# Anonymity & Autonomy:
# Evidence from Google's 2012 Privacy Policy Change

James C. Cooper*

Antonin Scalia Law School

George Mason University

**PRIVACY**CON

# Anonymity and Autonomy

- Anonymity as an important dimension of privacy:
  - Engage in the world without actions being traced to identity

- Autonomy:
  - Zone to make private decisions free from observation or interference

- Reduction in ability to remain anonymous can reduce autonomy
- Harms:
  - Dignity
  - Personal development
  - Society
  - Privacy protective behavior

# Hypothesis

- 2012 Google announced it would combine data across platforms

- At the margin, this increased view into one's life will deter engagement in search behavior that one may want to keep private

- Measuring a reduction in autonomy due to loss in anonymity

PRIVACYCON

# Hypothesis

- Did people know and care?

- Examin... ...rch as a result of the 201...
- Google...
- "Differe...
  - Se...
  - No...
  - Lo... ...ore and after the cha...
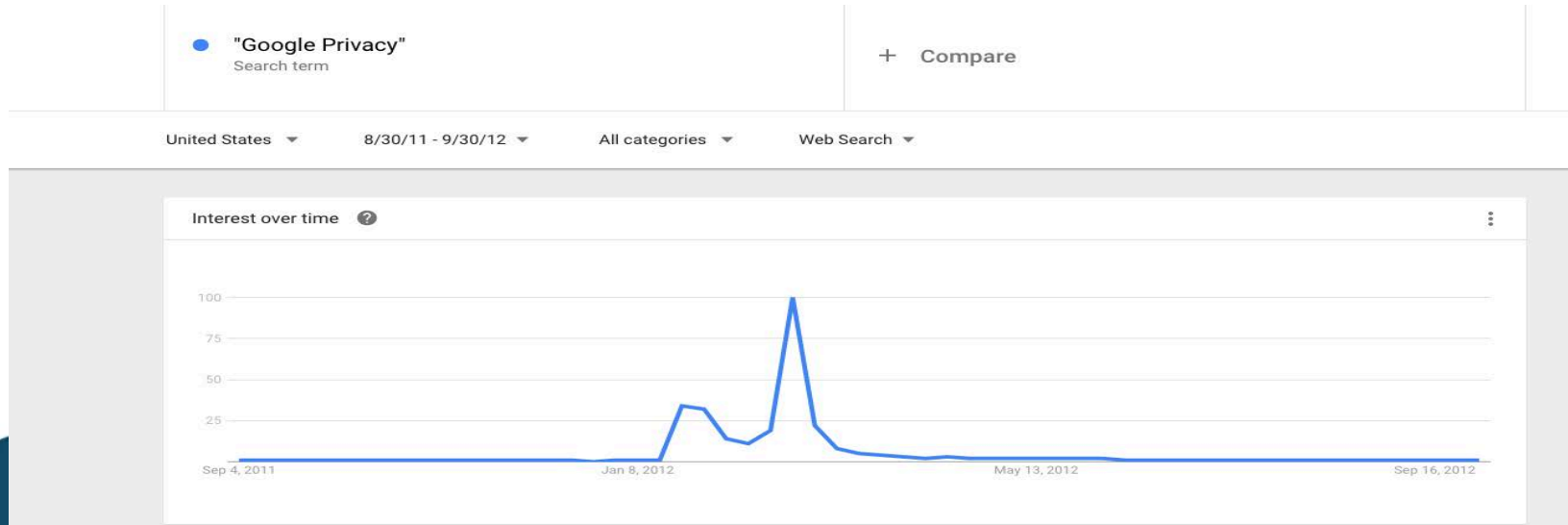- If mingl... ...see an increase in the c...

