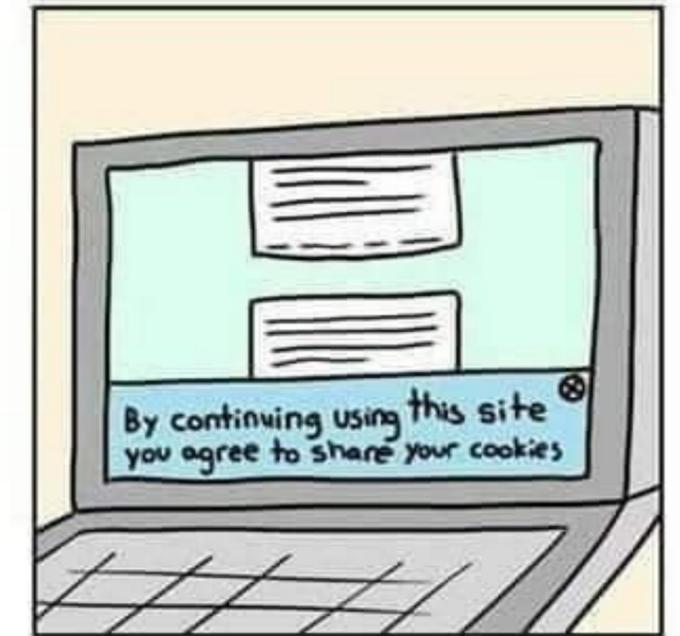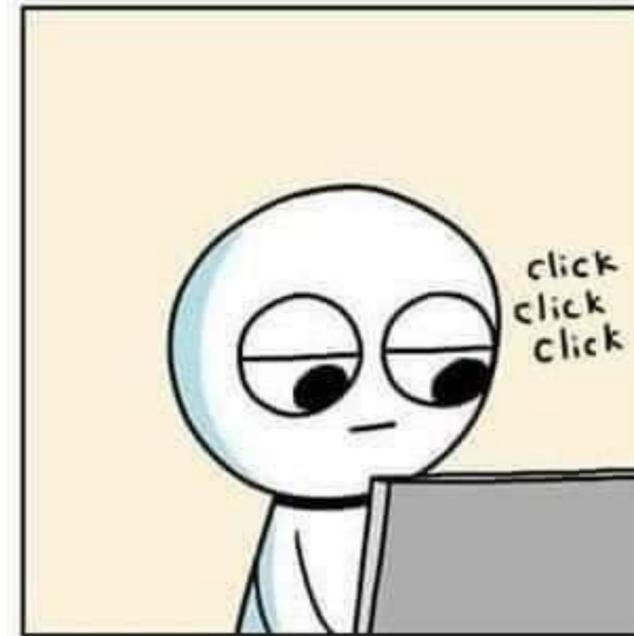# The Past, Present and Future of Online Tracking

**Dr Michael Veale**

Faculty of Laws, University College London

# A short history of Web tracking

- In the early 90s, the Web was 'stateless' — it had no *memory of its visitors*.

- **Cookies** were invented to solve this problem: they are simply **text placed on your browser by a web server that a server can look at later**.

```
Syntax of the Set-Cookie HTTP Response Header:
Set-Cookie: NAME=OPAQUE\_STRING \
    [; expires= ] \
    [; path=] \
    [; domain=] \
    [; secure]

Syntax of the Cookie HTTP Request Header:
Cookie: NAME=OPAQUE\_STRING *[; NAME=OPAQUE\_STRING]
```

First proposal for state management on the
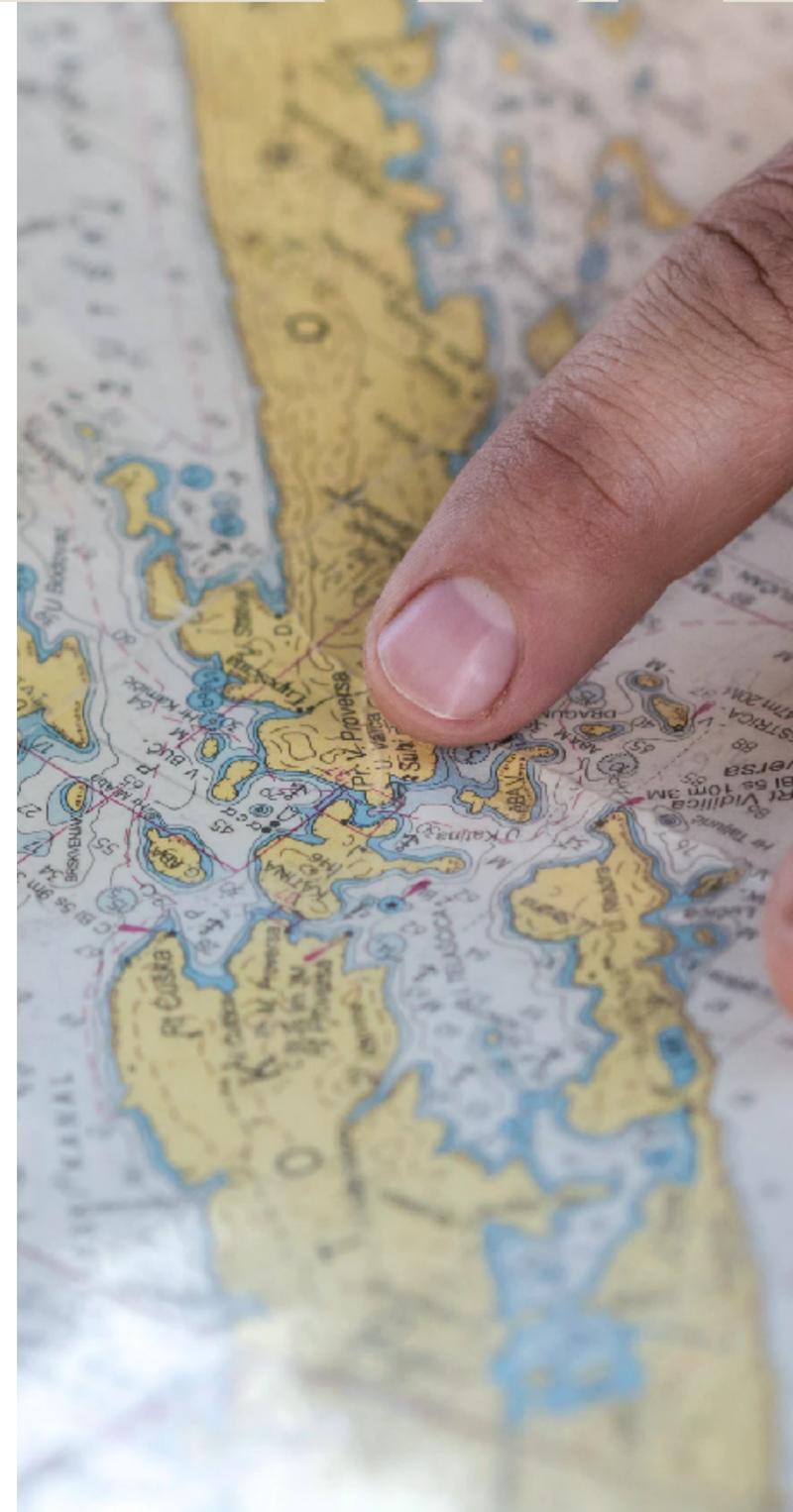web (Apr. 18, 1995)

# Webpage complexity grows

- In the early days of the Web, all content on a webpage came from the same **server**.

