



HAL
open science

Tracking the Pixels: Detecting Web Trackers via Analyzing Invisible Pixels

Imane Fouad, Nataliia Bielova, Arnaud Legout, Natasa Sarafijanovic-Djukic

► **To cite this version:**

Imane Fouad, Nataliia Bielova, Arnaud Legout, Natasa Sarafijanovic-Djukic. Tracking the Pixels: Detecting Web Trackers via Analyzing Invisible Pixels. 2019. hal-01943496v2

HAL Id: hal-01943496

<https://inria.hal.science/hal-01943496v2>

Preprint submitted on 23 Jul 2019 (v2), last revised 17 Dec 2019 (v4)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Imane foudad*, Nataliia Bielova, Arnaud Legout, and Natasa Sarafijanovic-Djukic

Tracking the Pixels: Detecting Unknown Web Trackers via Analysing Invisible Pixels

Abstract: Web tracking has been extensively studied over the last decade. To detect tracking, previous studies and user tools rely on filter lists. However, there has always been a suspicion that lists miss many trackers. In this paper, we propose an alternative method to detect trackers inspired by analyzing behavior of invisible pixels. By crawling 84,658 webpages from 8,744 domains, we detect that third-party invisible pixels are widely deployed: they are present on more than 94.51% of domains and constitute 35.66% of all third-party images. We propose a fine-grained behavioral classification of tracking based on the analysis of invisible pixels. We use this classification to detect new categories of tracking and uncover new collaborations between domains on the full dataset of 4,216,454 third-party requests. We demonstrate that two popular methods to detect tracking, based on EasyList&EasyPrivacy and on Disconnect lists respectively miss 25.22% and 30.34% of the trackers that we detect. Moreover, we find that if we combine all three lists, 379,245 requests originated from 8,744 domains still track users on 68.70% of websites.

Keywords: online tracking; ad-blocker; cookie syncing; invisible pixels

DOI Editor to enter DOI

Received ..; revised ..; accepted ...

1 Introduction

The Web has become an essential part of our lives: billions are using Web applications on a daily basis and while doing so, are placing *digital traces* on millions of websites. Such traces allow advertising companies, as well as data brokers to continuously profit from collecting a vast amount of data associated to the users. Recent works have shown that advertising networks and data brokers use a wide range of techniques to track users across the Web [2, 8, 21, 23, 29, 35, 36, 38, 39, 42], from standard stateful cookie-based tracking [24, 39], to advanced cross-browser device fingerprinting [12, 35]. In

the last decade, numerous studies measured prevalence of third-party trackers on the Web [2, 10, 11, 23, 29–31, 35, 39, 45]. Web Tracking is often considered in the context of targeted behavioral advertising, but it's not limited to ads. Third-party tracking has become deeply integrated into the Web contents that website owners include

But what makes a tracker? How to recognize that a third-party request is performing tracking? To detect trackers, the research community applied a variety of methodologies, that has never been compared until now.

The most known Web tracking technique is based on *cookies*, but only some cookies contain unique identifiers and hence are capable of tracking the users. Some studies detect trackers by analysing cookie storage, and third-party requests and responses that set or send cookies [29, 39], while other works measured the mere presence of third-party cookies [30, 31]. To measure *cookie syncing*, researchers applied various heuristics to filter cookies with unique identifiers [1, 23, 24]. However, this approach has never been applied to detect tracking at large scale. Overall, previous works demonstrate that there is no unified method to identify third-party requests that are responsible for tracking.

Detection of identifier cookies and analysing behaviors of third-party domains is a complex task. Therefore, most of the state-of-the-art works that aim at measuring trackers at large scale rely on *filter lists*. In particular, EasyList [19] and EasyPrivacy [20] (EL&EP) and Disconnect [16] lists became the *de facto* approach to detect third-party tracking requests in privacy and measurement communities [9–11, 22, 23, 26–28, 38]¹. EasyList and EasyPrivacy are the most popular publicly maintained blacklist of know advertising and tracking requests, used by the popular blocking extensions AdBlock Plus [4] and uBlockOrigin [43]. Disconnect is another very popular list for detecting domains known for tracking, used in Disconnect browser extension [15] and in tracking protection of Firefox browser [25].

Nevertheless, filter lists detect only known tracking and ad-related requests, therefore a tracker can easily

¹ We summarize the usage of filter lists in security, privacy and web measurement community in Table 12 in the Appendix.

avoid this detection by registering a new domain. Third parties can also incorporate tracking behavior into functional website content, which is never blocked by filter lists because blocking functional content would harm user experience. Therefore, it is interesting to evaluate how effective are filter lists at detecting trackers, how many trackers are missed by the research community in their studies, and whether filter lists should still be used as the *default tools* to detect trackers at scale.