TABLE A1
SENSITIVE AND NON-SENSITIVE SEARCH TERMS

| Sensitive Terms | Average Trends Score (Jan.1, 2011 – Dec. 31, 2013) | Non-Sensitive Terms | Average Trends Score (Jan.1, 2011 – Dec. 31, 2013) |
|---|---|---|---|
| Abortion | 49.1 | Amazon | 53.0 |
| Acne | 69.5 | Apple | 44.6 |
| Adultery | 37.9 | Calculator | 77.7 |
| AIDS | 48.0 | CNN | 30.3 |
| Bankruptcy | 54.8 | Craigslist | 73.7 |
| Coming out | 54.7 | Ebay | 81.0 |
| Depression | 66.8 | Espn | 56.3 |
| Divorce | 29.0 | Facebook | 77.1 |
| Erectile Dysfunction | 49.4 | Games | 63.0 |
| Escort | 65.6 | Google | 69.0 |
| Gay | 60.3 | Iphone | 37.4 |
| Herpes | 64.0 | Mail | 83.5 |
| HIV | 42.6 | Maps | 71.6 |
| KKK | 36.7 | Netflix | 58.9 |
| Liposuction | 45.5 | News | 53.8 |
| Porn | 85.8 | Obama | 11.7 |
| Sexual Addiction | 36.6 | Target | 40.6 |
| Strip Club | 54.7 | Walmart | 37.1 |
| Suicide | 49.1 | Weather | 39.7 |
| Therapist | 66.3 | Yahoo | 84.3 |
| White power | 43.2 | Youtube | 75.9 |
| Total | 52.8 | | 58.1 |

# Results



Figure 2
Mean GT Scores: Before & After Google Policy Change

# Results

- Regression analysis
  - Setup:
    - Unit of observation:  GT score for search $i$, during week $t$, in state $j$
    - 13k – 108k observations, depending on window
    - Week, term, and state effects in all specifications
  - Main Findings:
    - 5-10% reduction in sensitive search with +/-1 and +/- 3 month windows
    - No measurable impact with +/- 6 month window
    - Impact does not vary by state-level privacy demand

# Results

- Results robust to different mixes of sensitive terms
- Falsification check for 2011 fails, but check for 2013 OK



Figure 3
Random Sample Estimates: +/-6 Month Window



Figure 5
Random Sample Estimates: +/-1 Month WIndow

# Conclusion

- Change may have induced a small drop in sensitive search, but this faded quickly

- Limitations/Future Work:
  - Trends not volume
  - Don't have universe of sensitive search
  - Other unmeasured margins may be more important:
    - E.g., content in Gmail; viewing in YouTube

# Insights from a 1-million-site Measurement of Online Tracking

**Steven Englehardt**
**@s_englehardt**

Dillon Reisman
@dillonthehuman

Arvind Narayanan
@random_walker

PRIVACYCON

# Visiting 2 websites results in 84 third parties contacted

# Open Web Privacy Measurement (OpenWPM)



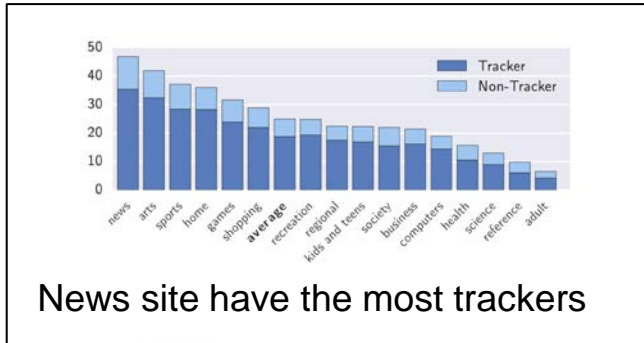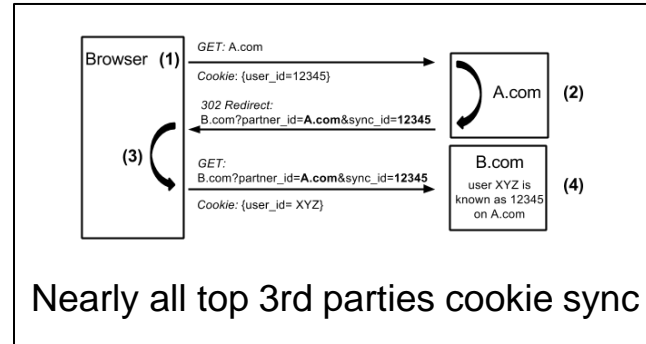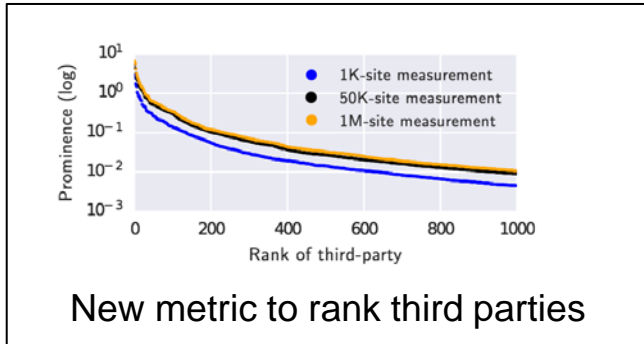https://github.com/citp/OpenWPM