- An early, popular browser, **Netscape Navigator**, introduced the function of rendering two webpages in a single browsing window in 1996 (frames).

- This created a problem: could the second website access the cookies the first had laid?
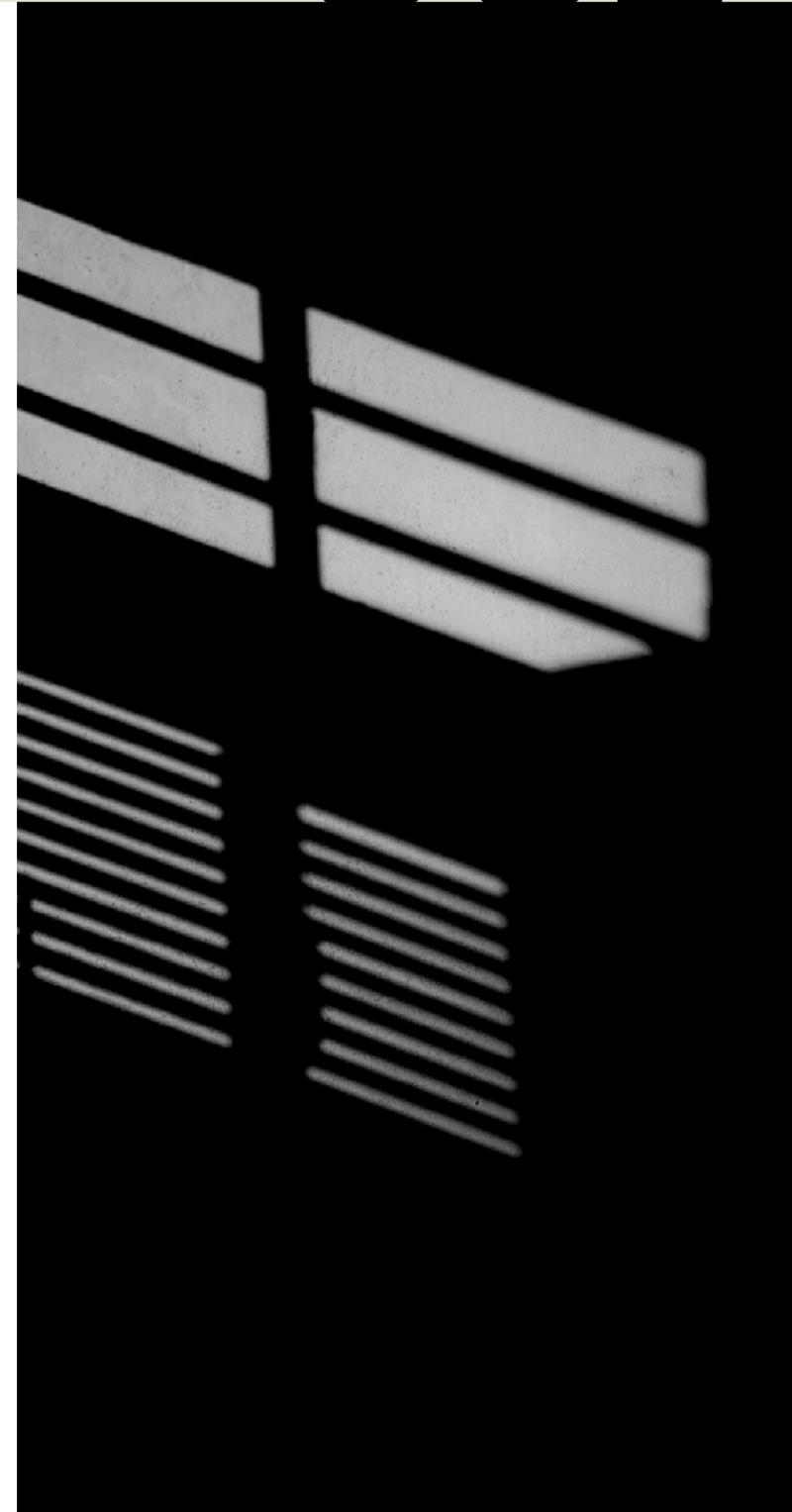
# The same origin policy

- The solution — the Same Origin Policy.

- Cookies only accessible by servers that share features (particularly the domain) of the one that laid them.

- A user visiting ucl.ac.uk should expect only ucl.ac.uk cookies to be read — not kcl.ac.uk cookies.

# Crafty workarounds



- Didn't fix the problem for long:
  - Websites started calling many distinct servers. Used to be 1, now 100s — because a website would instruct your computer to query many domains.
  - These many domains collaborate to share information about users' Web usage and more — called Cookie Syncing.

- Google calls home with unique identifiers for at least 28% of all web page loads, while Facebook does the same for approximately 15%. The proportion is significantly higher in certain sectors, such as news, compared to others, such as banking.*

