# Duumviri: Detecting Trackers and Mixed Trackers with a Breakage Detector

**He (Shawn) Shuang**, University of Toronto

Lianying Zhao, University of Toronto/Carleton University

David Lie, University of Toronto

# Tracker Blockers

Introduction     Overview     Breakage Detector     Conclusion / Discussions

# Filter Lists

Human-written rules

```
_adobe_analytics.js
_track_pixel.gif?
://analytics.*/collect
/clicklog.js
&action=js_stats&
&ev=PageView&
```

**Adblock Plus**

> Ads 1/3
> ☑ EasyList 👁 🏠 🕐 77,579 used out of 78,174
> Privacy 1/4
> ☑ EasyPrivacy 👁 🏠 🕐 52,618 used out of 53,268

Introduction          Overview          Breakage Detector          Conclusion / Discussions

# Filter Lists Problem



Human-written rules

_adobe_analytics.js
_track_pixel.gif?
://analytics.*/collect
/clicklog.js
&action=js_stats&
&ev=PageView&

> Ads 1/3
  ☑ EasyList 👁 🏠 🕐 77,579 used out of 78,174
> Privacy 1/4
  ☑ EasyPrivacy 👁 🏠 🕐 52,618 used out of 53,268

Problems:
1. Scalability
2. Accuracy

# Previous Works

Feature sources

1. **HTML features**
   a. E.g., whether iframe element is the ancestor of a third party script element
2. **Network layer features**
   a. E.g., whether the request contains over three query string parameters
3. **JavaScript execution features**
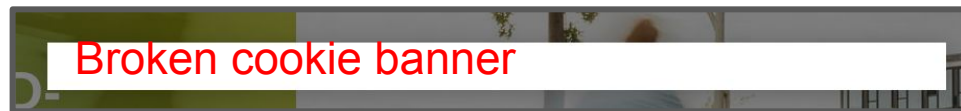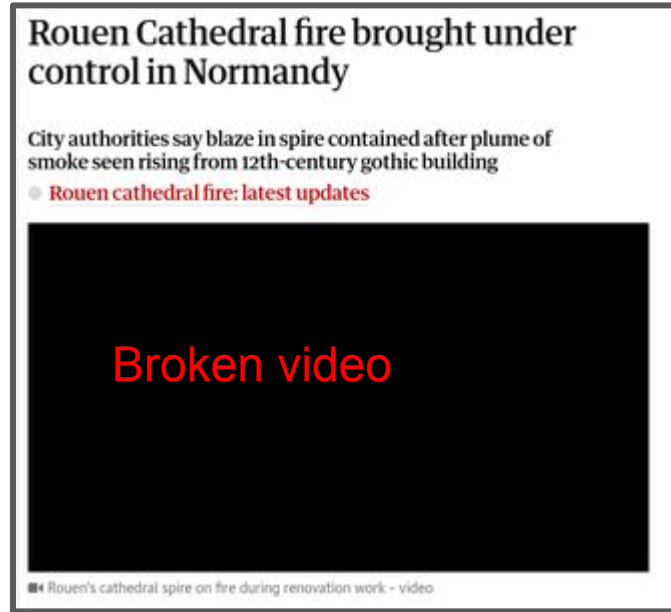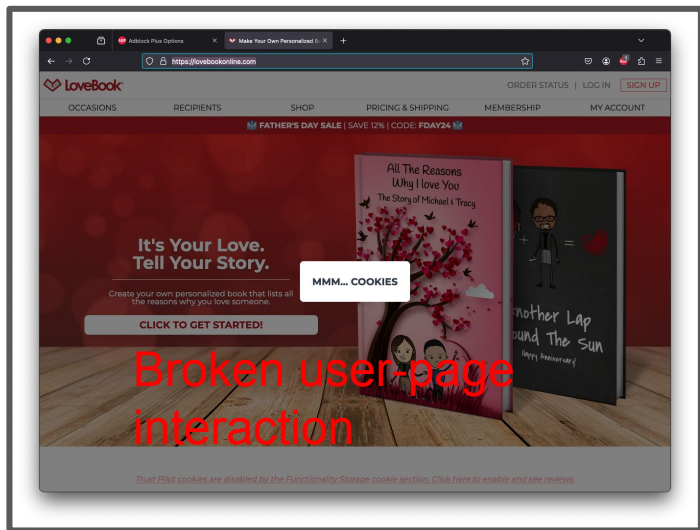   a. E.g., whether a JavaScript called *Canvas* API before sending a request
4. **Graph features** involves the interactions between HTML elements, network requests and JavaScript executions
   a. E.g., flow of information from cookiejar to requests
   b. [AdGraph SP 2020], [WebGraph USENIX 2022]

