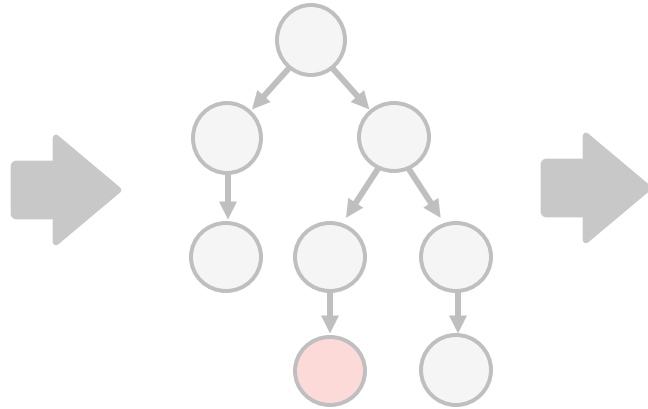
**Random forest classifier**     Non-ATS     ATS

# Challenge #2: HTML Manipulation

- Perturbed features should be **reflected in an HTML format**.



**Webpage**

**Graph representation**

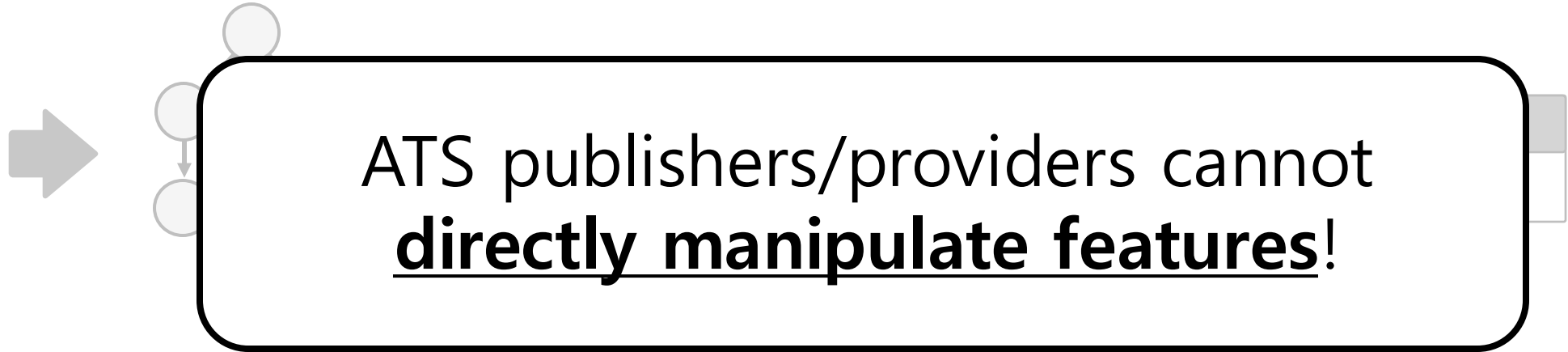ATS publishers/providers cannot
**directly manipulate features**!

Non-ATS

ATS

**Random forest classifier**

# Challenge #2: HTML Manipulation

- Perturbed features should be **reflected in an HTML format**.

o https://news.com
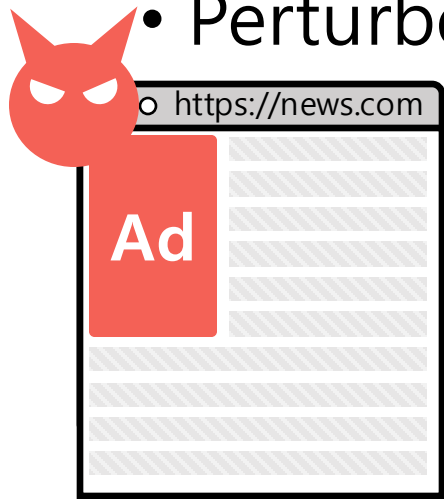
**Ad**

**Webpage
(Manipulated)**

**Graph representation**

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 97 | 21 | ... |

**Extracted features**

**Random forest classifier**

**Non-ATS**

**ATS**

# Challenge #2: HTML Manipulation

- Perturbed features should be **reflected in an HTML format**.

https://news.com

Ad

**Webpage
(Manipulated)**

**Graph representation
(Manipulated)**

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 97 | 21 | ... |

**Extracted features**

**Random forest classifier**

Non-ATS

ATS

KAIST  Web Security & Privacy Lab  Oregon State University  SAIL

# Challenge #2: HTML Manipulation

- Perturbed features should be **reflected in an HTML format**.

