# No boundaries: Data exfiltration by third-party tracking scripts

*Steven Englehardt*

<u>Joint work with:</u>
Güneş Acar, Jeffrey Han, and Arvind Narayanan
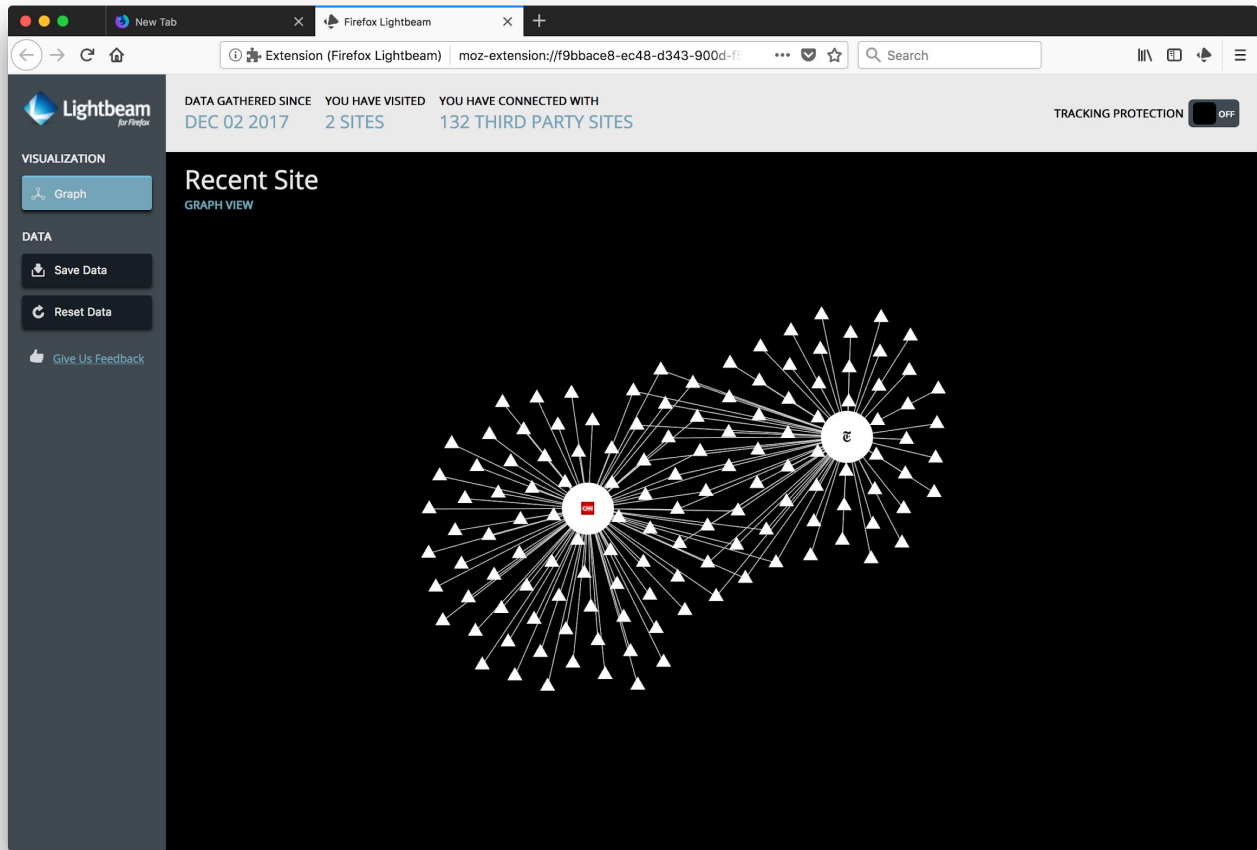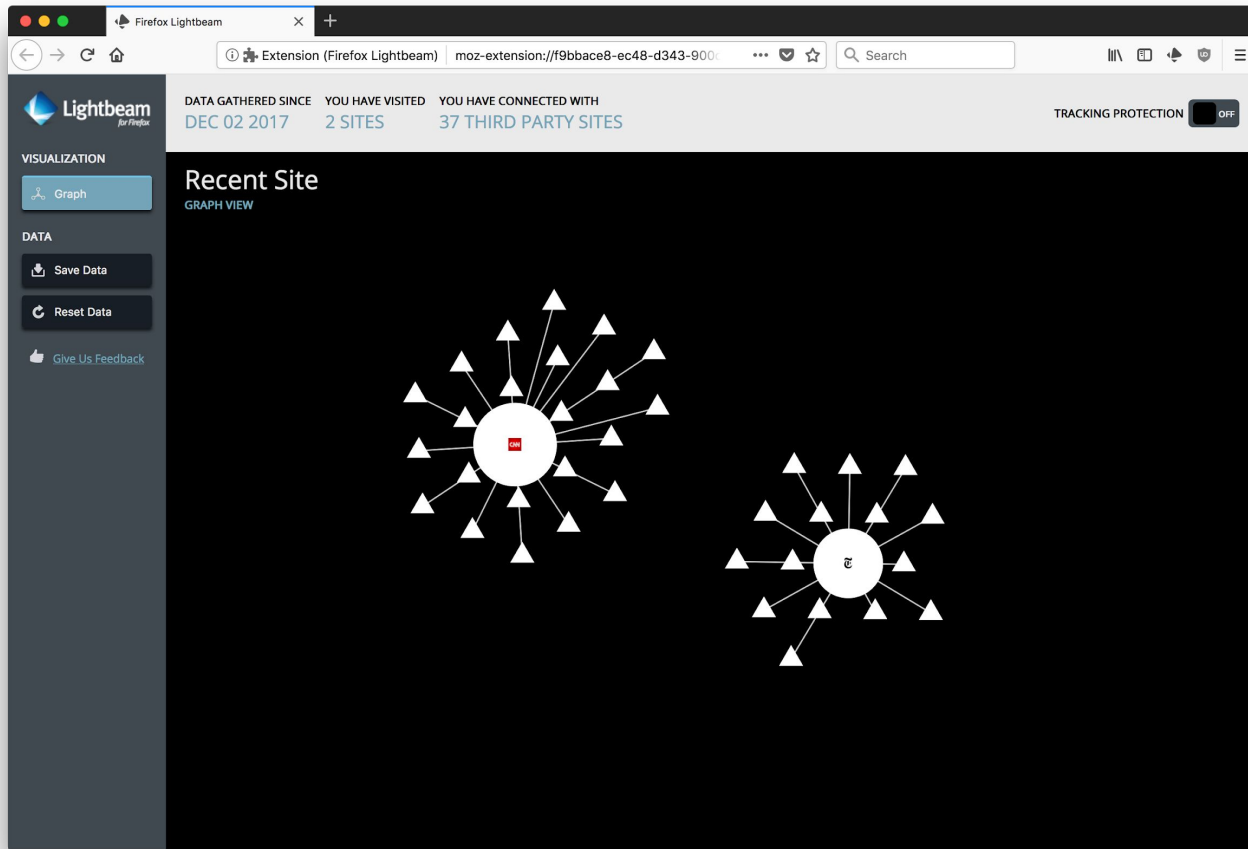
**I'm now at...**