Introduction          Overview          Breakage Detector          Conclusion / Discussions

# High Breakage Rate

Problem: current tracker detector can cause high breakage

- Previous work (e.g., [AdGraph SP 2020]) breaks 15% of pages

Rouen Cathedral fire brought under control in Normandy

City authorities say blaze in spire contained after plume of smoke seen rising from 12th-century gothic building

● Rouen cathedral fire: latest updates

Broken video

Rouen's cathedral spire on fire during renovation work – video

Broken user-page interaction

Broken cookie banner

Introduction     Overview     Breakage Detector     Conclusion / Discussions

# Reasons for Breakage

**Reason #1:** functional request is misidentified as tracking request

- Trained from an imperfect data source
- Hard to improve the quality as developers working on it

**Reason #2:** functional information cannot be separated from tracking requests

- Mixed tracker

For functionality

For Tracking

# Mixed Trackers Example

Example mixed tracker

URL: https://www.EXAMPLE.com/landingpage-a-psurl.html?
goods_id=601099526089385&
sku_id=17592258865022&
_x_ns_msclkid=eeec99c83e911b00583ffc4bc3e34060

Consequences

- Blocking the request causes the redirect to fail
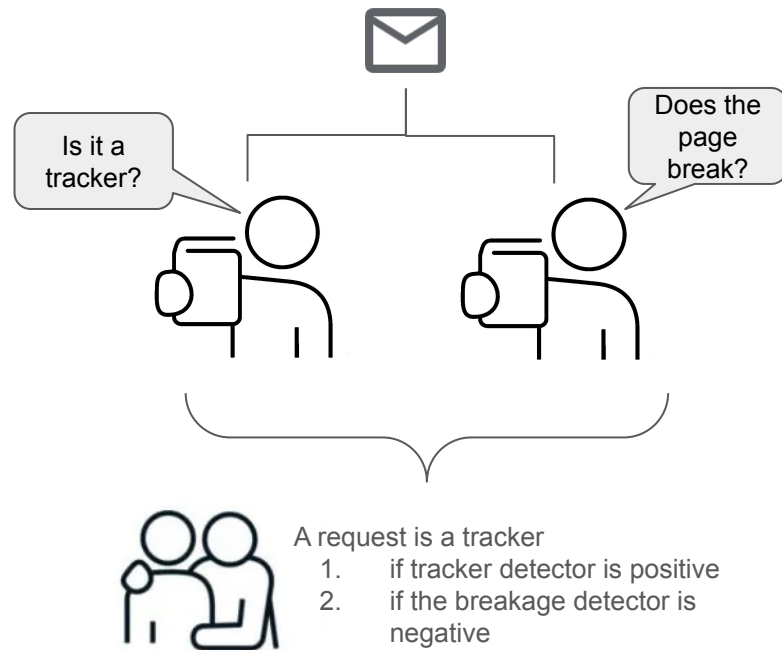- Allowing the request hurts privacy

# Duumviri

Duumrivi: addressing the high breakage rate of previous works
- Maximizing privacy: identifying trackers
- Minimizing breakage: not breaking web pages

Contributions

1. A two-modeled approach for tracker detection with a dedicated breakage detector
2. Detecting mixed trackers automatically

Is it a tracker?

Does the page break?