https://news.com

**Ad**

**Webpage**
**(Manipulated)**

**Graph representation**
**(Manipulated)**

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 98 | 21 | ... |

**Perturbed features** $(x + \delta)$

**Random forest classifier**

**Non-ATS**

**ATS**

# Challenge #2: HTML Manipulation

- Perturbed features should be **reflected in an HTML format**.

https://news.com

**Ad**

**Webpage
(Manipulated)**

**Graph representation
(Manipulated)**

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 98 | 21 | ... |

**Perturbed features** $(x + \delta)$

**Random forest classifier**

**Non-ATS**

**ATS**

# Challenge #2: HTML Manipulation

- Perturbed features should be **reflected in an HTML format**.
  - ➢ We implement **46 HTML update functions** that reflect different feature changes to a target HTML webpage.

Webpage
(Manipulated)

Graph representation
(Manipulated)

| Request URL | # of Nodes | URL Length | ... |
|---|---|---|---|
| http://ad.com/show_ad | 98 | 21 | ... |

Perturbed features $(x + \delta)$

Non-ATS

Random forest classifier

ATS

# Challenge #3: Preserving Functionalities

- Each feature has **a different functionality breakage risk** of manipulating it.

| HTML |
|------|
| <body><br>  <iframe src = 'http://ad.com/show_ad'><br></body> |

**Original webpage**

| HTML |
|------|
| <body><br>  <iframe src = 'http://ad.com/show_ad?1234'><br></body> |

**Increased the URL length**

☺ **Preserves the functionality!**

# Challenge #3: Preserving Functionalities

- Each feature has **a different functionality breakage risk** of manipulating it.
  - ➤ We **designed a cost model** that prioritizes which features to manipulate first.

| HTML |
|---|
| **\<body\>**<br>  \<iframe src = 'http://ad.com/show_ad'\><br>**\</body\>** |

**Original webpage**

| HTML |
|---|
| **\<script\>**<br>  \<iframe src = 'http://ad.com/show_ad'\><br>**\</script\>** |

**Changed its parent tag name**

⚠️ **Breaks the functionality!**

# YOPO Overview

One-time preparation step!



Webpages

Graph representation

| # of Nodes | URL Length | Semicolon | Label |
|---|---|---|---|
| 97 | 21 | FALSE | ATS |
| 97 | 27 | FALSE | ATS |
| 108 | 24 | TRUE | Non-ATS |
| 81 | 29 | TRUE | Non-ATS |

Data instances

Surrogate model

| # of Nodes | URL Length | Semicolon | Label |
|---|---|---|---|
| 97 | 21 | FALSE | ATS |
| 97 | 27 | FALSE | ATS |

Sampled ATS instances

**1. Data collection**    **2. Surrogate training**    **3. UAP generation**

Non-ATS

Misclassified result

Target ATS Blocker

Ad

Manipulated webpages

| # of Nodes | URL Length | Semicolon |
|---|---|---|
| 99 | 26 | TRUE |
| 99 | 32 | TRUE |

UAP-injected ATS instances

| # of Nodes | URL Length | Semicolon |
|---|---|---|
| +2 | +5 | TRUE |

Universal adversarial perturbation

**5. ATS classification**    **4. HTML/JS manipulation**

KAIST    Web Security & Privacy Lab    Oregon State University    SAIL

# Data Collection for Training a Surrogate Model

Tranco's Top-10K
Webpages

Graph representation

| Request URL | # of Nodes | URL Length | Semicolon | # of Cookie Read |
|---|---|---|---|---|
| http://ad.com/show_ad | 97 | 21 | FALSE | 0 |
| http://ad.com/track_user.js | 97 | 27 | FALSE | 1 |
| https://sec.com/logo.png | 108 | 24 | TRUE | 0 |
| https://acsac.org/favicon.ico | 81 | 29 | TRUE | 3 |

**Extracted features**

ATS blockers classify
**network request nodes**!