moz://a

PRINCETON UNIVERSITY

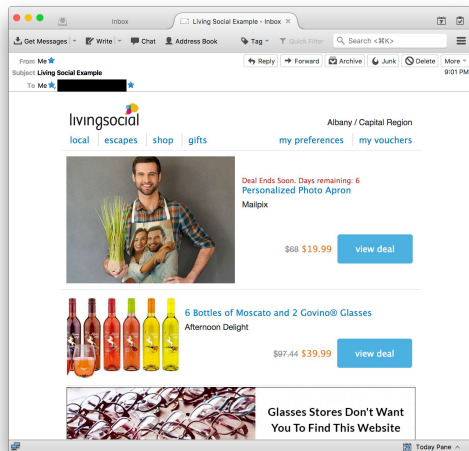CENTER FOR INFORMATION TECHNOLOGY POLICY
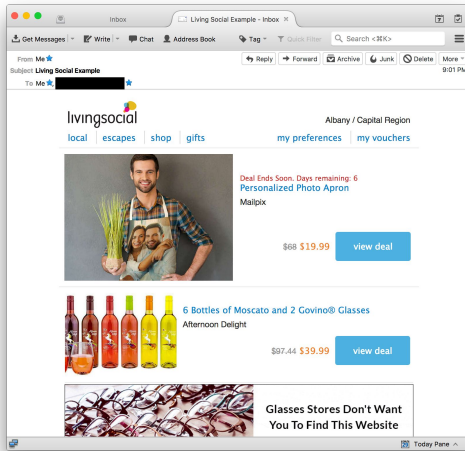PRINCETON UNIVERSITY

Just **two page visits** cause requests to **132 distinct hostnames**.

With uBlock Origin enabled, the number of hostnames requested is down to 37. **Nearly 100 of the hosts loaded were ads, trackers, and analytics.**

# What happens when you load remote content in an email?

Your device contacts 24 companies
→ 20 can track you (if supported)
→ 10 receive an MD5 hash of your email address

**Sets a Cookie**

OpenX (openx.net)
comScore (scorecardresearch.com, voicefive.com)
Oracle (bluekai.com)
Google (doubleclick.net)
Realtime Targeting Aps (mojn.com)

MediaMath (mathtag.com)
TapAd (tapad.com)
IPONWEB (bidswitch.net)
AOL (advertising.com)
Centro (sitescout.com)
The Trade Desk (adsrvr.org)
Adobe (demdex.net)

**Receives MD5(email address) & Sets a Cookie**

American List Counsel (alcmpn.com)
LiveIntent (liadm.com)
Oracle (nexac.com)
Acxiom (rlcdn.com, pippio.com, acxiom-online.com)
Criteo (criteo.com)
Conversant Media (dotomi.com)
V12 Data (v12group.com)
VideoAmp (videoamp.com)
<Unknown> (alocdn.com)

**Receives MD5(email addr.)**

Criteo (emailretargeting.com)
Neustar (agkn.com)

**Receives Bare Request**

LiveIntent (licasd.com)
Google (2mdn.net)
Akamai (akamai.net)

*Englehardt, Han, and Narayanan, "I never signed up for this! Privacy implications of email tracking" (PETS 2018)*

# A user's email address is the perfect identifier!

- It's unique

- It rarely changes

- It's the same across devices

- Consumers freely provide it to stores

- There's a lot of associated data

**PII-based tracking**

```
UUID = {
  MD5(bob@example.com),
  SHA1(bob@example.com),
  SHA256(bob@example.com)
}
```

Are trackers also collecting PII on the web?

# We can use web crawls to detect PII collection:

- **Crawl 50K sites with OpenWPM**
  - main page and 5 inner pages
- **Monitor access to PII sources**
  - Autofilled credentials
  - Mutation events to monitor form insertion
  - HTMLInputElement instrumentation to intercept access to form input fields
- **Search for PII in network traffic**
  - Request and response headers
  - POST payloads

OpenWPM: https://github.com/citp/OpenWPM

# Challenge: Measurements require the automated submission of PII to sites

## Mailing list sign-ups

**Email Address**

*

**Birthdate**

📅 ( mm / dd / yyyy ) *

**Your Country / Territory**

United States of America ⌄ *

**State**

Not Specified ⌄

**Zip Code**

**Your Gender**

Prefer not to say ⌄

☐ By checking this box you agree to the TaylorSwift.com Terms of Use and Privacy Policy.

**Subscribe**

## Login Forms

**Sign in**

Email address:

Password:

I forgot my password.