# The Princeton Web Census

Monthly
1 Million Site Crawl

Collecting:

- Javascript Calls
- All javascript files
- HTTP Requests and Responses
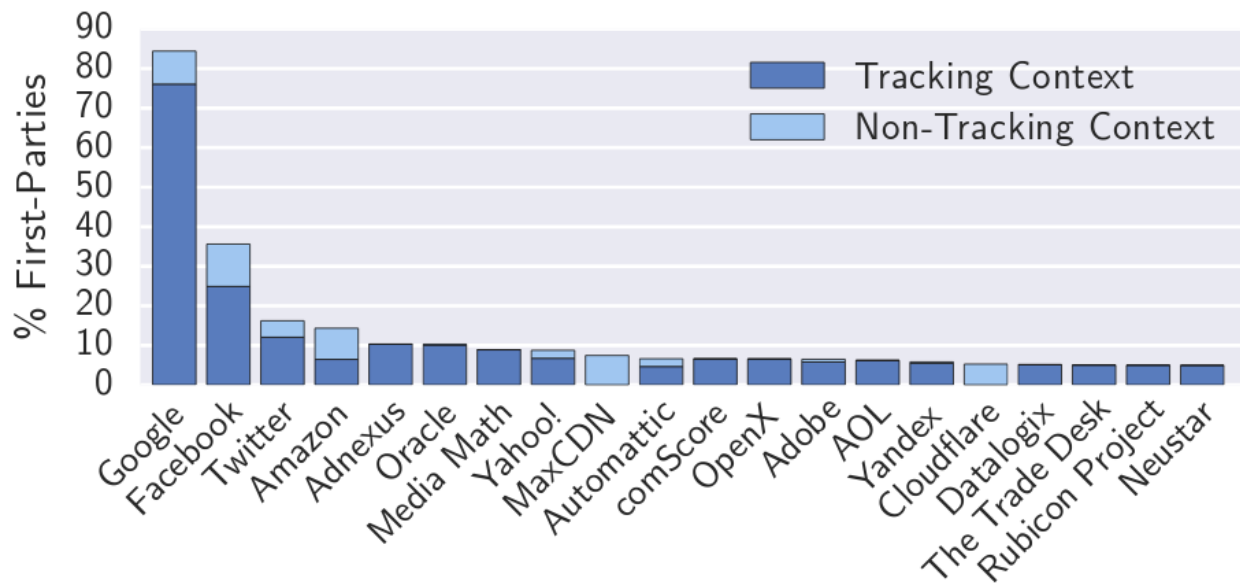- Storage (cookies, Flash, etc)
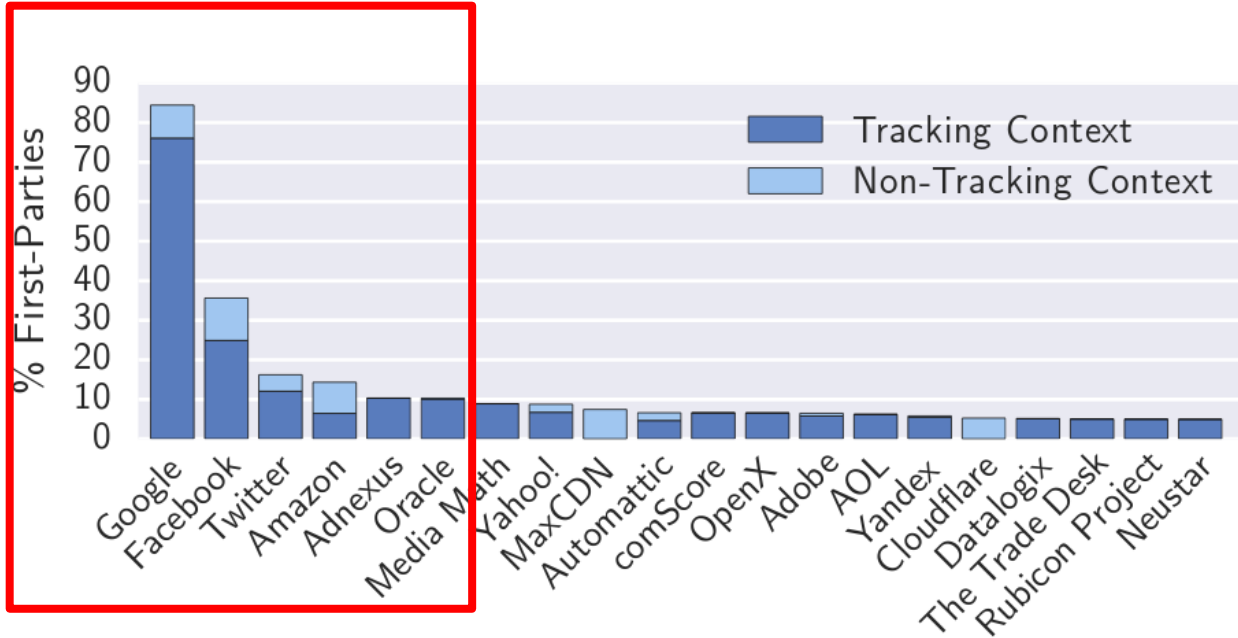
# Results of the Princeton Web Census