A request is a tracker
1. if tracker detector is positive
2. if the breakage detector is negative

Introduction          Overview          Breakage Detector          Conclusion / Discussions
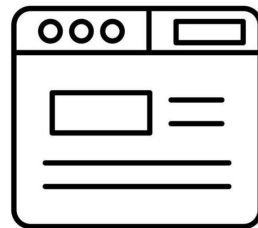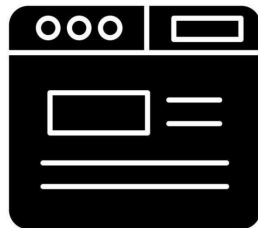
# Breakage Detector

Goal: a model that can predict when a web page is broken due to tracker misidentification

1. What is a broken page in the context of tracker detection? How do we determine breakage?
   a. We define a breakage to be changes to web pages from the origin page subjectively determined by the human user
2. How do we collect samples of breakages?
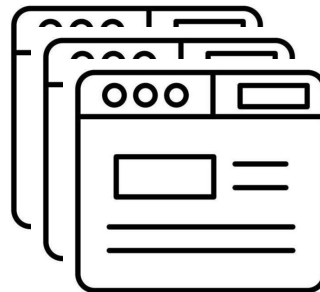
**Unbroken page**

**Changed page**

Compare and draw **differential features**

# Breakage Detector - Data Samples

**Unbroken page**

**Changed page**

Problem: need instances of unbroken page and changed pages to draw differential features

- We cannot crawl the web for breakages
    - Breakage reports are often fixed quickly, live sites with breakages are rare
    - Even if we could find live breakages, we do not have an unbroken version of the page (as the page is currently broken)
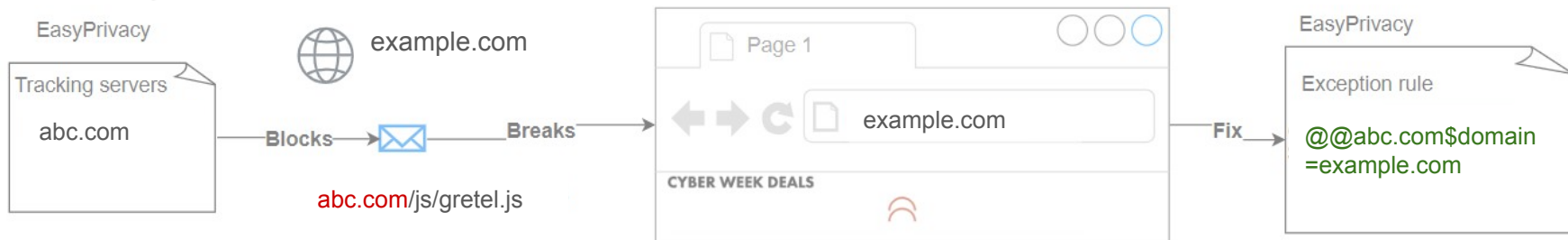- We do not want crafted breakages