- Collaboration between trackers means that even under conservative estimates, 53 firms observe more than 91% of users' browsing behaviour.**
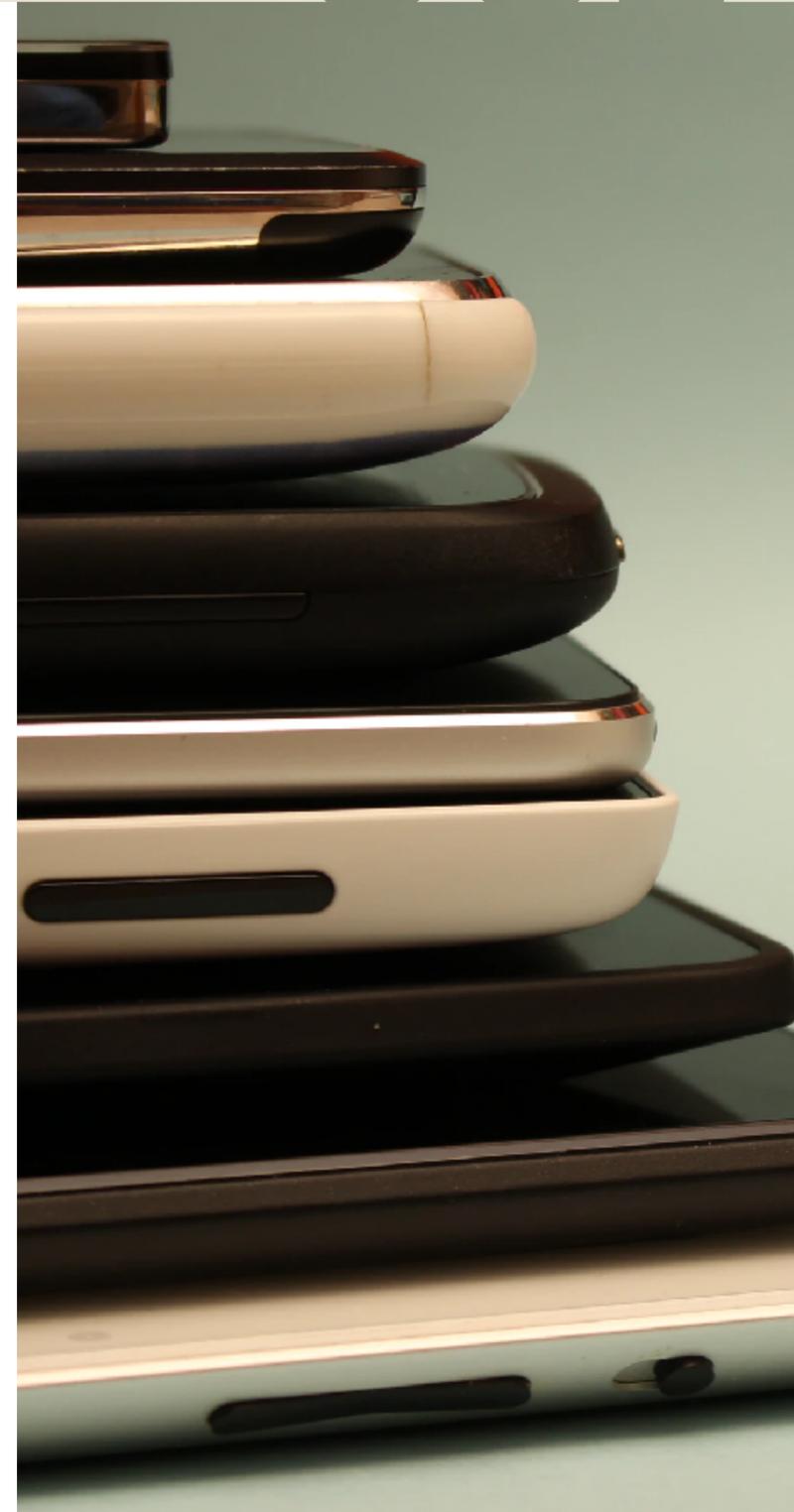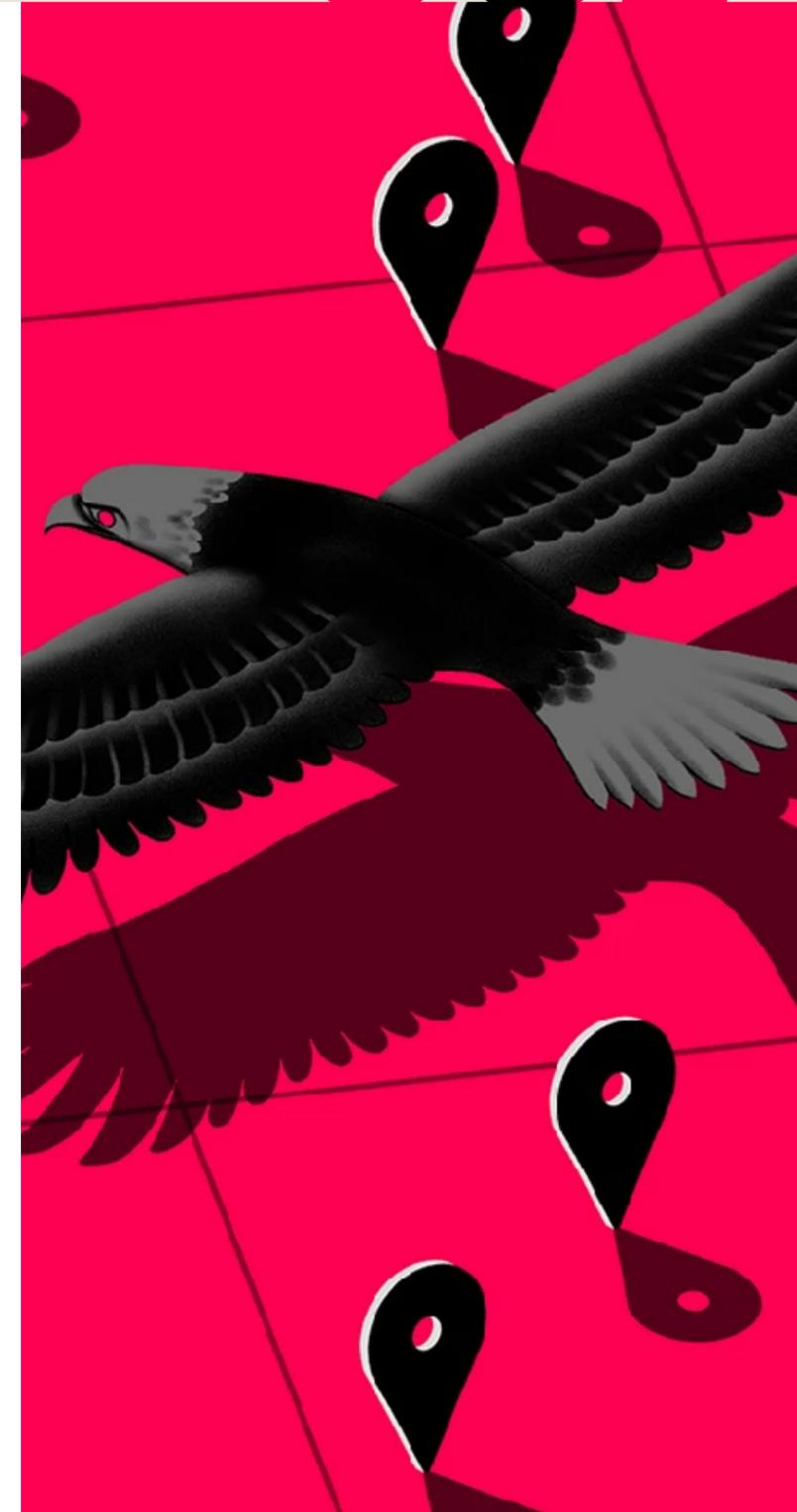
@mikarv

# Apps are pretty bad

- Users have more ability to control the Web, through browsers. For apps, they have almost none.