New metric to rank third parties



Nearly all top 3rd parties cookie sync



News site have the most trackers



Tracking protection misses less popular 3rd parties and techniques

https://webtransparency.cs.princeton.edu/webcensus/
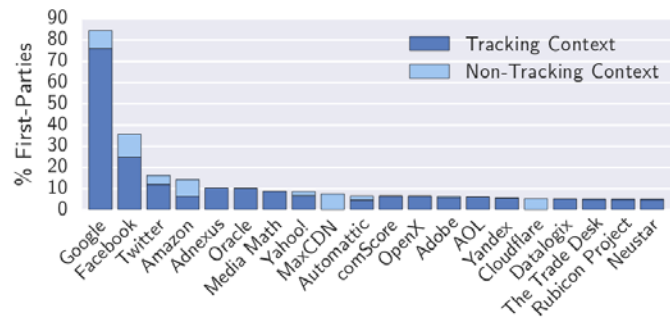
PRIVACYCON

# Consolidation of top trackers

# Only 6 organizations are present on >10% of sites

# Takeaways of consolidation

(1) Enforcement efforts can target large players, proactively set tracking norms.

(2) Large trackers can quickly deploy technique to a massive number of sites.

(3) Acquisitions can quickly shift tracking capability

# Trackers Impede HTTPS Adoption

# Trackers Impede HTTPS Adoption



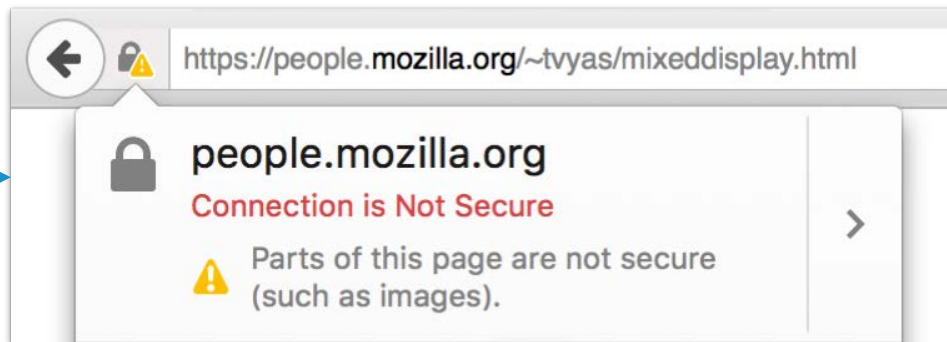**Mixed content downgrades security indicator!**
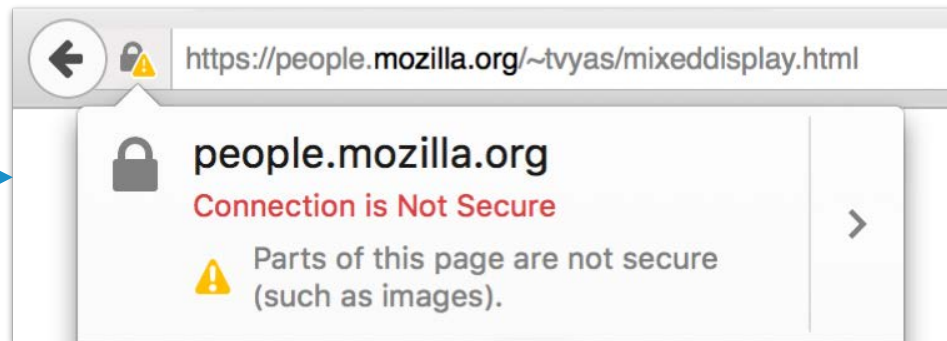
# Trackers Impede HTTPS Adoption



Firefox 47 ⓘ 🔒 | https://w

Chrome 47 🔒 https://w

https://people.**mozilla**.org/~tvyas/mixeddisplay.html

🔒 **people.mozilla.org**
**Connection is Not Secure**
⚠ Parts of this page are not secure (such as images).