Our contributions: To evaluate the effectiveness of filter lists, we propose a new, fine-grained behavior-based tracking detection. Our results are based on data collected 4M third-party requests collected in a stateful crawl of more than 800K pages from 8K domains. We make the following contributions:

- *We analyse all the requests and responses that lead to invisible pixels (By “invisible pixels” we mean 1x1 pixel images or images without content.) Pixels are routinely used by trackers to send information or third-party cookies back to their servers: the simplest way to do it is to create a URL containing useful information, and to dynamically add an image tag into a webpage. This makes invisible pixels the perfect suspects for tracking. and propose a new classification of tracking behaviors. Our results show that pixels are still widely deployed: they are present on more than 94% of domains and constitute 35.66% of all third-party images. We found out that pixels are responsible only for 23.34% of tracking requests, and the most popular tracking content are scripts: a mere loading of scripts is responsible for 34.36% of tracking requests.*
- *We uncover hidden collaborations between third parties. We applied our classification on more than 4M third-party requests collected in our crawl, we have detected new categories of tracking and uncovered hidden collaborations between domains. We have discovered a new type of cookie syncing, called *first to third party cookie syncing*, that allows parties to merge profiles collected in first- and third-party context. This tracking appears on 67.96% of websites.*
- *We show that filter lists miss a significant number of cookie-based tracking. Our evaluation of the effectiveness of EasyList&EasyPrivacy and Disconnect lists shows that they respectively miss 25.22% and 30.34% of the trackers that we detect. Moreover, we find that if we combine all three lists, 379,245 requests originated from 8,744 domains still track users on 68.70% of websites.*
- *We explain why missed content is tracking the users. We identify two major reasons why missed content*

performs tracking: (i) the tracking cookies were set in a first-party context (as first-party cookies) and subsequently sent in a third-party request (as third-party cookies); (ii) a cookie set with a 2^{nd} -level TLD domain, is sent when content is loaded from any subdomain: for example, a third party `a.site.com` sets a cookie with `site.com` as its domain. The browser sends these cookies when fetching content from any other subdomain, such as `b.site.com`.

2 Methodology

To track users, domains deploy different mechanisms that have a different impacts on the user’s privacy. While some domains are only interested in tracking the user within the same website, others are recreating her browsing history by tracking her across sites. In our study, by “Web tracking” we refer to both within-site and cross-sites tracking.

To detect Web tracking, we first collect data from Alexa top 10,000 domains, then by analyzing the invisible pixels we define a new classification of Web tracking behaviors that we apply to the full dataset. In this section, we explain the data collection process and the criteria we used to detect identifier cookies and cookie sharing.

2.1 Data collection

Two stateful crawls: We performed passive Web measurements using the OpenWPM platform [23]. OpenWPM uses the Firefox browser, and provides browser automation by converting high-level commands into automated browser actions. We launched *two stateful crawls on two different machines*. For each crawl, we used one browser instance and saved the state of the browser between websites. In fact, measurement of Web tracking techniques such as cookie syncing is based on re-using cookies stored in the browser, and hence it is captured more precisely in a stateful crawl.

Full dataset: We performed a stateful crawl of the top 10,000 domains according to Alexa ranking in February 2019 in France [6] in two different machines. Due to the dynamic behavior of the websites, the content of a same page might differs every time this page is visited. To reduce the impact of this dynamic behavior and reduce the difference between the two crawls, we launched the two crawls at the same time. For each domain, we visited the home page and the first 10 links

pointing to pages in the same domain. The timeout for loading a homepage is set up to 90s, and the timeout for loading a link on the homepage is set up to 60s. Out of 10,000 Alexa top domains, we successfully crawled 8,744 domains with a total of 84,658 pages.

For every page we crawl, we store the HTTP request (url, method, header, date, and time), the HTTP response (url, method, status code, header, date, and time), and the cookies (both set/sent and a copy of the browser cookie storage) to be able to capture the communication between the client and the server. We also store the body of the HTTP response if it's an image with a *content-length* less than 100 KB. We made this choice to save storage space. Moreover, in addition to HTTP requests, responses and cookies, we were only interested in the storage of invisible pixels. In our first dataset, named *full dataset*, we capture all HTTP requests, responses, and cookies.

Prevalence of invisible pixels: As a result of our crawl of 84,658 pages, we have collected 2,297,716 images with a *content-length* less than 100 KB that represents 89.83% of the total number of delivered images. Even though we didn't store all the images, we were able to get the total number of delivered images using the content-type HTTP header extracted from the stored HTTP responses.

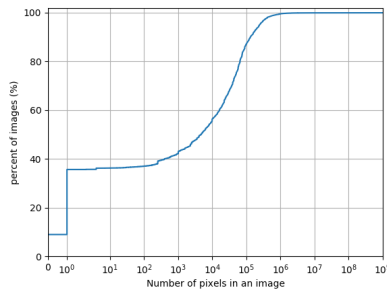


Fig. 1. Cumulative function of the number of pixels in images with a *content-length* less than 100 KB. 35.66% of images are invisible pixels: 9.00% have no content (they are shown as zero-pixel images) and 26.66% are of size 1x1 pixels.

Figure 1 shows the distribution of the number of pixels in all collected images. We notice that invisible pixels (1x1 pixels and images with no content) represent 35.66% of the total number of collected images.