- One recent study identified 2,121 separate advertising tracking services in apps in the Android ecosystem, which can be grouped by ownership into approximately 292 parent organisations.*

- Another study found that 88.4% of apps contained a tracker owned by Alphabet (Google), 42.6% by Facebook, 33.9% by Twitter, 26.3% by Verizon and 22.2% by Microsoft. 30% of News apps, 28% of Family apps, and 25% of Gaming & Entertainment apps contain trackers from more than ten distinct tracker companies.**

*Abbas Razaghpanah and others, 'Apps, Trackers, Privacy, and Regulators: A Global Study of the Mobile Tracking Ecosystem' (2018) 13–14 <http://eprints.networks.imdea.org/1744/>; **Reuben Binns and others, 'Third Party Tracking in the Mobile Ecosystem', Proceedings of the 10th ACM Conference on Web Science (ACM 2018) 27

@mikarv

# This data isn't just used for advertising



- The US military, Customs and Border Protection, the Secret Service and Homeland Security buy location data from adtech firms as a way of avoiding obtaining warrants.

  - 'X-Mode' trackers in apps: e.g. a Quar'an/prayer time app downloaded by 98m Muslims around the world; a Craigslist searching app and a spirit level app.

  - Bidstream data from real-time bidding (Venntel, Babel Street)

- Many vendors found selling data that can specifically be used to locate women who have visited abortion clinics (Kochava, SafeGraph, Placer.ai)

# This infrastructure isn't just used for advertising

- NSA and GCHQ utilise the uniquely identifying Google "PREF" cookies to single out a user's computer and allow it to be remotely exploited using hacking tools developed by these state actors.

- DoubleClick cookies to reidentify Tor Browser users

- Major Belgacom hack achieved by GCHQ and CSEC through MUTANT BROTH, a system bulk-storing a range of cookies laid on popular websites, including those from Google and Facebook, in order to both identify users and build up a pattern of their daily habits and routines.

Ashkan Soltani and others, 'NSA Uses Google Cookies to Pinpoint Targets for Hacking', *Washington Post* (11 December 2013); Ryan Gallagher, 'Profiled: From Radio to Porn, British Spies Track Web Users' Online Identities' (*The Intercept*, 25 September 2015); Huib Modderkolk, 'Waarom kwam de Britse geheime dienst zo makkelijk weg met het hacken van Belgacom?' (*de Volkskrant*, 17 February 2018)

@mikarv

# But when they are used for ads, it's pretty bad too.

# Real-Time Bidding

From about 2010, automated auctions for your eyeballs.



Visitor | Publisher | Supply-side platform | Ad exchange | Demand-side platforms | Data management platforms

Request page
Get page (before ads)
Bid request for advert
Browser/cookie data
Bid request
Transmitted to 100s of would-be bidders
Enrichment of bid with tracking data
Winning advert is served
Auction for highest bid
Cookie sync with some bidders/DMPs

# Data sent to bidders each time this happens

- Site
  - URL of the site being visited
  - Site category or topic
- Device
  - Operating system
  - Browser software and version
  - Device manufacturer, model
  - Mobile provider
  - Screen dimensions

- User
  - Unique identifiers set by vendor and/or buyer.
  - Advertising exchange's cookie ID.
  - A demand-side platform's user identifier
  - Year of Birth
  - Gender
  - Interests
  - Metadata reporting on consent provided
- Geography
  - Longitude and latitude
  - Postal/ZIP code

| | |
|---|---|
| | ...es/Migraines |
| IAB7-24 | Heart Disease |
| IAB7-25 | Herbs for Health |
| IAB7-26 | Holistic Healing |
| IAB7-27 | IBS/Crohn's Disease |
| IAB7-28 | Incest/Abuse Support |
| IAB7-29 | Incontinence |
| IAB7-30 | Infertility |
| IAB7-31 | Men's Health |
| IAB7-32 | Nutrition |
| IAB7-33 | Orthopedics |
| IAB7-34 | Panic/Anxiety Disorders |
| IAB7-35 | Pediatrics |
| IAB7-36 | Physical Therapy |

# ... and retained

- Bid requests go to hundreds or thousands of companies; little oversight.

- **Vectaury** in France — small company ,with only 3.5m€ annual turnover — retained 68m bid request records (and fined by the French data regulator, CNIL) in 2018.

- Their website even claimed that they discarded 70% of all data, and only kept any of it for 12 months meaning that this small company was possibly sent 1/4 billion bid requests in just a single year.

SOLUTIONS          CMP          TECHNOL

PRIVACY IS HARDCO
VECTAURY'S D

We strive for the creation of a constr
sustainable ecosystem, serving all
stakeholders

# Data at scale for real-time bidding (RTB)

## Leading RTB exchanges, daily bid request estimates

| | |
|---|---|
| Index Exchange | 50 billion[ii] |
| OpenX | 60 billion+[i] |
| Rubicon Project | Unknown. Claims to reach 1 billion people's devices.[iii] |
| PubMatic | 70 billion+[iv] |
| Oath/AOL | 90 billion[v] |
| AppNexus | 131 billion[vi] |
| Smaato | 214 billion[vii] |
| Google DoubleClick | Unknown. DoubleClick is the dominant exchange. |

i. "Tour IX's Amsterdam and Frankfurt Data Centers", Index Exchange, 2 July 2018 (URL: https://www.indexexchange.com/tour-ix-amsterdam-frankfurt-data-centers/).
ii. "OpenX Ad Exchange", OpenX (URL: https://www.openx.com/uk_en/products/ad-exchange/).
iii. "Buyers", Rubicon Project, (URL: https://rubiconproject.com/buyers/).
iv. "How PubMatic Is Learning Machine Learning", PubMatic, 25 January 2019 (URL: https://pubmatic.com/blog/learning-machine-learning/)
v. "Maximize yield with Oath's publisher offerings", Oath, 3 April 2018 (URL: https://www.oath.com/insights/maximize-yield-with-oath-s-publisher-offerings/)

vi. 500 Billion / 29.6 = 18.6 billion impressions per day. Using AppNexus 1:11.5 ratio, this is 214 auctions per day. 500+ impressions figure cited in "Optimize your mobile strategy", Smaato, (URL: https://www.smaato.com/).
vii. "Transacting at a peak of 11.4 billion daily impressions, our marketplace handles more traffic each day than Visa, Nasdaq, and the NYSE combined" at https://www.appnexus.com/sell. Note that in 2017, AppNexus said in "AppNexus Scales with DriveScale", 2017, (URL: http://go.drivescale.com/rs/451-ESR-800/images/DRV_Case_Study_AppNexus-final.v1.pdf) that 10.7 billion "impressions transacted" came as a result of running 123 billion auctions. The impressions transacted to auctions ratio appears to be roughly 1:11.5. Therefore, the 11.4 daily impressions reported in 2018 equates to 131 billion auctions per day.