**SIGN IN** **Cancel**

# Injecting PII into the web: bait technique

# Third parties collect PII for tracking

**Autofill abuse**

**Social Login**

**Session Recording**

# Login manager abuse for web tracking

*Acar, Englehardt and Narayanan, "No boundaries for user identities: Web trackers exploit browser login managers" (Freedom to Tinker)*

# Built-in login managers

- Remembers username & passwords (opt-in)

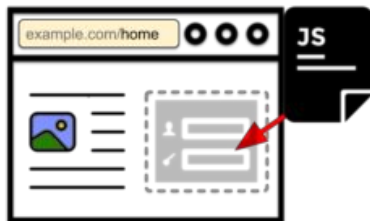- Autofills login forms

- Different than CC and address autofill

**User submits a login or registration form, clicks "Save" to store the credentials.**

example.com/**login**

**Third-party script is not present** on the login page

Would you like Firefox to save this login for princeton.edu?

username@example.com

●●●●●●●

☐ Show password

Don't Save | Save

---

**User visits a non-login page on the same site; this time the third party script is present**

example.com/**home** | JS

example.com/**home** | username@p... | JS

example.com/**home** | username@p... | JS

- MD5(email)
- SHA1(email)
- SHA256(email)

*1.* Third-party script injects an invisible login form

*2.* Login manager fills in user's email and password

*3.* The script reads the email address from the form and sends it hashes to third-party servers

# Findings

| Company | Script address | No of sites |
|---------|----------------|-------------|
| Adthink | https://static.audienceinsights.net/t.js | 1047 |
| OnAudience | http://api.behavioralengine.com/scripts/be-init.js | 63 |

# Social Login abuse for web tracking



Login with Facebook

Firstparty.com will receive: your public profile and email address

Continue

Englehardt, Acar, and Narayanan, "No boundaries for Facebook data: third-party trackers abuse Facebook Login" (Freedom to Tinker)

# Findings

| Company | Script Address | Facebook Data Collected |
|---------|----------------|-------------------------|
| OnAudience* | http://api.behavioralengine.com/scripts/be-init.js | User ID (hashed), Email (hashed), Gender |
| Augur | https://cdn.augur.io/augur.min.js | Email, Username |

...as well as several others grabbing user ID

# Session recording scripts scoop up sensitive information

# What are session recording scripts?

- Session recording scripts create a "video" of all of a user's actions on a site.

    ○ Key presses

    ○ Mouse clicks, mouse movements

    ○ Scrolling behavior...

- Publishers can later review the videos.

# Why use session recording scripts?



Answer questions like:

- Who are my most valuable customers?

- Who added items to the cart but didn't convert?

- Where do users leave the onboarding flow?

- Where are users frustrated?

**The problem:** recordings require **a ton** of data

Full page source and text

Mouse movements & clicks

Keypresses

# Companies support redaction



**Easily protect your user's privacy.**

Exclude sensitive customer data from ever leaving your customer's browser by using our in-app point and click system.

# How can things go wrong?

# Redactions miss sensitive information

- Name
- CC #
- CVV

# Walgreens misses fields during redaction



Walgreens makes thorough use of redaction

# Walgreens misses fields during redaction



But prescription information is missed!

(the user's full name was not redacted on the previous page)

Walgreens makes thorough use of redaction

# Session recordings are widespread

- 14+ analytics company offer recording services

    - Present on 99,174 of the top 1 million sites


- Evidence of recording on 7,918 sites.

    - Likely a lower bound as recording scripts sample users

Session recording present on ~1 - 10% of the top 1 million sites. We found several severe PII leaks after manually reviewing ~30 sites.