Target
ATS Blocker

| Request URL | # of Nodes | URL Length | Semicolon | # of Cookie Read | Label |
|---|---|---|---|---|---|
| http://ad.com/show_ad | 97 | 21 | FALSE | 0 | ATS |
| http://ad.com/track_user.js | 97 | 27 | FALSE | 1 | ATS |
| https://sec.com/logo.png | 108 | 24 | TRUE | 0 | Non-ATS |
| https://acsac.org/favicon.ico | 81 | 29 | TRUE | 3 | Non-ATS |

**Labeled features**

KAIST  Web Security & Privacy Lab  Oregon State University  SAIL

# Surrogate Model Training

- We selected a four-layer feed-forward neural network as our surrogate.

| Request URL | # of Nodes | URL Length | Semicolon | # of Cookie Read | Label |
|---|---|---|---|---|---|
| http://ad.com/show_ad | 97 | 21 | FALSE | 0 | ATS |
| http://ad.com/track_user.js | 97 | 27 | FALSE | 1 | ATS |
| https://sec.com/logo.png | 108 | 24 | TRUE | 0 | Non-ATS |
| https://acsac.org/favicon.ico | 81 | 29 | TRUE | 3 | Non-ATS |

**Data instances**

Train

**Surrogate model (DNN)**

# UAP Generation

| # of Nodes | URL Length | Semicolon | # of Cookie Read | Label |
|---|---|---|---|---|
| 97 | 21 | FALSE | 0 | ATS |
| 97 | 27 | FALSE | 1 | ATS |

**Sampled 40K ATS instances ($x$)**

**+**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|---|---|---|---|
| +0 | +0 | FALSE | +0 |

**Universal Adversarial Perturbation ($\delta$)**

We optimize the perturbation over
**all these instances** to make it universal!

**Surrogate model
(DNN)**

ATS

# UAP Generation

- **Optimization goal**

1) Bypass the target ATS blocker with a single perturbation

---

$$\underset{\delta}{\mathrm{argmax}} \; \boldsymbol{E}_{(x,y) \sim D_{ATS}}[L_{CE}(f'(\theta, x + \delta), ATS)]$$

$\delta$: Perturbation $\qquad\qquad$ $f'(\theta)$: Surrogate model (DNN)

$D_{ATS}$: Data instances labeled as ATS $\quad$ $L_{CE}$ : Cross-Entropy Loss

# UAP Generation

- **Optimization goal**

1) Bypass the target ATS blocker with a single perturbation

2) Minimize the breakage risk of manipulating each feature

How can we achieve this goal?

$$\underset{\delta}{\operatorname{argmax}}\ \boldsymbol{E}_{(x,y)\sim D_{ATS}}[L_{CE}(f'(\theta, x+\delta), ATS)]$$

$\delta$: Perturbation

$D_{ATS}$: Data instances labeled as ATS

$f'(\theta)$: Surrogate model (DNN)

$L_{CE}$ : Cross-Entropy Loss

# UAP Generation

- **Optimization goal**

  1) Bypass the target ATS blocker with a single perturbation

  2) Minimize the breakage risk of manipulating each feature

Simply minimizing the perturbation size is **insufficient**!

$$\underset{\delta}{\arg\max} \; \boldsymbol{E}_{(x,y)\sim D_{ATS}}[L_{CE}(f'(\theta, x + \delta), ATS)] - \|\delta\|$$

$\delta$: Perturbation      $f'(\theta)$: Surrogate model (DNN)

$D_{ATS}$: Data instances labeled as ATS      $L_{CE}$ : Cross-Entropy Loss

KAIST   Web Security & Privacy Lab   Oregon State University   SAIL

# Cost Model

- **Optimization goal**

  1) Bypass the target ATS blocker with a single perturbation

  2) Minimize the breakage risk of manipulating each feature

Considered **a web-specific cost**!

$$\underset{\delta}{\mathrm{argmax}} \; \boldsymbol{E}_{(x,y)\sim D_{ATS}}[L_{CE}(f'(\theta, x + \delta), ATS)] - C \cdot \|\delta\|$$

$\delta$: Perturbation                                  $f'(\theta)$: Surrogate model (DNN)

$D_{ATS}$: Data instances labeled as ATS    $L_{CE}$ : Cross-Entropy Loss

KAIST    Web Security & Privacy Lab    Oregon State University    SAIL

# Cost Model

Prioritize with **a lower cost value**!

| Perturbation | Assigned Cost |
|---|---|
| URL_LENGTH | 0.2 |
| | |

Cost model

| HTML |
|---|
| <body><br>  <iframe src = 'http://ad.com/show_ad'><br></body> |

**Original webpage**

| HTML |
|---|
| <body><br>  <iframe src = 'http://ad.com/show_ad?1234'><br></body> |