# "Nanotargeting"

- Some ad infrastructures facilitate messages to be targeted to the level of specific individuals.

- Studies on Facebook's ad infrastructure have modelled that this "nanotargeting" is possible using either the 4 rarest interests of an individual or the 22 random interests from the interests set Facebook assigns — both options make users unique on Facebook with a 90% probability (González-Cabañas et al., 2021).

- Various "war stories" of this occurring, from intimate partner abuse to Labour Party HQ targeting Jeremy Corbyn.

@mikarv

# 'Cookie Banners'

# No cookie banner??

LONG READ

We asked five menstruation apps for our demand but this what we found...

We asked five menstruation apps to give us access to our data. We got a dizzying dive into the most intimate information about us.

FR

For a world where technology

- Personal data in data protection law: information relating to an identified or identifiable natural person.

- Data protection law (e.g. the GDPR) requires a 'lawful basis' for all personal data processing

  - **Not a consent-first law**: If you are doing something aligned with the user, not using sensitive categories of data like ethnicity or health, you typically won't need consent.

- However, ePrivacy Directive: consent for storing or retrieving data for terminal devices (history of rootkits, tracking) if not necessary for the requested service.

*Article 4*
### Definitions

11. 'consent' of the data subject means any **freely given, specific, informed and unambiguous indication** of the data subject's wishes by which he or she, **by a statement or by a clear affirmative action**, signifies agreement to the processing of personal data relating to him or her;

*Article 7*
### Conditions for consent

1. Where processing is based on consent, the controller shall be able to **demonstrate** that the data subject has consented to processing of his or her personal data.

2. If the data subject's consent is given in the context of a written declaration which also concerns other matters, the request for consent shall be presented in a manner which is **clearly distinguishable from the other matters**, in an intelligible and easily accessible form, using clear and plain language. Any part of such a declaration which constitutes an infringement of this Regulation shall not be binding.

3. The data subject shall have the **right to withdraw his or her consent at any time**. The withdrawal of consent shall not affect the lawfulness of processing based on consent before its withdrawal. Prior to giving consent, the data subject shall be informed thereof. **It shall be as easy to withdraw as to give consent.**

4. **When assessing whether consent is freely given, utmost account shall be taken of whether, inter alia, the performance of a contract, including the provision of a service, is conditional on consent to the processing of personal data that is not necessary for the performance of that contract.**
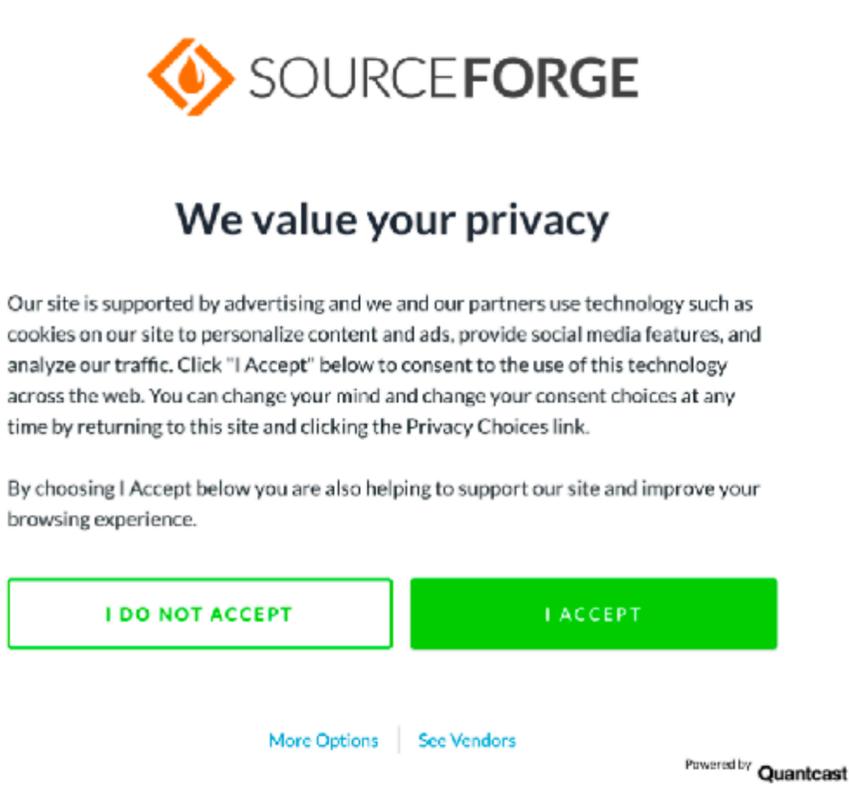
*Recital 32*