→ **How many more leaks are out there?**

# What can we do?

**1. Just keep measuring?**

*Will public backlash be enough? (Probably not)*

**2. Try to plug holes in browsers?**

*Limit autofill? → Sure*

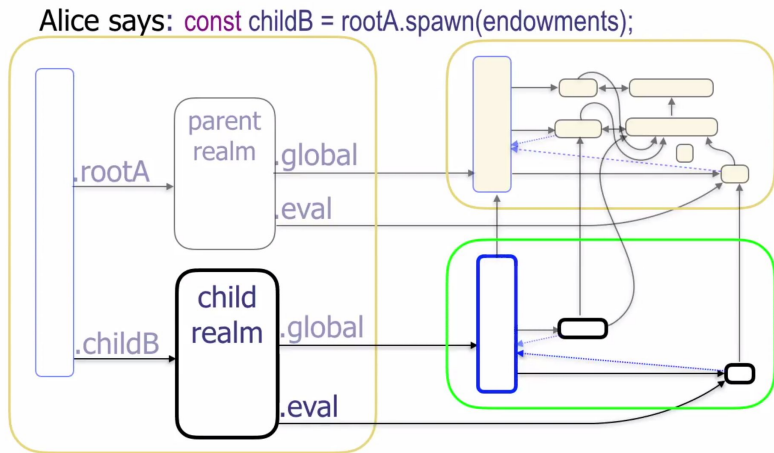*Limit social login sniffing? → How?*

**3. Push for regulation?**

*Hopeful in Europe, but what about the rest of the world?*

# **Possible direction:** Better JS confinement
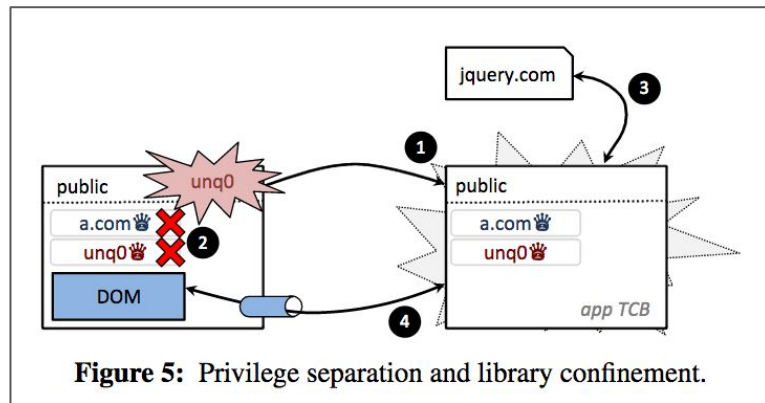
## Frozen Realms



(https://github.com/tc39/proposal-frozen-realms)

## COWL



**Figure 5:** Privilege separation and library confinement.

(https://www.usenix.org/node/186158)

→ Problem: Requires the cooperation of sites ←

# **Possible direction:** Better JS confinement



Insert the Javascript code directly on your website

Here's the code you need to put on your website. Copy and paste it into Google Tag Manager. Or you can paste it between the <head> and </head> tags on the pages you want to track visitors on.

```
<script type="text/javascript">
    window.smartlook||(function(d) {
    var o=smartlook=function(){ o.api.push(arguments)},h=d.getElementsByTagName('head')[0];
    var c=d.createElement('script');o.api=new Array();c.async=true;c.type='text/javascript';
    c.charset='utf-8';c.src='https://rec.smartlook.com/recorder.js';h.appendChild(c);
    })(document);
    smartlook('init', ██████████████████████);
</script>
```

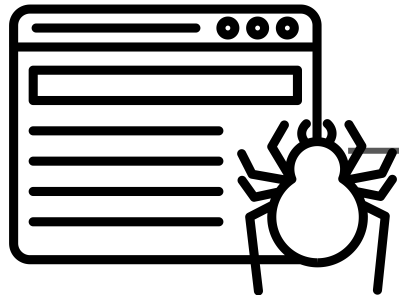COPY THE CODE    Or send it to your developer via email

→ Problem: Requires the cooperation of sites ←

# **Possible direction:** Measurement + Blocking

## Detect invasive scripts

Real users

Crawlers

## Build better blocklists

**Possible direction:** Measurement + Blocking

Problems:
- Broken sites
- Obfuscation
- User privacy concerns

# How can we stop invasive web tracking?

1. **Just keep measuring?**

2. **Try to plug holes in browsers?**

3. **Push for regulation?**

4. **Work on confinement solutions?**

5. **Detect and block invasive scripts?**

**Research:** https://freedom-to-tinker.com/tag/noboundaries/

**Me:** https://senglehardt.com | senglehardt@mozilla.com