Of sites with mixed content:
   half is caused solely by third parties (10% by trackers)

# Trackers Impede HTTPS Adoption



Firefox 47 ⓘ 🔒 | https://w

Chrome 47 🔒 https://w

https://people.**mozilla**.org/~tvyas/mixeddisplay.html

🔒 people.mozilla.org
**Connection is Not Secure**

⚠ Parts of this page are not secure
(such as images).

Of sites with mixed content:
> half is caused solely by third parties (10% by trackers)

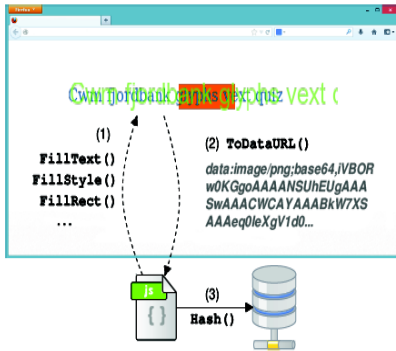Half of all third-parties are HTTP-only

# Takeaway: Tracking may have second-order privacy impacts

- (1) Slow the adoption of encryption
- (2) Identifier leakage in requests to
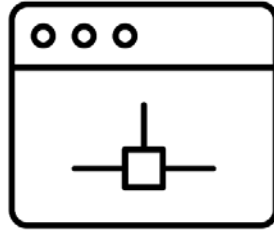- (3) Can aid network surveillance efforts
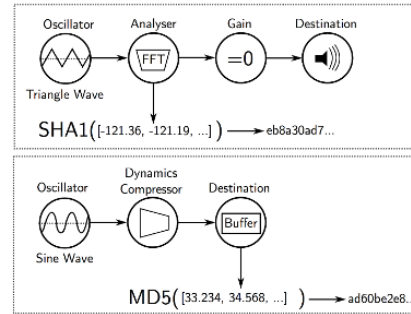
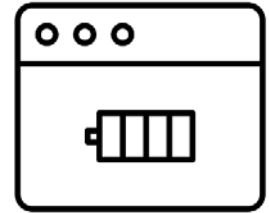# New Browser Features Used for Fingerprinting
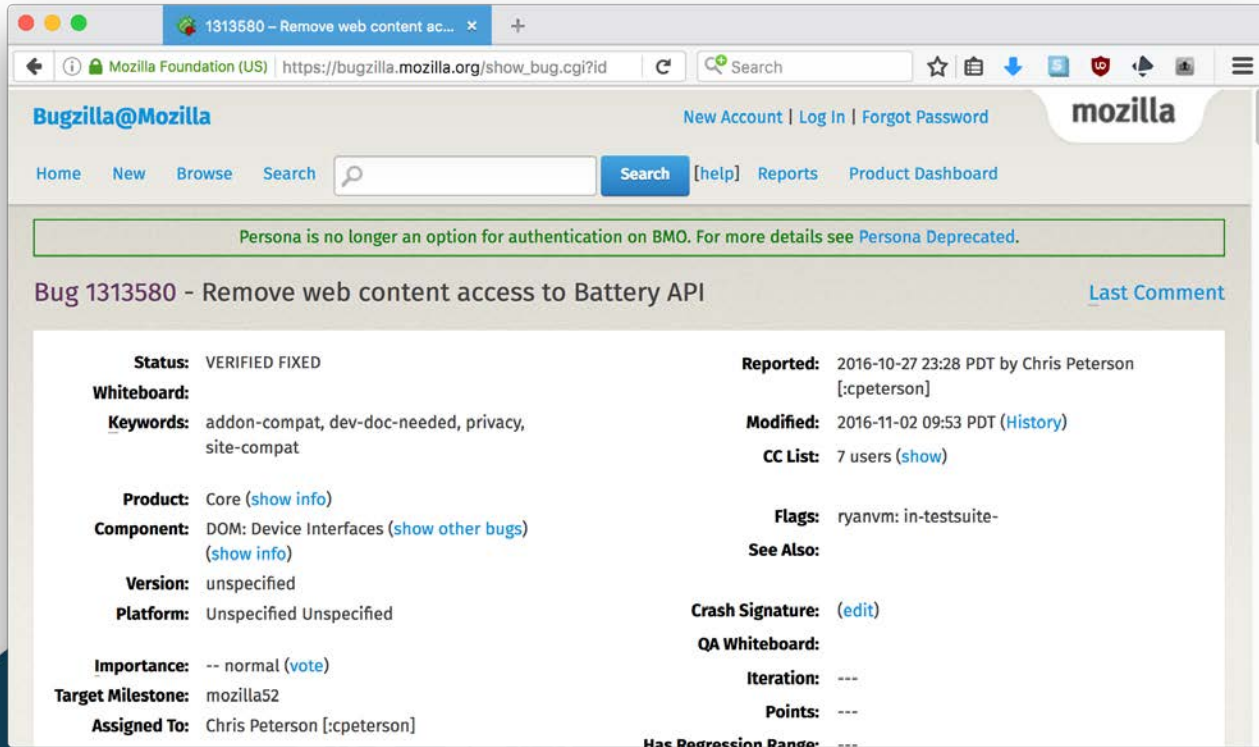
**Canvas**  **WebRTC**  **Audio**  **Battery**