Consent should be given by a **clear affirmative act** establishing a **freely given, specific, informed and unambiguous indication** of the data subject's agreement to the processing of personal data relating to him or her, such as by a written statement, including by electronic means, or an oral statement. This could include **ticking a box when visiting an internet website**, choosing technical settings for information society services or another statement or conduct which clearly indicates in this context the data subject's acceptance of the proposed processing of his or her personal data. **Silence, pre-ticked boxes or inactivity should not therefore constitute consent.** Consent should cover all processing activities carried out for the same purpose or purposes. When the processing has multiple purposes, consent should be given for all of them. **If the data subject's consent is to be given following a request by electronic means, the request must be clear, concise and not unnecessarily disruptive to the use of the service for which it is provided.**
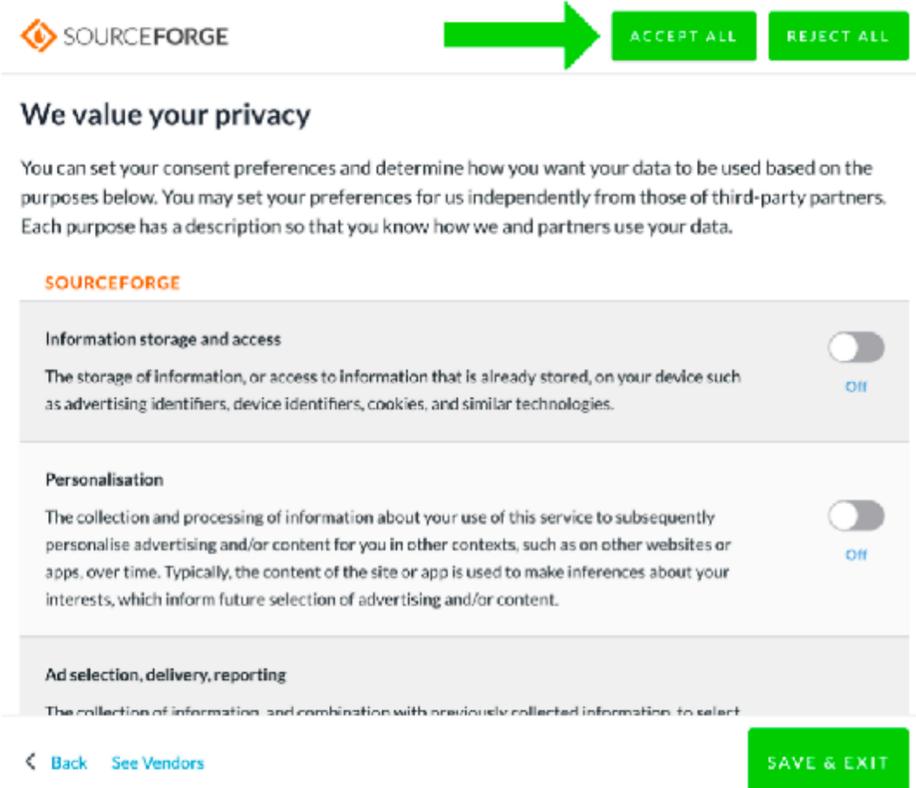
SIGN HERE
←

# With hundreds of trackers… how?
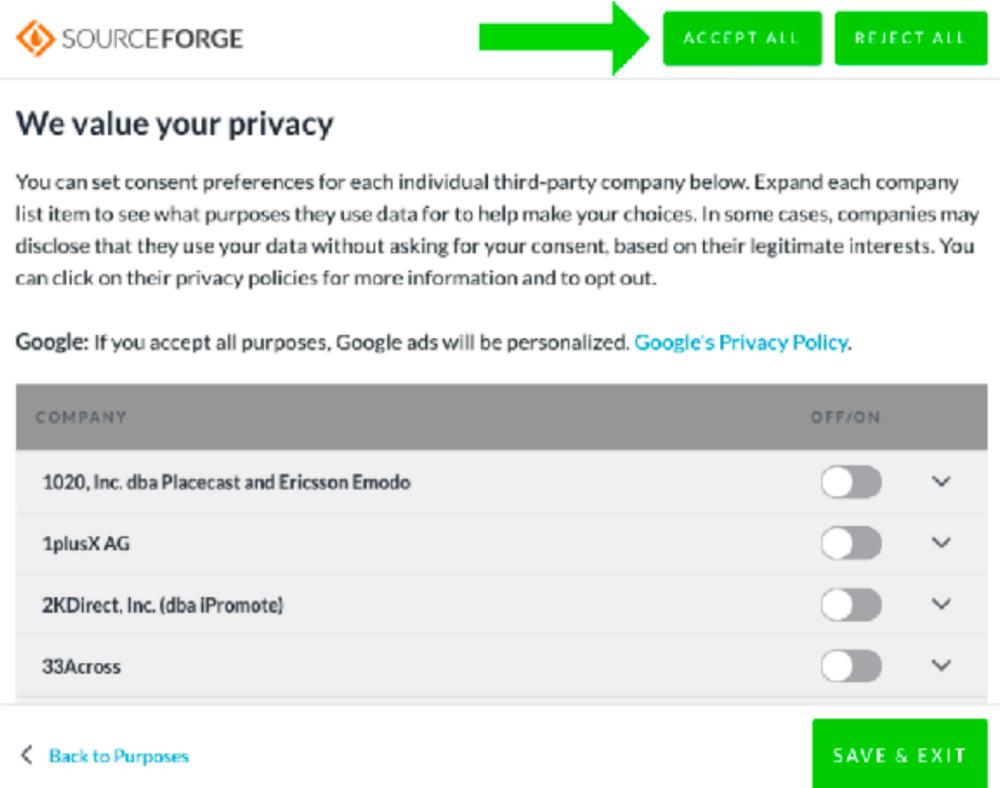
Consent management platforms emerge

# Legal entrepreneurship of an unsavoury kind



(a) First page    (b) Categories and purposes    (c) Vendors/third-parties

Figure 1. The three components of the QuantCast CMP on https://sourceforge.net as of September 2019.

# Many vendors: but are they compliant with the law?

UCL

**OneTrust**



**CrownPeak**



**CookieBot**



**Quantcast**



**TrustArc**



**Cookieinformation**

# Empirical, computational legal analysis to find out 🏛 **UCL**

**List of URLs to check**



```
● ● ●                    📄 url-list.txt ∨
au.dk
dtu.dk
ekstrabladet.dk
tv2.dk
aau.dk
bt.dk
sdu.dk
politiken.dk
um.dk
```

- Together with Aarhus University and MIT, we investigated whether these interfaces were providing valid consent under EU law.

- Built a bespoke *web scraper* and fed it the top 10K UK websites in 2019. We coded it to be able to analyse the top 5 CMPs to see how they were configured.

**Software analyses pop-ups**



```
17  //accept+reject buttons: cookiebot has two divs wi
18  const dialogBodyButtonsDisplay = document.getElemen
19  const dialogBodyLevelWrapperDisplay = document.getE
20
21  if (dialogBodyButtonsDisplay !== 'none') {
22      const acceptBtn = document.getElementById( eleme
23      const rejectBtn = document.getElementById( eleme
```
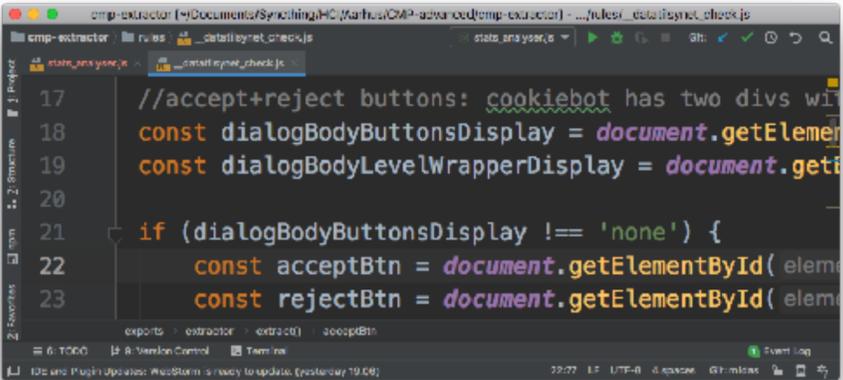
**Returns data on compliance**

| CMP | Explicit/implicit consent | Banner/barrier | Preticked options | Minimum compliance |
|---|---|---|---|---|
| Cookiebot | 45/40 | 78/7 | 64 (75.3%) | 2 (5.6%) |
| Crownpeak | 46/37 | 52/31 | 67 (80.7%) | 0 (0%) |
| OneTrust | 47/118 | 158/7 | 108 (65.4%) | 3 (1.8%) |
| QuantCast | 279/0 | 132/147 | 90 (32.3%) | 73 (26.2%) |
| TrustArc | 42/26 | 26/42 | 53 (77.9%) | 2 (2.9%) |
| **all** | **459/221** | **446/234** | **382 (56.2%)** | **80 (11.8%)** |