We found that out of 8,744 successfully crawled domains, 8,264 (94.51%) domains contain at least one page with one invisible pixel. By analyzing webpages independently, we found that 92.85% out of 84,658 visited pages include at least one invisible pixel.

Invisible pixels subdataset: The invisible pixels do not add any content to the Web pages. However, they are widely used in the web. They generally allow the third party to send some information using the requests sent to retrieve the images. Moreover, the user is totally unaware of their existence. Hence, every invisible pixel represents a threat to the user privacy. We consider the set of requests and responses used to serve the invisible pixels as a ground-truth dataset that we call *invisible pixels dataset*. The study of this *invisible pixels dataset* allow us to excavate the tracking behaviors of third party domains in the web.

2.2 Detecting identifier cookies

Cookies are a classical way to track users in the web. A key task to detect this kind of tracking is to be able to detect cookies used to store identifiers. We will refer to these cookies as *identifier cookies*. In order to detect identifier cookies we analyzed data extracted from the two simultaneous crawls from different machines. We refer to the owner of the cookie as host, and we define cookie instance as (host, cookie-name, cookie-value)

We compare cookies instances between the two crawls: A tracker associates different identifiers to different users in order to distinguish them. Hence, an identifier cookie should be unique per user (user specific). We analyzed the 8,744 crawled websites where we have a total 607,048 cookies instances belonging to 179,580 cookies (host, name). If the same cookie instance appears in the two crawls, we consider that it's is not used for tracking. We refer to such cookies as *safe cookies*. We found 108,252 safe cookies instances. They represent 17.83% of cookies instances.

Due to the dynamic behavior of websites, not all cookies appear in both crawls. We mark as unknown cookies, cookies (host,name) that appear only in one crawl. In total we found 15,386 unknown cookie (8.56%). We exclude these cookies from our study.

We don't consider the cookie lifetime: The lifetime of the cookie is used to detect identifier cookies in related works [1, 23, 24]. Only cookies that expire at least a month after being placed are considered as identifier cookies. In our study, we don't put any boundary on the cookie lifetime because domains can continuously update cookies with a short lifetime and do the mapping of these cookies on the server side which will allow a long term tracking.

Detection of cookies with identifier cookie as name: We found that some domains stores the iden-

Host	# cookies instances
lpsnmedia.net	583
i-mobile.co.jp	223
rubiconproject.com	83
justpremium.com	72
juicyads.com	64
kinoafisha.info	64
aktualne.cz	63
maximonline.ru	61
sexad.net	47
russian7.ru	45

Table 1. Top 10 domains storing the identifier as cookie name.

tifier cookie as part of the cookie name. To detect this kind of behavior, we analyzed the cookies with the same host, cookie value and a different cookie name across the two crawls. we found 5,295 (0.87%) cookies instances with identifier cookie as name, This behavior was performed by 966 different domains. Table 1 present the top 10 domains involved. The cookies with identifier cookie as name represent only 0.87% of the total number of cookies. Therefore, we will exclude them from our study.

2.3 Detecting identifier sharing

Third party trackers not only collect data about the users, but also exchange it to build richer users profiles. Cookie syncing is a common technique used to exchange user identifier stored in cookies. To detect such behaviors we need to detect the *identifier cookies* shared between domains. A cookie set by one domain cannot be accessed by another domain because of the cookie access control and Same Origin Policy[41]. Therefore, trackers need to pass identifiers through the URL parameters.

Identifier sharing can be done in different ways: it can be sent in clear as a URL parameter value, or in a specific format, encoded or even encrypted. To detect identifiers, we take inspiration from [1, 23]. We split cookies and URL parameter values using as delimiters any character not in `[a-zA-Z0-9,'-','_',';']`.

We have deployed six different techniques to detect identifier sharing, described in Figure 2. The first three methods are generic: either the identifier is sent as the parameter value (**Direct sharing**), as part of the parameter value (**ID as part of the parameter**) or it's stored as part of the cookie value and sent as parameter value (**ID as part of the cookie**).

We noticed that the requests for invisible images, where we still didn't detect any cookie sharing originate mostly from `google-analytics.com` and `doubleclick.net`. Indeed, these domains are prevalent

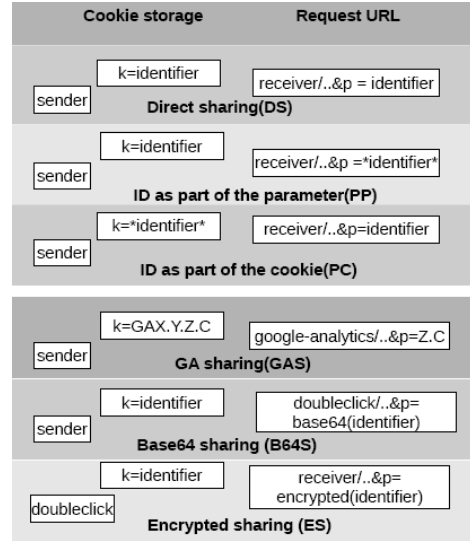


Fig. 2. Detecting identifier sharing. "Sender" is the domain that owns the cookie and triggers the request, "receiver" is the domain that receives this request and "identifier" is a cookie value that we detected as *identifier cookie*. "*" represents any string.

in serving invisible pixels across websites (see Figure 12 in Appendix). We therefore base the next techniques on these two use cases.