https://webtransparency.cs.princeton.edu/webcensus/

PRIVACYCON

# Browsers remove BatteryStatus API citing privacy

# Browsers remove BatteryStatus API citing privacy



Browser window 1 — Bugzilla@Mozilla

- URL: https://bugzilla.mozilla.org/show_bug.cgi?id
- Tab: 1313580 – Remove web content ac...
- Home | New | Browse | Search

Persona is no longer an option for a...

**Bug 1313580 - Remove web content acces...**

- **Status:** VERIFIED FIXED
- **Whiteboard:**
- **Keywords:** addon-compat, dev-doc-needed, privacy site-compat
- **Product:** Core (show info)
- **Component:** DOM: Device Interfaces (show other bugs) (show info)
- **Version:** unspecified
- **Platform:** Unspecified Unspecified
- **Importance:** -- normal (vote)
- **Target Milestone:** mozilla52
- **Assigned To:** Chris Peterson [:cpeterson]

Browser window 2 — WebKit Bugzilla

- URL: https://bugs.webkit.org/show_bug.cgi?id=164213
- Tab: Bug 164213 – Remove Battery Stat...
- Home | New | Browse | Search | [?] | Reports | Requests | Help | New Account | Log In | Forgot Password

|« First Last »| « Prev Next »   This bug is not in your last search results.

**Bug 164213 - Remove Battery Status API from the tree**

Bug 164213: Remove Battery Status API from the tree

- **Status:** RESOLVED FIXED
- **Product:** WebKit
- **Component:** WebKit Misc.
- **Version:** WebKit Nightly Build
- **Platform:** Unspecified Unspecified
- **Importance:** P2 Normal
- **Assigned To:** Alex Christensen
- **URL:**
- **Keywords:**
- **Depends on:**
- **Blocks:**
- **Reported:** 2016-10-30 20:26 PDT by Brady Eidson
- **Modified:** 2016-11-02 14:32 PDT (History)
- **CC List:** 8 users (show)
- **See Also:** 129040

Show dependency tree / graph

PRIVACYCON

# Takeaway: Expect any new API to be analyzed for its fingerprintability

1. Early detection of abuse can stem adoption
2. Browsers **view fingerprinting as abuse**
   a. Mitigate fingerprinting during standardization
   b. Remove APIs due to fingerprinting use

# Our data is available!

The data is available as bzipped PostgreSQL dumps. The schema file used in all of the datasets is available here.

| Dataset | Comments |
|---|---|
| 1 Million Site Stateless | Parallel Stateless Crawl |
| 100k Site Stateful | Parallel Stateful Crawl -- 10,000 site seed profile |
| 10k Site ID Detection (1) | Sequential Stateful Crawl -- Flash enabled -- Synced with ID Detection (2) |
| 10k Site ID Detection (2) | Sequential Stateful Crawl -- Flash enabled -- Synced with ID Detection (1) |
| 55k Site Stateless with cookie blocking | Parallel Stateless Crawl -- Firefox set to block all third-party cookies |
| 55k Site Stateless with Ghostery | Parallel Stateless Crawl -- Ghostery extension installed and set to block all possible trackers |
| 55k Site Stateless with HTTPS Everywhere | Parallel Stateless Crawl -- HTTPS Everywhere installed |