# Exception Rules



EasyPrivacy

Tracking servers
abc.com

example.com

Blocks → ✉ → Breaks

abc.com/js/gretel.js

Page 1

example.com

CYBER WEEK DEALS

Fix →

EasyPrivacy

Exception rule

@@abc.com$domain
=example.com

# Breakage Reconstruction

EasyPrivacy

Tracking servers
abc.com

🌐 example.com

**Blocks** → ✉ **Breaks** →

abc.com/js/gretel.js

Page 1

example.com

CYBER WEEK DEALS

**Fix** →

EasyPrivacy

Exception rule

@@abc.com$domain
=example.com

---

Exception rules in EasyPrivacy

@@abc.com$domain=
example.com
…

🌐 example.com

Flipped rule:
abc.com

**Blocks** →

✉

abc.com/js/gretel.js

**Breaks**

Page 1

example.com

CYBER WEEK DEALS

Introduction    👥 Overview    👤 Breakage Detector    Conclusion / Discussions

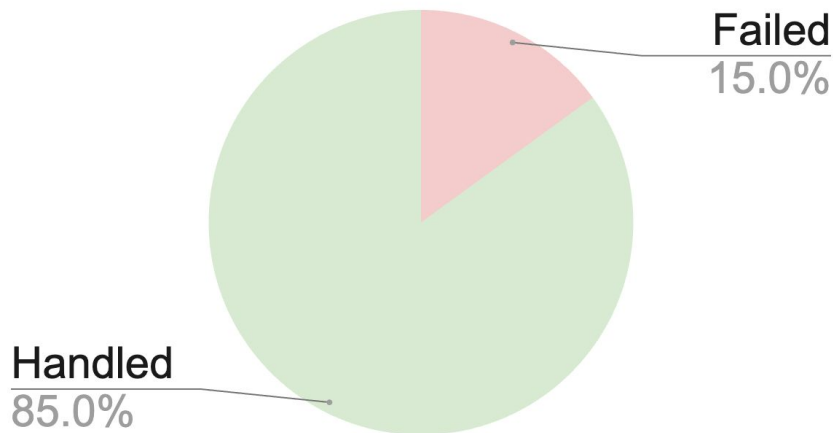# Breakage Reconstruction Evaluation

Our breakage reconstruction is accurate and reliable.
We manually look at 40 user reports of web breakages.

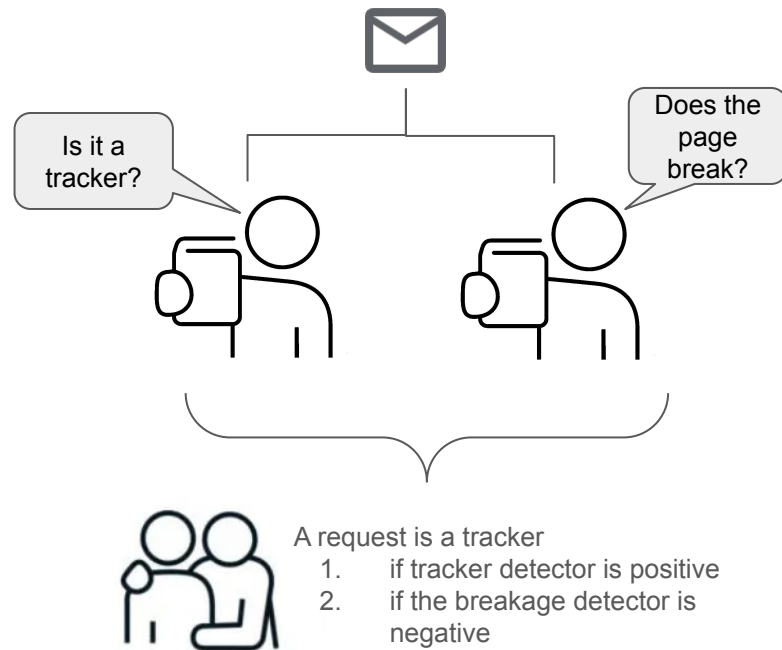- Compare the reconstructed breakage to user reports