**Increased the URL length**

☺ **Preserves the functionality!**

# Cost Model

| Perturbation | Assigned |
|---|---|
| URL_LENGTH | 0.2 |
| PARENT_TAG_NAME | 3 |

**Deprioritize with a higher cost value!**

**Cost model**

| HTML |
|---|
| **\<body\>**<br>  \<iframe src = 'http://ad.com/show_ad'\><br>**\</body\>** |

**Original webpage**

| HTML |
|---|
| **\<script\>**<br>  \<iframe src = 'http://ad.com/show_ad'\><br>**\</script\>** |

**Changed its parent tag name**

⚠️ **Breaks the functionality!**

# Cost Model

| Perturbation | Assigned Cost |
|---|---|
| URL_LENGTH | 0.2 |
| PARENT_TAG_NAME | 3 |
| … | … |

**Cost model**

The cost model represents **a relative risk of manipulating it**.

| HTML |
|---|
| **<body>**<br>  <iframe src = 'http://ad.com/show_ad'><br>**</body>** |

**Original webpage**

| HTML |
|---|
| **<script>**<br>  <iframe src = 'http://ad.com/show_ad'><br>**</script>** |

**Changed its parent tag name**

⚠️ **Breaks the functionality!**

# HTML/JS Manipulation

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|:---:|:---:|:---:|:---:|
| +2 | +5 | TRUE | +4 |

**Universal Adversarial Perturbation ($\delta$)**

**+**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|:---:|:---:|:---:|:---:|
| 97 | 21 | FALSE | 0 |
| 97 | 27 | FALSE | 1 |

**Target ATS instances ($x$)**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|:---:|:---:|:---:|:---:|
| 99 | 26 | TRUE | 4 |
| 99 | 32 | TRUE | 5 |

**UAP-injected ATS instances ($x + \delta$)**

Ad

**Webpages**

Ad

**Manipulated Webpages**

# HTML/JS Manipulation

| HTML |
|------|
| **<body>** |
|   **<iframe** src = 'http://ad.com/show_ad'> |
|   **...** |
|   **<script** src = 'http://ad.com/track_user.js'> |
|   **...** |
| |
| **</body>** |

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|:----------:|:----------:|:---------:|:----------------:|
| +2 | +5 | TRUE | +4 |

**UAP**

**Target network request!**

| track_user.js |
|---------------|
| **...** |
| **// Tracking Users** |
| **user_cookie = document.cookie;** |
| |
| |

**Webpage**

# HTML/JS Manipulation

| HTML |
|---|
| `<body>`<br>  `<iframe src = 'http://ad.com/show_ad'>`<br>  ...<br>  `<script src = 'http://ad.com/track_user.js'>`<br>  ...<br><br>`</body>` |

| track_user.js |
|---|
| ...<br>*// Tracking Users*<br>`user_cookie = document.cookie;` |

**Webpage**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|:---:|:---:|:---:|:---:|
| +2 | +5 | TRUE | +4 |

**UAP**

`<body>`

`<iframe>`     `<script>`

http://ad.com/track_user.js

...

Cookie storage

**Graph representation**

# HTML/JS Manipulation

| HTML |
|---|
| `<body>` |
|   `<iframe src = 'http://ad.com/show_ad'>` |
|   ... |
|   `<script src = 'http://ad.com/track_user.js'>` |
|   ... |
| |
| `</body>` |

| track_user.js |
|---|
| ... |
| // Tracking Users |
| user_cookie = document.cookie; |
| |
| |

**Webpage**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|---|---|---|---|
| +2 | +5 | TRUE | +4 |

**UAP**

`<body>`

`<iframe>`  `<script>`

http://ad.com/track_user.js

...

Cookie storage

**Graph representation**

# HTML/JS Manipulation

| HTML |
|---|
| `<body>` |
|   `<iframe` src = 'http://ad.com/show_ad'> |
|   ... |
|   `<script` src = 'http://ad.com/track_user.js'> |
|   ... |
|   **`<div hidden=""></div>`** |
|   **`<p hidden=""></p>`** |
| `</body>` |

| track_user.js |
|---|
| ... |
| // Tracking Users |
| user_cookie = document.cookie; |
| |

**Webpage**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|---|---|---|---|
| +2 | +5 | TRUE | +4 |

**UAP**

`<body>`

`<iframe>`  `<script>`

http://ad.com/track_user.js

...