https://webtransparency.cs.princeton.edu/webcensus/index.html#data

PRIVACYCON

# Getting third-party responses from our data

```python
tp_query = "SELECT r.url, h.value FROM http_responses_view AS r " \
           "LEFT JOIN http_response_headers_view as h ON h.response_id = r.id " \
           " WHERE r.top_url LIKE %s AND " \
           "url not LIKE %s and h.name = 'Content-Type'"
cur = connection.cursor()
cur.itersize = 100000
try:
    top_ps = utils.get_host_plus_ps(top_url)
except AttributeError:
    print("Error while finding public suffix of %s" % top_url)
    return None

cur.execute(tp_query, (top_url, top_ps))

el_parser = BlockListParser('easylist.txt')
ep_parser = BlockListParser('easyprivacy.txt')
response_data = defaultdict(dict)

for url, content_type in cur:
    if utils.should_ignore(url):
        continue

    url_data = dict()

    url_ps = utils.get_host_plus_ps(url)
    if url_ps == top_ps:
        continue
    url_data['url_ps'] = url_ps

    is_js = utils.is_js(url, content_type)
    is_img = utils.is_img(url, content_type)
    is_el_tracker = utils.is_tracker(url,
                        is_js=is_js,
                        is_img=is_img,
                        first_party=top_url,
                        blocklist_parser=el_parser)
    is_ep_tracker = utils.is_tracker(url,
                        is_js=is_js,
                        is_img=is_img,
                        first_party=top_url,
                        blocklist_parser=ep_parser)
    is_tracker = is_el_tracker or is_ep_tracker

    url_data['is_js'] = is_js
    url_data['is_img'] = is_img
    url_data['is_tracker'] = is_tracker
    response_data[url] = url_data
```

```python
def is_js(url, content_type):
    if get_top_level_type(content_type) == 'script':
        return True
    if urlparse(url).path.split('.')[-1].lower() == 'js':
        return True
    return False

def is_img(url, content_type):
    if get_top_level_type(content_type) == 'image':
        return True
    extension = urlparse(url).path.split('.')[-1]
    if extension.lower() in IMAGE_TYPES:
        return True
    return False
```

```python
def get_host_plus_ps(url):
    """Strip the URL down to just a hostname+publicsuffix.

    If the provided url contains an IP address, the IP address is returned.
    """

    hostname = urlparse(url).hostname
    try:
        ip_address(hostname)
        return hostname
    except ValueError:
        return psl.get_public_suffix(hostname)
```

```python
def get_trackers(url_list, first_party, blocklist_parser=None, blocklist="easylist.txt"):
    """Identify domains that are identified as trackers from list of URLs.

    Returns set of domains/IPs filtered by the given blocklist_parser.
    TODO: Better to return set of domains/IPs, or list of filtered urls?
    """

    if not blocklist_parser:
        blocklist_parser = BlockListParser(blocklist)

    filtered_domains = set()
    for url in url_list:
        if is_tracker(url, first_party, blocklist_parser):
            filtered_domains.add(get_host_plus_ps(url))

    return filtered_domains
```

# Getting third-party responses from our data

# Getting third-party responses with Census.py

```
census.get_third_party_responses_by_domain(
    database_connection,
    "http://nytimes.com"
)
```

# Getting third-party responses with Census.py

- `get_third_party_responses_by_domain`
- `get_third_party_responses_by_domain`
- `get_cookie_syncs_on_domain`
- `is_tracker`
- `get_trackers`

**PRIVACY**CON

# Getting third-party responses with Census.py

- `get_third_party_responses_by_domain`
- `get_third_party_responses_by_domain`
- `get_cookie_syncs_on_domain`
- `is_tracker`
- `get_trackers`

Contact us for access to "alpha" analysis server and library!

# Thanks for listening!

**Full Paper:**
senglehardt.com/papers/ccs16_online_tracking.pdf

**Data and Analysis:**
webtransparency.cs.princeton.edu/webcensus/

**Collaborate:**
webtap.princeton.edu/research/

**Contact Me**

**Email:** ste@cs.princeton.edu

**Twitter:** @s_englehardt

**Web:** senglehardt.com

Image Assets from the Noun Project:
Browser Network and Browser Battery by Aybige

**PRIVACY**CON

# Detection and Circumvention of Anti Ad-Blockers:
# A New Arms Race on the Web

Zubair Shafiq (University of Iowa)
Zhiyun Qian (UC-Riverside)

**PRIVACY**CON

PRIVACYCON

# (Self) Regulation

FTC

AdChoices

Do Not Track

FCC

# Privacy Preserving Tools

- EFF's Privacy Badger
  - Target publishers that do not respect DNT

- Ghostery
  - Proprietary

- Ad-block (open-source, public filter lists)
  - Adblock Plus
  - uBlock Origin

# Popularity of Ad-blockers



GLOBAL MONTHLY ACTIVE USERS
(mobile adblocking browsers)

198M

237M

270M

334M

376M

Jan 2015   April 2015   July 2015   Oct 2015   Jan 2016

PageFair
in partnership with PRIORI DATA

PRIVACYCON

# But…

## Here's The Thing With Ad Blockers

**We get it:** Ads aren't what you're here for. But ads help us keep the lights on. So, add us to your ad blocker's whitelist or pay $1 per week for an ad-free version of WIRED. Either way, you are supporting our journalism. We'd really appreciate it.