First, we notice that first party cookies set by `google-analytics.com` have the format `GAX.Y.Z.C`, but the identifier sent to it are of the form `Z.C`. We therefore detect this particular type of cookies, that were not detected in previous works that rely on delimiters (**GA sharing**). Second, by base64 decoding the value of the parameter sent to `doubleclick.net`, we detect the encoded sharing (**Base64 sharing**). Finally, by relying on Doubleclick documentation [18] we infer that encrypted cookie was shared (**Encrypted sharing**). For more details see the Section 9.1 in the Appendix.

2.4 Limitations

We detected six different techniques used to share the identifier cookie. However, trackers today use various techniques to share their identifier cookie, One example of such techniques, is to encrypt the cookie before sharing it. We only detected encrypted cookies when it's shared following a specific semantic set by `doubleclick` [18]. Moreover, in our study we don't inspect the payload of POST requests that could be used to share the identifier cookie. For example, it's known that `google-analytics.com` sends the identifier cookie as part of the URL parameters with GET requests or in the payload of the POST requests [7]. In our analysis we will not be detecting the second case.

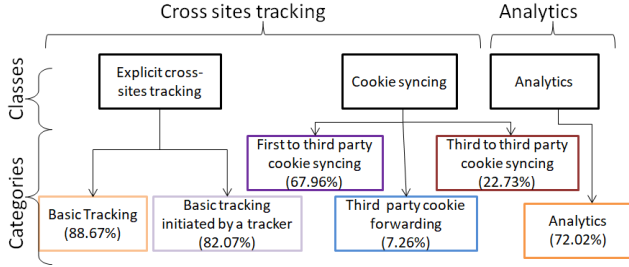


Fig. 3. Classification overview. (%) represents the percentage of domains out of 8,744 where we detected the tracking behavior. A tracking behavior is performed in a domain if it's detected in at least one of its pages.

To detect the sender of the request in case of inclusion, we use the refer field. Therefore, we may miss interpreting who is the effective initiator of the request, it can be either the first party or an included script.

3 Overview of tracking behaviors

In section 2.1, we detected that invisible pixels are widely present on the Web and are perfect suspects for tracking. In this section, we detect the different tracking behaviors by analyzing *invisible pixels dataset*

In total, we have 747,816 third party requests leading to invisible images. By analyzing these requests, we detected 6 categories of different tracking behaviors in 636,053 (85.05%) requests that lead to invisible images.

We further group the categories into three main classes: explicit cross-sites tracking (Section 4.1), cookie syncing (Section 4.2), and analytics (Section 4.3).

After defining our classification using the *invisible pixels dataset*, we apply it on the *full dataset* where we have a total of 4,216,454 third-party requests collected from 84,658 pages on the 8,744 domains successfully crawled. By analyzing these requests, we detected the 6 tracking behaviors in 2,724,020 (64.60%) requests.

Figure 3 presents an overview of all classes (black boxes) and categories of tracking behaviors and their prevalence in the full dataset. Out of 8,744 crawled domains, we identified at least one form of tracking in 91.92% domains. We have further analyzed prevalence of each tracking category that we report in Section 4. We found out that *first to third party cookie syncing* (see Sec. 4.2.3) that has never been detected in the previous works, actually appears on 67.96% of the domains!

In addition, we analyzed the most prevalent domains involved in either cross-sites tracking, analyt-

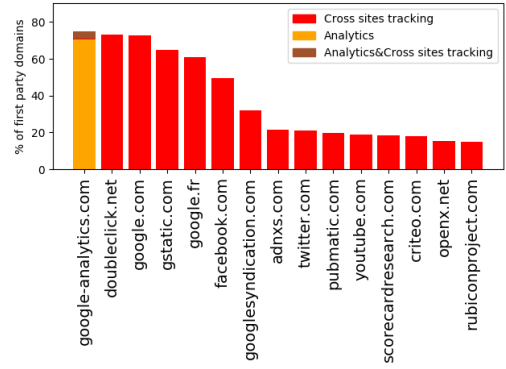


Fig. 4. Top 15 third parties involved in analytics, cross-sites tracking or both behaviors on the same first-party domain.

ics, or both behaviors. Figure 4 demonstrated that a third party domain may have several behaviors. For example, we detect that `google-analytics.com` exhibits both cross-sites tracking and analytics behavior.

Content type	% requests
Script	34.36%
Invisible images	23.34 %
Text/html	20.01%
Big images	8.54 %
Application/json	4.32%

Table 2. Top 5 types of content used in the 2,724,020 third party tracking requests.