Table 1. Key statistics on scraped CMPs.

**@mikarv**

# And what did we find?

**Turned case law into three legal tests**

1. No optional boxes preticked

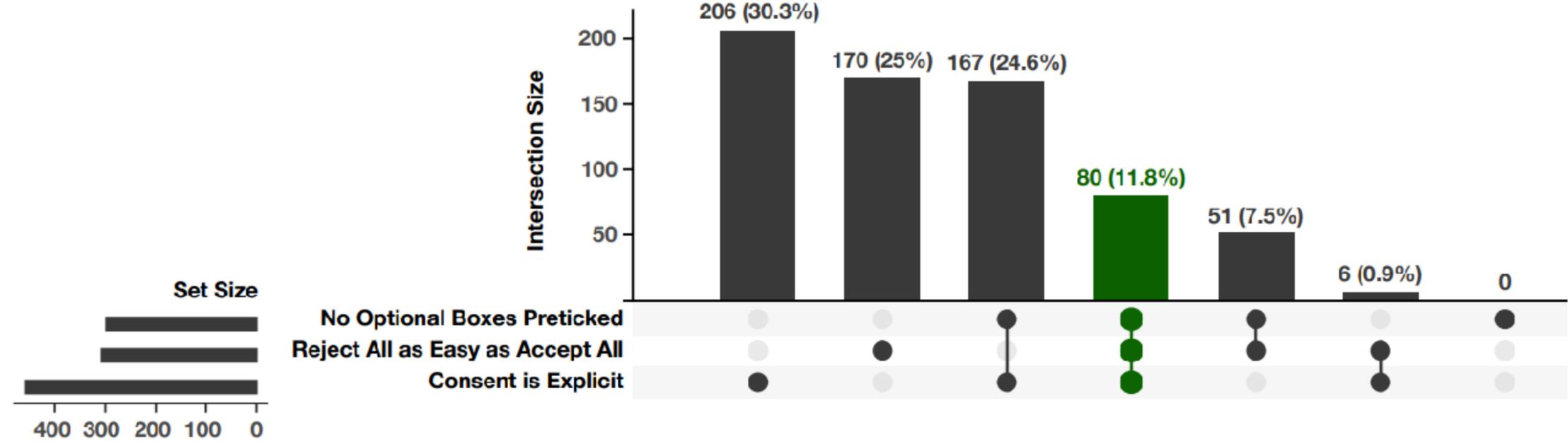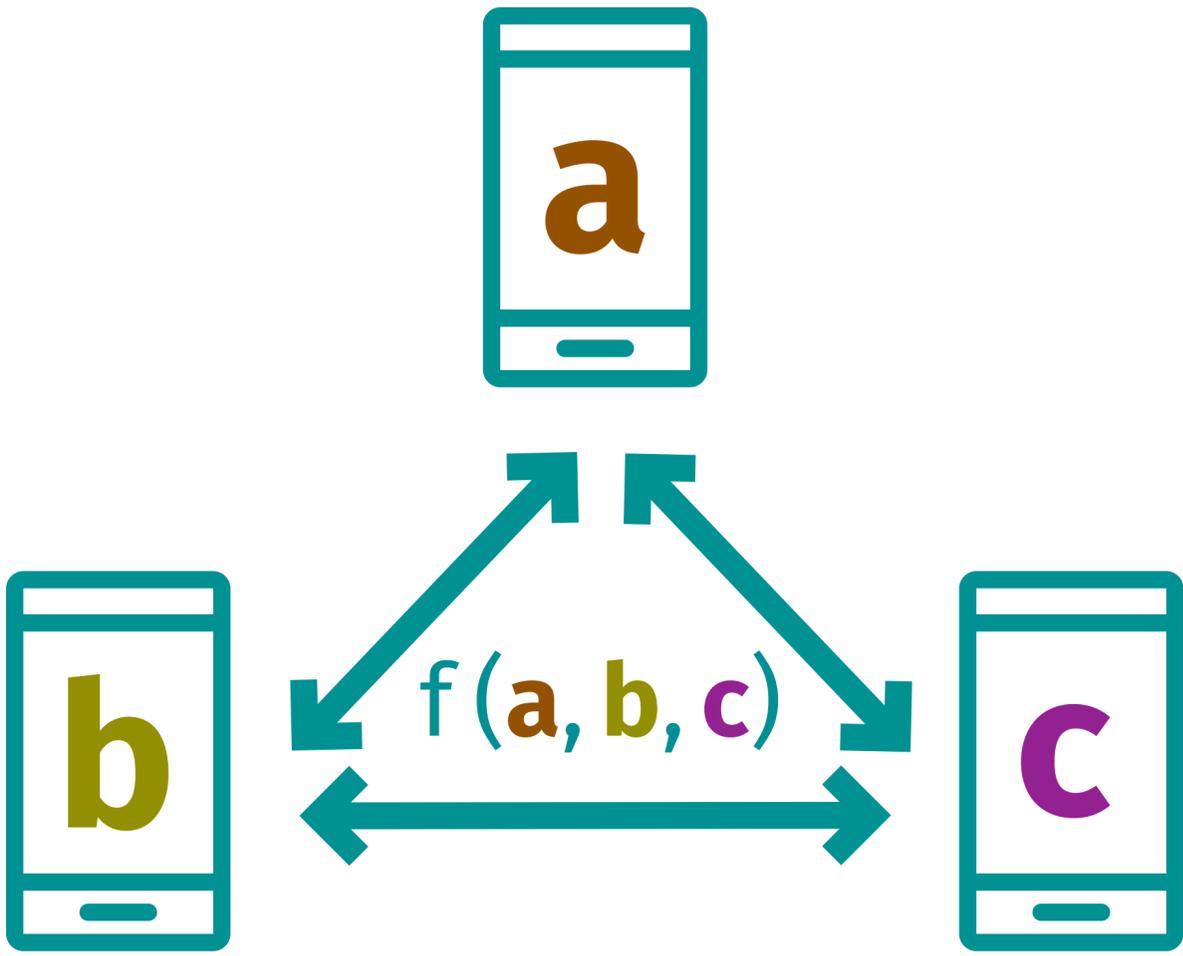2. Reject all as easy as Accept all