**Sign Up**

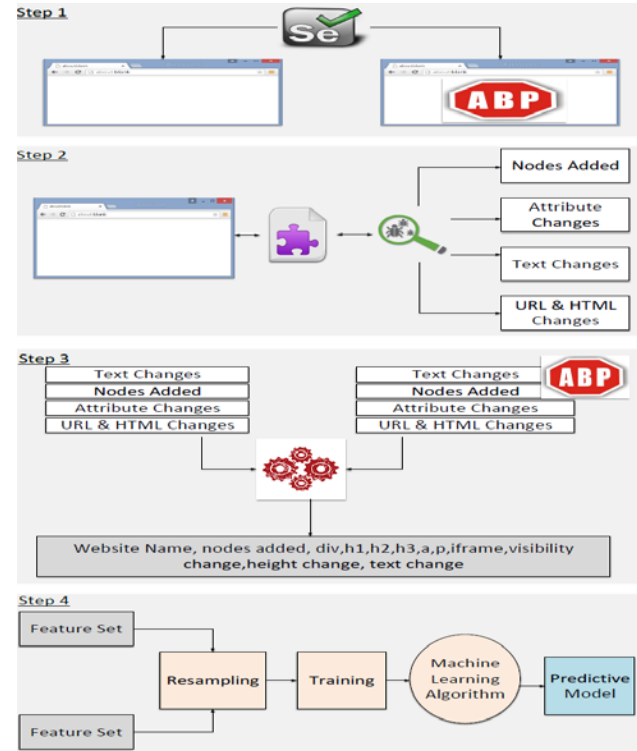Already a member? Log in

# Project Goals

- Measure Anti Ad-blocking in the Wild
  - Number of anti ad-blocking publishers
  - Third-party anti ad-blocking services
  - Ad-block detection techniques

- Develop a Stealthy Ad-blocker
  - Automatically block anti ad-block scripts

# Measuring Anti Ad-blockers

- Crawl Alexa top 100K websites
- A/B testing
  - with and without ad-block
- Feature extraction
  - nodes, attributes, text
- Machine learning models
  - Random forest, SVM, Bayesian
- Results
  - 95% precision, 93% recall
  - 1100 websites use anti ad-blockers

# How Anti Ad-blockers Work?

Websites employ anti
ad-block scripts

– Identify leaked extension
information

– Verify ads (active or
passive)

```javascript
//step 1: set timeout
var myVar = setInterval(function() {
  myFunc()
}, 2000);

function myFunc() {

  // step 2: condition check
  if (window.iExist === undefined ||
    (!$("#XUinXYCfBvqpyDHOrOAVClxoWJemrlPpfYCdWfiyAzNY").is(
      ":visible") && (($(".vip_052x003").height() < 100 && !$(
      "#vipchat").length) && $(".vip_09x827").height() < 25))) {

    //step 3: response
    $("#XUinXYCfBvqpyDHOrOAVClxoWJemrlPpfYCdWfiyAzNY").css(
      "width:100%;height:100%;position:fixed;z-index:999999;top:0");
    $("#XUinXYCfBvqpyDHOrOAVClxoWJemrlPpfYCdWfiyAzNY").show();
  }
  else if ($("#XUinXYCfBvqpyDHOrOAVClxoWJemrlPpfYCdWfiyAzNY").is(
    ":visible") && $(".vip_052x003").height() > 249) {

    $("#XUinXYCfBvqpyDHOrOAVClxoWJemrlPpfYCdWfiyAzNY").hide()
  }
}
```

# Towards a Stealthy Ad-blocker

- Remove anti ad-block scripts through filter lists
- Crowd-sourced, manually populated

# Takeaway

- Users, Society, Economics

- "Ad blocking is like garlic.
  You hang it on your door to keep Dracula away from sucking your blood.
  Ad blocking is not the enemy.
  Just stop being Dracula."

  – Doc Searls