We found that not all the tracking detected in the *full dataset* is based on invisible pixels. We extracted the type of the content served by the tracking requests using the HTTP header *Content-Type*. Table 2 presents the top 5 types of content used for tracking. Out of the 2,724,020 requests involved in at least one tracking behavior in the full dataset, the top content delivered by tracking requests is scripts (34.36%), while the second most common content is invisible pixels (23.34%). We also detected other content used for tracking purposes such as visible images.

4 Classification of tracking

In this section, we explain all the categories of tracking behaviors presented in Figure 3 that we have uncovered by studying the *invisible pixels dataset*. For each category, we start by explaining the tracking behavior, we then give its privacy impact on the user's privacy, and finally we present the results from the *full dataset*.

4.1 Explicit cross-sites tracking

Explicit cross-sites tracking class includes two categories: *basic tracking* and *basic tracking initiated by a tracker*. In both categories, we do not detect cookie syncing that we analyze separately in Section 4.2.

4.1.1 Basic tracking

Basic tracking is the most common tracking category as we see from Figure 3.

Tracking behavior: Basic tracking happens when a third party domain, say A. com, sets an identifier cookie in the user’s browser. Upon a visit to a webpage with content from A. com, a request is sent to A. com with its cookie. Using this cookie, A. com identifies the user across all websites that include content from A. com.

Privacy impact: *Basic tracking* is the best known tracking technique that allows third parties to track the user across websites, hence to recreate her browsing history. However, third parties are able to track the user only on the websites where their content is present.

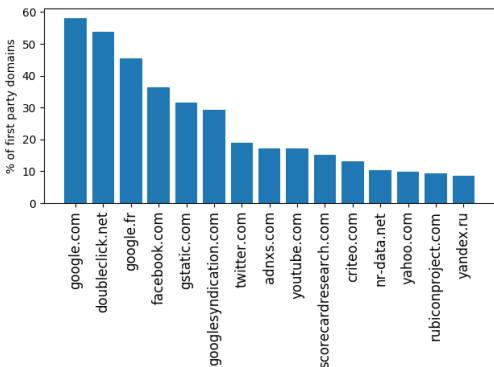


Fig. 5. Basic tracking: Top 15 cross-sites trackers in 8,744 domains.

Results: We detected basic tracking in 88.67% of visited domains. In total, we found 5,421 distinct third parties making basic tracking. Figure 5 shows the top domains involved in basic tracking. We found that google.com alone is tracking the user on over 5,079 (58.08%) domains. By only using the *basic tracking*, google.com recreates 46.17% of the user’s browsing history. We define the browsing history as the 84,658 visited pages. This percentage becomes more important if we consider the company instead of the domain. In fact, by only using the *basic tracking* Google recreates 54.42% of the user’s browsing history.

4.1.2 Basic tracking initiated by a tracker

When the user visits a website that includes content from a third party, the third party can redirect the request to a second third party tracker or include it. The second tracker will associate his own identifier cookie to the user. In this case the second tracker is not directly embedded by the first party and yet it can track her.

Tracking behavior: *Basic tracking initiated by a tracker* happens when a basic tracker is included in a website by another basic tracker.

Privacy impact: By redirecting to each other, trackers trace the user activity on a larger number of websites. They gather the browsing history of the user on websites where at least one of them is included. The impact of these behavior on the user’s privacy could be similar to the impact of cookie syncing. In fact, by mutually including each other on websites, each tracker can recreate the combination of what both partners have collected using basic tracking. Consequently, through *basic tracking initiated by a tracker* trackers share the user’s browsing history without syncing cookies.

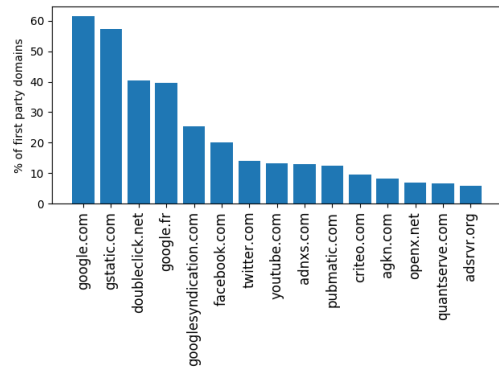


Fig. 6. Basic tracking initiated by a tracker: Top 15 trackers included in 8,744 domains.

Results: We detected Basic tracking initiated by a tracker in 82.07% of domains. From Figure 6 we can notice that google.com is the top tracker included by other third parties. It’s included in over 5,374 (61.45%) domains. By only relying on it’s partners, without being directly included by the developer, google.com can recreate 32.51% of the user’s browsing history.

Note that a domain can implement multiple kinds of tracking on the same page. Google.com is included by 295 different third party trackers in our dataset. In Table 3 we present the top 10 partners that are mutually including each other on websites.