Cookie storage

**Graph representation**

# HTML/JS Manipulation

| HTML |
|---|
| `<body>` |
|   `<iframe src = 'http://ad.com/show_ad'>` |
|   ... |
|   `<script src = 'http://ad.com/track_user.js'>` |
|   ... |
|   **`<div hidden=""></div>`** |
|   **`<p hidden=""></p>`** |
| `</body>` |

| track_user.js |
|---|
| ... |
| *// Tracking Users* |
| `user_cookie = document.cookie;` |

**Webpage**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|---|---|---|---|
| +2 | +5 | TRUE | +4 |

**UAP**

`<body>`

`<iframe>`    `<script>`    `<div>`    `<p>`

http://ad.com/track_user.js

...

Cookie storage

**Graph representation**

# HTML/JS Manipulation

| HTML |
|---|
| `<body>` |
|   `<iframe src = 'http://ad.com/show_ad'>` |
|   `...` |
|   `<script src = 'http://ad.com/track_user.js?123;'>` |
|   `...` |
|   `<div hidden=""></div>` |
|   `<p hidden=""></p>` |
| `</body>` |

| track_user.js |
|---|
| `...` |
| `// Tracking Users` |
| `user_cookie = document.cookie;` |
| `for (let i = 1; i <= 4; i++) {` |
|   `getCookie();` |
| `}` |

**Webpage**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|---|---|---|---|
| +2 | +5 | TRUE | +4 |

**UAP**

`<body>`

`<iframe>`    `<script>` `<div>`    `<p>`

http://ad.com/track_user.js?123;

...

Cookie storage

**Graph representation**

# HTML/JS Manipulation

| HTML |
|---|
| <body> |
|   <iframe src = 'http://ad.com/show_ad'> |
|   … |

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|---|---|---|---|
| +2 | +5 | TRUE | +4 |

UAP

## YOPO implements 46 HTML/JS manipulation functions!

```
…
// Tracking Users
user_cookie = document.cookie;
for (let i = 1; i <= 4; i++) {
  getCookie();
}
```

…

Cookie storage

Webpage

Graph representation

# Experimental Setup

- Target ATS blockers
    1) **AdGraph** [S&P '20]

    2) **WebGraph*** [Security '22]

    3) **AdFlush** [WWW '24]

    4) **PageGraph*** [WWW '20]

    * We used all content, structural, and flow features for WebGraph.

    ** We revised PageGraph to support all ATS resource types.

- We measured **attack success rate (ASR)** against 2,000 target ATS requests.

# Attack Success Rate

- ASRs measured against target ATS blockers

| ATS blockers | Attack success rate (%) |
|:---:|:---:|
| AdGraph | 89.27 |
| WebGraph | 71.21 |
| AdFlush | 61.91 |
| PageGraph | 84.16 |

Recent ML-based ATS blockers are **vulnerable to universal attacks** using a single perturbation!

# Attack Success Rate

- ASRs measured against target ATS blockers

| | |
|---|---|
| | |
| **Adversaries can launch attacks against these ATS blockers at scale!** | |
| PageGraph | 84.16 |

Recent ML-based ATS blockers are **vulnerable to universal attacks** using a single perturbation!

# Attack Success Rate

- ASRs measured against target ATS blockers

| | |
|---|---|
| **Where does this vulnerability stem from?** | |
| PageGraph | 84.16 |

Recent ML-based ATS blockers are **vulnerable to universal attacks** using a single perturbation!

# Top-5 Most Influential Features

| Features | Type | ASR drop (↓) |
|---|---|---|
| PARENT_ATTR_ASYNC | Binary | **-19.87%** |
| SEMICOLON_IN_URL | Binary | **-8.28%** |
| PARENT_ATTR_DEFER | Binary | **-5.83%** |
| DOMAIN_NAME_IN_QS | Binary | **-5.52%** |
| URL_LENGTH | Numerical | **-1.89%** |

**Top-5 most influential features for attacking AdGraph**

# Top-5 Most Influential Features

| Features | Type | ASR drop (↓) |
|---|---|---|
| PARENT_ATTR_ASYNC | **Binary** | -19.87% |
| SEMICOLON_IN_URL | **Binary** | -8.28% |
| PARENT_ATTR_DEFER | **Binary** | -5.83% |
| DOMAIN_NAME_IN_QS | **Binary** | -5.52% |
| URL_LENGTH | Numerical | -1.89% |

When applying a UAP, YOPO overwrites binary features to have **a specific combination of values**!