3. Consent is explicit



Figure 2. UpSet diagram [16, 36] of sites by adherence to three core conditions of EU law. Sites meeting all three in green.

| CMP | Sites | Median vendors (low./upp. quartiles) | Explicit/implicit consent | Banner/barrier | Preticked options | Minimum compliance |
|---|---|---|---|---|---|---|
| Cookiebot | 12.5% (85) | 104 (61, 232) | 45/40 | 78/7 | 64 (75.3%) | 2 (5.6%) |
| Crownpeak | 12.2% (83) | 38.5 (18.8, 132.3) | 46/37 | 52/31 | 67 (80.7%) | 0 (0%) |
| OneTrust | 24.3% (165) | 58 (26.5, 104.5) | 47/118 | 158/7 | 108 (65.4%) | 3 (1.8%) |
| QuantCast | 41% (279) | 542 (542, 542) | 279/0 | 132/147 | 90 (32.3%) | 73 (26.2%) |
| TrustArc | 10% (68) | 87 (38, 152) | 42/26 | 26/42 | 53 (77.9%) | 2 (2.9%) |
| **all** | **680** | **315 (58, 542)** | **459/221** | **446/234** | **382 (56.2%)** | **80 (11.8%)** |

Table 1. Key statistics on scraped CMPs.

Midas Nouwens and others, 'Dark Patterns after the GDPR: Scraping Consent Pop-Ups and Demonstrating Their Influence' in (ACM 2020) Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI 2020).

# Future of Targeting

UCL

a

b

$f(a, b, c)$

c

**multi-party
computation
among servers**

**secret
sharing**

- Can we target users in as before, but without data leaving devices?

- Train shared models with privacy enhancing data analysis so tracking data never leaves the phone.
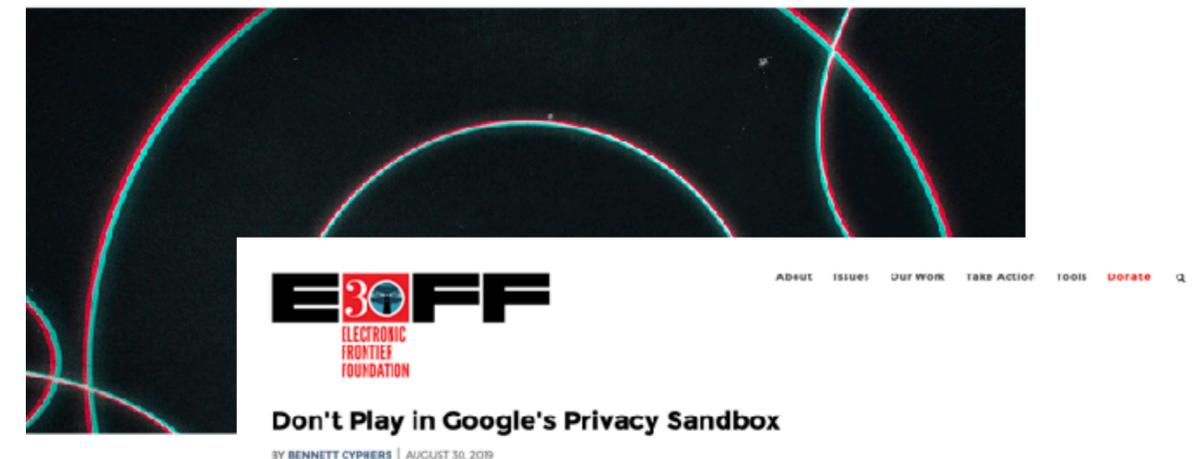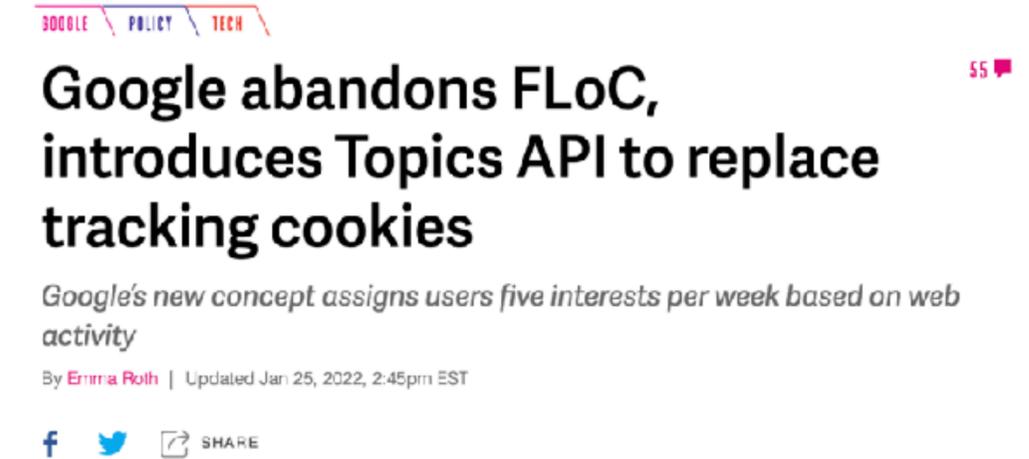
- Use other technologies such as 'zero-knowledge proofs' — or just locked down platforms — to check individuals are profiling themselves sufficiently, seeing the adverts.

Michael Veale, 'Future of online advertising: Adtech's new clothes might redefine privacy more than they reform profiling' (netzpolitik.org, 25 February 2022) <https://netzpolitik.org/2022/future-of-online-advertising-adtechs-new-clothes-might-redefine-privacy-more-than-they-reform-profiling-cookies-meta-mozilla-apple-google/> accessed 14 March 2022.

@mikarv

# Emerging moves



- Google's Privacy Sandbox
  - Investigation by UK competition and markets authority
- Others: Apple's hires in AdTech; Meta and Mozilla's proposals in the IETF.

# Questions for the future of online advertising

- **Interaction between input data and confidential computing**: what is the theoretical, ethical, and legal basis not to use e.g. blood pulse, eye-tracking, etc data when the output is confidential?

- **Is your device betraying you?** Rights for people facilitating computing, rather than just rights related to how data about you is used. Links to research ethics.