Partners	# requests
googlesyndication.com ↔ doubleclick.net	176,295
doubleclick.net ↔ google.com	53,371
gstatic.com ↔ google.com	22,547
google.com ↔ google.fr	12,990
google.com ↔ youtube.com	5,289
pubmatic.com ↔ doubleclick.net	4,392
criteo.com ↔ doubleclick.net	2,258
googlesyndication.com ↔ adnxs.com	1,508
googlesyndication.com ↔ openx.net	1,344
doubleclick.net ↔ adnxs.com	1,199

Table 3. Basic tracking initiated by a tracker: Top 10 partners that include each other. (↔) both ways inclusion.

4.2 Cookie syncing

To create a more complete profile of the user, third party domains need to merge profiles they collected on different websites. One of the most known techniques to do so is cookie syncing. We separate the previously known technique of cookie syncing [1, 23] into two distinct categories, *third to third party cookie syncing* and *third party cookie forwarding*, because of their different privacy impact. We additionally detect a new type of cookie syncing that we call *first to third party cookies syncing*

4.2.1 Third to third party cookie syncing

When two third parties have an identifier cookie in user’s browser and need to merge user profile, they use third to third party cookie syncing.

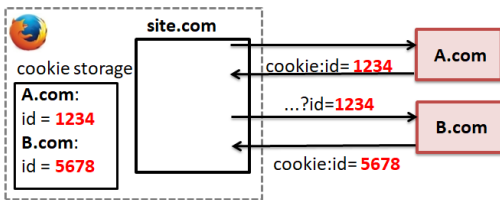


Fig. 7. Third to third party cookie syncing behavior.

Tracking explanation: Figure 7 demonstrates cookie syncing². The first party domain includes a content having as source the first third party A.com. A request is then sent to A.com to fetch the content. Instead of sending the content, A.com decides to redirect

² Notice that in figures that explain a tracking behavior, we show cookies only in the response, and never in a request. This actually represents both cases when cookies are sent in the request and also set in the response.

to B.com and in the redirection request sent to B, A.com includes the identifier it associated to the user. In our example B.com will receive the request B.com?id=1234 where 1234 is the identifier associated by A.com to the user. Along with the request, B.com will receive its cookie *id = 5678* which will allow B.com to link the two identifiers to the same user.

Privacy impact: *Third to third party cookie syncing* is one of the most harmful tracking techniques that impacts the user’s privacy. In fact, third party cookie syncing can be seen as set of trackers performing *basic tracking* and then exchanging the data they collected about the user. It’s true that a cross sites tracker recreates part of the user’s browsing history but this is only possible on the websites on which it was embedded. Using cookie syncing a tracker not only log the user’s visit to the websites where it’s included but it can also log her visits to the websites where it’s partners are included. What makes this practice even more harmful is when a third party has several partners with whom it syncs cookies One example of such behavior is rubiconproject.com that syncs its identifier cookie with 7 partners: tapad.com, openx.net, imrworldwide.com spotxchange.com, casalemedia.com, pubmatic.com and bidswitch.net.

Partners	# re-requests	Sharing technique
adnxs.com → criteo.com	1,962	→DS
doubleclick.net → facebook.com	789	→DS
casalemedia.com → adsvr.org	778	→DS
mathtag.com ↔ adnxs.com	453	→DS
trafficjunky.net ↔ traf- ficjunky.net	415	↔ DC, PPS
pubmatic.com → lijit.com	321	→DS
lijit.com ↔ lijit.com	303	↔ DC
media.net ↔ media.net	302	↔ DC, PCS
adobedtm.com → facebook.com	269	→DS
doubleclick.net ↔ criteo.com	250	→ DC, PCS; ← DS

Table 4. Third to third party cookie syncing: Top 10 partners.

The arrows represent the flow of the cookie synchronization, (→) one way matching or (↔) both ways matching. DS (Direct sharing), B64S (Base64 sharing), ES (Encrypted sharing), PCS (ID as part of the cookie), PPS (ID sent as part of the parameter) are different sharing techniques described in Figure 2.

Results: We detected third to third party cookie syncing in 22.73% websites. We present in Table 4 the top 10 partners that we detect as performing cookie syncing. In total we have detected 1,263 unique partners performing cookie syncing. The syncing could be done

in both ways as it's the case for doubleclick.net and criteo.com or in one way as it's the case for adnxs.com and criteo.com. In case of two ways matching we noticed that the two partners can perform different identifier sharing techniques. We see the complexity of the third to third cookie syncing that involves a large variety of sharing techniques. We also noticed that cookie syncing can be done between two subdomains from the same domain as it's the case for trafficjunky.

4.2.2 Third party cookie forwarding

The purpose of the collaboration between third party domains in *third party cookie forwarding* is to instantly share the browsing history. Cookie forwarding has always been called “syncing” while instead it simply enables a third party to reuse an identifier of a tracker, without actually syncing its own identifier.

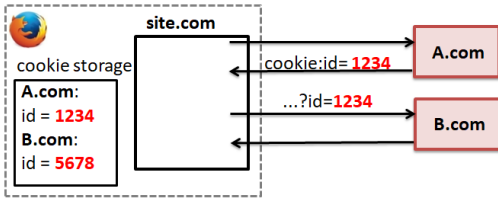


Fig. 8. Third party cookie forwarding behavior.