# Top-5 Most Influential Features

| Features | UAP Values | ASR drop (↓) |
|---|---|---|
| PARENT_ATTR_ASYNC | **TRUE** | **-19.87%** |
| SEMICOLON_IN_URL | **TRUE** | **-8.28%** |
| PARENT_ATTR_DEFER | **FALSE** | **-5.83%** |
| DOMAIN_NAME_IN_QS | **FALSE** | **-5.52%** |
| URL_LENGTH | +9 | -1 |

Non-ATS (98.48%)

ATS (1.52%)

**98.48%** of network requests that have
**this combination** in our training set are **non-ATS**!

KAIST · W Web Security & Privacy Lab · Oregon State University · SAIL

# Top-5 Most Influential Features

| Features | UAP Values | ASR drop (↓) |
|---|:---:|:---:|
| PARENT_ATTR_ASYNC | TRUE | -19.87% |
| SEMICOLON_IN_URL | TRUE | -8.28% |
| PARENT_ATTR_DEFER | FALSE | -5.83% |
| DOMAIN_NAME_IN_QS | FALSE | -5.52% |
| URL_LENGTH | +9 | -1.89% |

Non-ATS (98.48%)

ATS (1.52%)

| PARENT_ATTR_ASYNC | SEMICOLON_IN_URL | PARENT_ATTR_DEFER | DOMAIN_NAME_IN_QS |
|---|---|---|---|
| FALSE | FALSE | FALSE | TRUE |

ATS

| PARENT_ATTR_ASYNC | SEMICOLON_IN_URL | PARENT_ATTR_DEFER | DOMAIN_NAME_IN_QS |
|---|---|---|---|
| TRUE | TRUE | FALSE | FALSE |

Non-ATS

# Top-5 Most Influential Features

| Features | UAP Values | ASR drop (↓) |
|---|---|---|
| PARENT_ATTR_ASYNC | TRUE | -19.87% |

Non-ATS (98.48%)

**This arises from the inherent imbalance of binary features in real-world webpages!**

| PARENT_ATTR_ASYNC | SEMICOLON_IN_URL | PARENT_ATTR_DEFER | DOMAIN_NAME_IN_QS | ATS |
|---|---|---|---|---|
| FALSE | FALSE | FALSE | TRUE | |

| PARENT_ATTR_ASYNC | SEMICOLON_IN_URL | PARENT_ATTR_DEFER | DOMAIN_NAME_IN_QS | Non-ATS |
|---|---|---|---|---|
| TRUE | TRUE | FALSE | FALSE | |

KAIST  Web Security & Privacy Lab  Oregon State University  SAIL
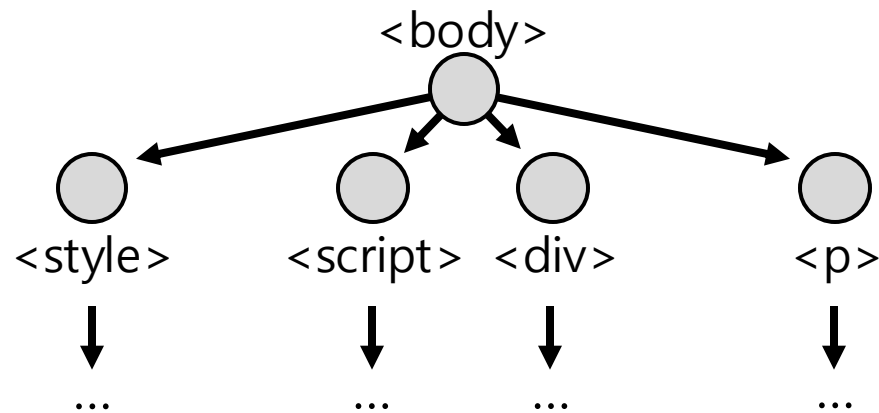
# Mitigation #1: Nullifying Binary Features

- We **nullified binary feat**s
when training each ATS

> Reduced the ASR **by at most 27.52%**
> without any performance drop!