Duumviri

Failed
0.0%

Handled
100.0%

Previous work

Failed
15.0%

Handled
85.0%

Introduction    👥 Overview    👤 Breakage Detector    Conclusion / Discussions

# Tracker Detection Evaluation

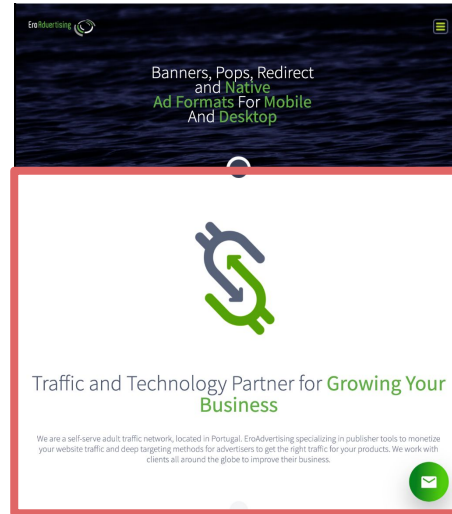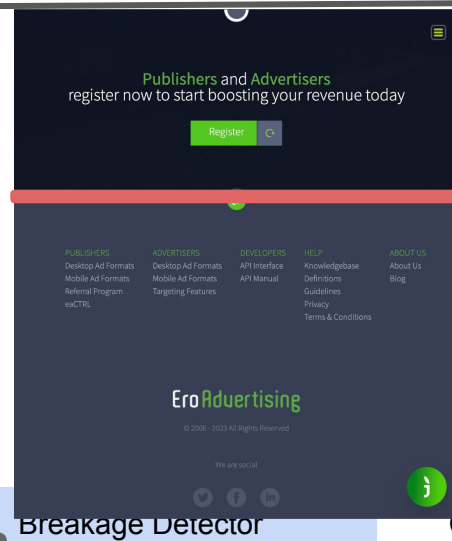**A large scale tracker evaluation** on 15K pages
- Accuracy 96.53% compared to filter lists
- Disagreement analysis shows
    - Duumviri is correct 55% of cases where Duumviri is positive

Is it a tracker?

Does the page break?

A request is a tracker
1. if tracker detector is positive
2. if the breakage detector is negative

# Tracker Detection Evaluation

**A large scale tracker evaluation** on 15K pages
- Accuracy 96.53% compared to filter lists
- Disagreement analysis shows
  - Duumviri is correct 55% of cases where Duumviri is positive
  - EasyPrivacy-caused breakages

**Unbroken page**

Web page body

**Broken page**

Web page body is missing

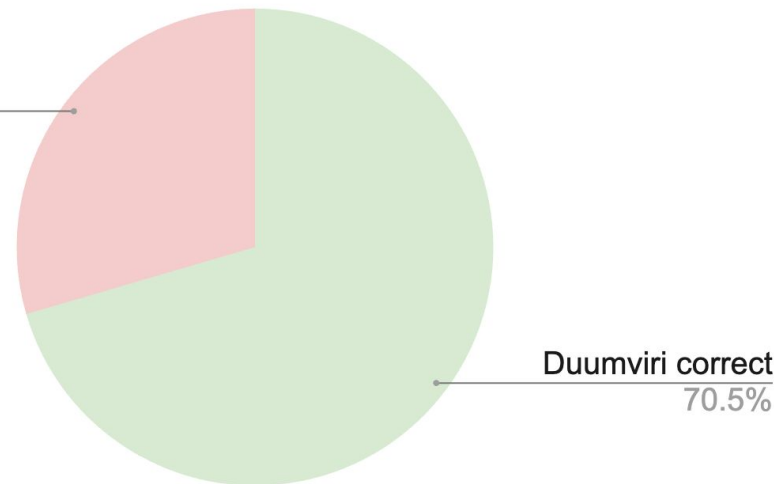https://github.com/easylist/easylist/issues/17829

# Comparison with previous work

Comparison with [AdGraph SP 2020]

- Similar accuracy using filter lists as ground truth
- Duumviri is correct in majority of times

| | AdGraph | **Duumviri** |
|---|---|---|
| Accuracy | 93.51 | 93.85 |
| Precision | 89.46 | 88.97 |
| Recall | 67.74 | 83.13 |
| AuROC | 0.9669 | 0.9682 |

Disagreement analysis



AdGraph correct
29.5%

Duumviri correct
70.5%

# Duumviri

- Proposes a two-modelled approach for tracker detection
    - A tracker detector
    - A breakage detector
        - Trained from exception rules
- Detects mixed trackers automatically

Code, data, models:

[github.com/dlgroupuoft/Duumviri-NDSS25](github.com/dlgroupuoft/Duumviri-NDSS25)

**UNIVERSITY OF TORONTO**

**Carleton University**

Artifact Evaluated
**NDSS**
SYMPOSIUM
Available
Functional
Reproduced