Tracking explanation: The first party domain site.com includes A.com’s content. To get the image a request is sent to A.com along with its cookie. A.com then redirects the request to its partner (B.com) and send as part of the URL parameters the identifier cookie it associated to the user (1234) (Figure 8).

Third party cookie forwarding differs from *Third to third party cookie syncing* depending on whether their is a cookie set by the receiver in the browser or not. This category is similar to Third party advertising networks in Roesner et al. and Lerner et al.’s work [39] [29], in the sense that we have a collaboration of third party advertisers. However, in our study we check that the second tracker do not use its own cookie to identify the user. This means that this tracker (B.com) is relying on the first one (A.com) to track the user. In fact B.com uses A.com’s identifier to recreate her browsing history.

Privacy impact: *Third party cookie forwarding* allows trackers to instantly share the browsing history of the user. A.com in Figure 8 does not only associate an identifier cookie to the user but it also redirect and share this identifier cookie with its partner. This practice allows both A.com and B.com to track the user across

websites. From a user privacy point of view, *third party cookie forwarding* is not as harmful as cookie syncing because the second tracker in this case does not contribute in the user’s profile creation but it passively receives the user’s browsing history from the first tracker.

Results: We detected Third party cookie forwarding in 7.26% of visited websites. To our surprise, the top domain receiving identifier cookie from third parties is google-analytics (Figure 13 in Appendix). Google-analytics is normally included by first party domains to get analytics of their website, it’s known as a *within domain tracker*. But in this case, google-analytics is used by the third party domains. The third party is forwarding its third party cookie to google-analytics on different websites, consequently google-analytics in this case is tracking the user across websites. This behavior was discovered by Roesner et al. [39]. They reported this behavior in only a few instances, but in our dataset we found 386 unique partners that forward cookies among which 271 are forwarding cookies to google-analytics.

Third parties	# requests
adtrue.com	298
google.com	123
architonic.com	120
bidgear.com	80
akc.tv	76
insticator.com	73
coinad.com	64
performgroup.com	52
chaturbate.com	47
2mdnsys.com	40

Table 5. Third party cookie forwarding; Top 10 third parties forwarding cookies to google-analytics.

In Table 5, we present the top 10 third parties forwarding cookies to google-analytics services.

4.2.3 First to third party cookie syncing

In this category, we detect that first party cookie get synced with third party domain.

Tracking explanation: Figure 9 demonstrates the cookie syncing of the first party cookie. The first party domain site.com includes a content from A.com?id=abcd, where A.com is a third party and abcd is the first party identifier cookie of the user set for site.com. A.com receives the first party cookie abcd in the URL parameters, and then redirects the request to

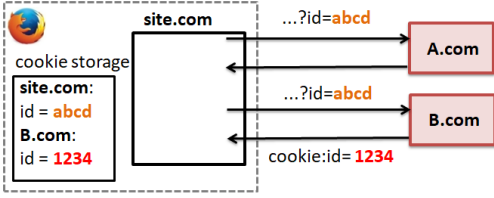


Fig. 9. First to third party cookie syncing behavior.

B.com. As part of the request redirected to B.com, A.com includes the first party identifier cookie. B.com sets its own identifier cookie *1234* in the user’s browser. Using these two identifiers (the first party’s identifier *abcd* received in the URL parameter and its own identifier *1234* sent in the cookie), B.com can create a matching table that allows B.com to link both identifiers to the same user.

The first party cookie can also be shared directly by the first party service (imagine Figure 9 where A.com is absent). In that case, site.com includes content from B.com and as part of the request sent to B.com, site.com sends the first party identifier cookie *1234*. B.com sets its own identifier cookie *1234* in the user’s browser. B.com can now link the two identifiers to the same user.

Privacy impact: In our study, we differentiate the case when cookie shared is a first party cookie and when it is a third party cookie. We made this distinction because, the kind and the sensitivity of the data shared differs in the two cases. Using this tracking technique, first party websites get to sync cookies with third parties. Moreover, pure analytic services allow to sync in-site history with cross-site history.

Partners	# requests
First party cookie synced through an intermediate service	
google-analytics.com → doubleclick.net	8,297
Direct first to third party cookie syncing	
hibapress.com → criteo.com	460
alleng.org → yandex.ru	332
arstechnica.com → condenastdigital.com	243
thewindowsclub.com → doubleclick.net	228
digit.in → doubleclick.net	224
misionesonline.net → doubleclick.net	221
wired.com → condenastdigital.com	219
newyorker.com → condenastdigital.com	218
uol.com.br → tailtarget.com	198

Table 6. First to third party cookie syncing: Top 10 partners.