| ATS blockers | ASR | Accuracy | Precision | Recall |
|---|---|---|---|---|
| AdGraph | 61.75 (27.52 ↓) | 92.15 (0.49 ↓) | 89.11 (0.79 ↓) | 84.43 (0.87 ↓) |
| WebGraph | 63.90 (7.31 ↓) | 95.39 (0.29 ↓) | 92.52 (0.74 ↓) | 92.08 (0.19 ↓) |
| AdFlush | 49.82 (12.09 ↓) | 95.79 (0.14 ↓) | 93.99 (0.35 ↓) | 90.77 (0.10 ↓) |
| PageGraph | 70.28 (13.88 ↓) | 95.78 (0.11 ↓) | 92.66 (0.24 ↓) | 93.06 (0.39 ↓) |

# Mitigation #2: Misleading Perturbation Directions

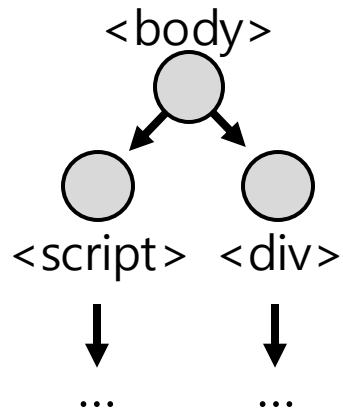- HTML manipulation **<u>decreasing feature values</u>** is more likely to **<u>break webpages</u>**.



**Graph representation**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|------------|------------|-----------|------------------|
| -2 | +5 | TRUE | +4 |

**UAP**

# Mitigation #2: Misleading Perturbation Directions

- HTML manipulation **<u>decreasing feature values</u>** is more likely to **<u>break webpages</u>**.



**Graph representation**

⚠️ **Breaks the functionality!**

| # of Nodes | URL Length | Semicolon | # of Cookie Read |
|:---:|:---:|:---:|:---:|
| -2 | +5 | TRUE | +4 |

**UAP**

# Mitigation #2: Misleading Perturbation Directions

- HTML manipulation **decreasing feature values** is more likely to **break webpages**.

- We **preprocessed input features** before training ATS blockers, thus **misleading perturbations to decrease feature values**.

- As a result, adversaries **cannot reflect such manipulation** at an HTML level.

# Mitigation #2: Misleading Perturbation Directions

Applying both mitigation strategies reduced the ASR **by at most 48.86%** without any performance drop!
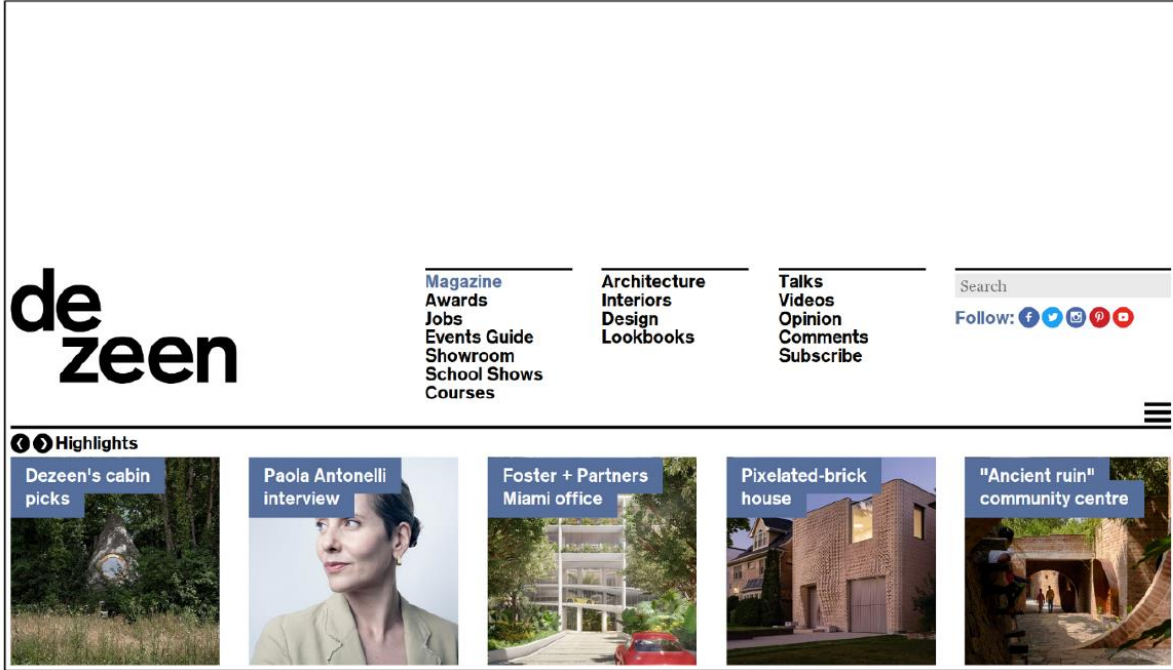
| ATS blockers | ASR | Accuracy | Precision | Recall |
|---|---|---|---|---|
| AdGraph | 40.41 (48.86 ↓) | 91.59 (1.05 ↓) | 85.49 (4.41 ↓) | 87.05 (1.75 ↑) |
| WebGraph | 48.55 (22.66 ↓) | 95.19 (0.49 ↓) | 91.64 (1.62 ↓) | 92.38 (0.11 ↑) |
| AdFlush | 42.74 (19.17 ↓) | 95.68 (0.25 ↓) | 95.00 (0.34 ↓) | 90.34 (0.53 ↓) |
| PageGraph | 64.51 (19.65 ↓) | 95.74 (0.15 ↓) | 92.59 (0.31 ↓) | 93.26 (0.19 ↓) |

KAIST · Web Security & Privacy Lab · Oregon State University · SAIL
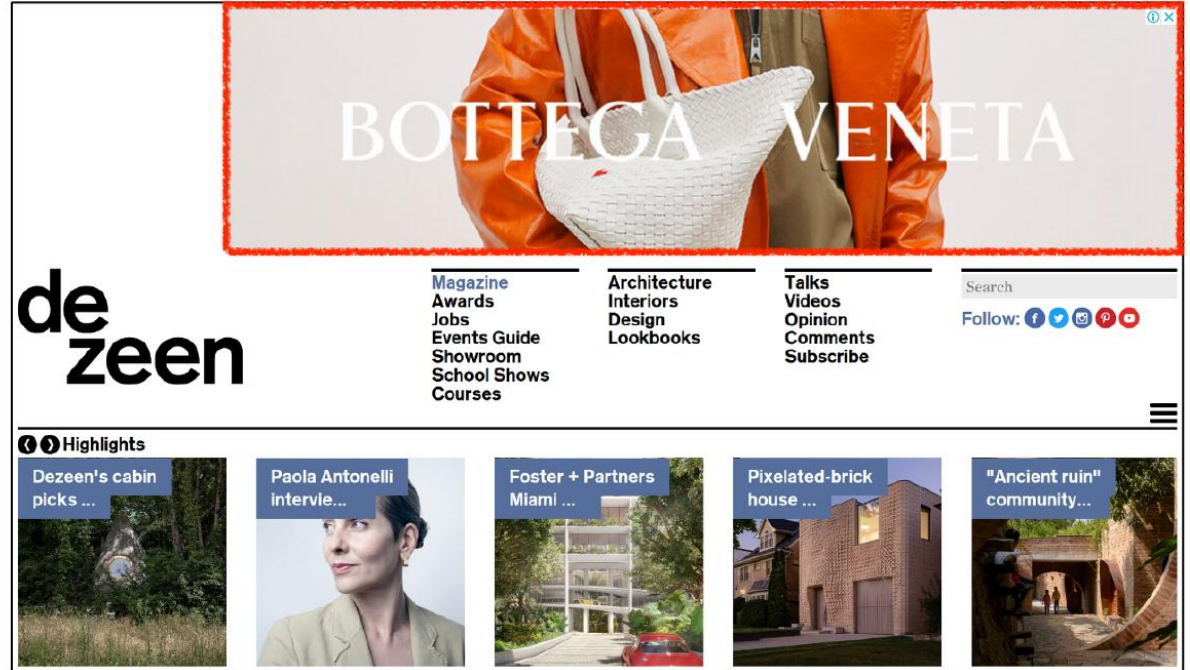
# Breakage Analysis

- We manually inspected 400 webpages manipulated by YOPO.

- We defined 4 breakage types following prior studies.

- **Only 14 webpages out of 400** exhibited functionality disruption.

KAIST  Web Security & Privacy Lab  Oregon State University  SAIL

# Breakage Analysis

Original webpage

Manipulated webpage

No functional breakage

# For More Details

- Case study

- Effect of the attack hyperparameters

- Attacking multiple requests

- Different cost models

- https://github.com/WSP-LAB/YOPO

# Question?