Results: We detected first to third party cookie syncing in 67.96% of visited websites. In Table 6, we present the top 10 partners syncing first party cookies. We differentiate the two cases: (1) first party

cookie synced through an intermediate service (as shown in Figure 9) and (2) first party cookie synced directly from the first party domain. In total we found 17,415 different partners involved. The top partners are google-analytics and doubleclick. We found that google-analytics first receives the cookie as part of the URL parameters. Then, through redirection process, google-analytics transfers the first party cookie to doubleclick that inserts or receives an identifier cookie in the user’s browser. We found out that google-analytics is triggering such first party cookie syncing on 38.91% of visited websites.

4.3 Analytics category

Instead of measuring website audience themselves, websites today use third party analytics services. Such services provide reports of the website traffic by tracking the number of visits, the number of visited pages in the website, etc. The first party website includes content from the third party service on the pages it wishes to analyze the traffic.

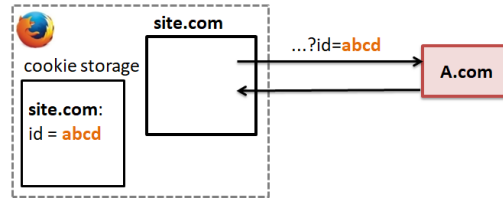


Fig. 10. Analytics behavior.

Tracking explanation: Figure 10 shows the analytics category where the domain directly visited by the user (site.com) owns a cookie containing a unique identifier in user’s browser. Such cookie is called first party identifier cookie. This cookie is used by the third party (A.com) to uniquely identify the visitors within site.com. The first party website makes a request to the third party to get the content and uses this request to share the first party identifier cookie.

Privacy impact: In analytic behavior, the third party domain is not able to track the user across websites because it does not set its own cookie in user’s browser. Consequently, for this third party, the same user will have different identifiers in different websites. However, using the first party identifier cookie shared by the first party, the third party can identify the user within the same website. From a user point of view, analytics behavior is not as harmful as the other tracking

methods because the third party domain can not recreate the user’s browsing history but it can only track her activity within the same domain which could be really useful for the website developer.

Results: We detected analytics in 72.02% of visited websites. We detect that google-analytics is the most common analytics service. It’s used on 69.25% of the websites. The next most popular analytics is alexametrics.com – it’s prevalent on 9.10% of the websites (see Figure 14 in the Appendix).

5 Are filter lists effective at detecting trackers?

Most of the state-of-the-art works that aim at measuring trackers at large scale rely on *filter lists*. In particular, EasyList [19] and EasyPrivacy [20] and Disconnect [16] lists became the *de facto* approach to detect third-party tracking requests in the privacy and measurement communities [9–11, 22, 23, 26–28, 38]. Nevertheless, filter lists detect only known tracking and ad-related requests, therefore a tracker can easily avoid this detection by registering a new domain. Third parties can also incorporate tracking behavior into functional website content, which could not be blocked by filter lists because blocking functional content would harm user experience. Therefore, it is interesting to evaluate how effective are filter lists at detecting trackers, how many trackers are missed by the research community in their studies, and whether filter lists should still be used as the *default tools* to detect trackers at scale.

In this section, we analyse how effective are filter lists at detecting third-party trackers. To do that, we compare all the cross-site tracking and analytics behavior reported in Section 4 (that we unite under one detecting method, that we call BehaviorTrack) with the third-party trackers detected by filter lists. EasyList and EasyPrivacy (EL&EP) and Disconnect filter lists in our comparison were extracted in April 2019. We use the python library *adblockparser* [3], to determine if a request would have been blocked by EL&EP. For Disconnect we compare to the domain name of the requests (Disconnect list contains full domain names, while EL&EP are lists of regular expressions that require parsing).

For the comparison, we used the *full dataset* of 4,216,454 third party requests collected from 84,658 pages of 8,744 successfully crawled domains.

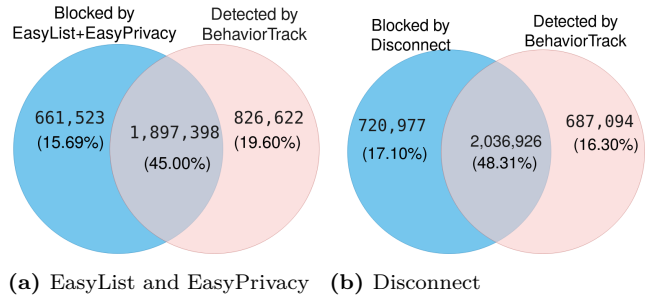


Fig. 11. Effectiveness of filter lists at detecting trackers on 4,216,454 third party requests from 84,658 pages.

Measuring tracking requests We apply filter lists on requests to detect which requests are blocked by the lists, as it has been done in previous works [23]. We then use the filter lists to classify follow-up third-party requests that would have been blocked by the lists. This technique has been extensively used in the previous works [22, 26–28] (for more details, see Table 12 in the Appendix). We classify a request as blocked if it matches one of the conditions:

- the requests directly matches the list;
- the request is a consequence of a redirection chain where an earlier request was blocked by the list.
- the request is loaded in a third-party content (an iframe) that was blocked by the list (we detect this case by analyzing the referrer header).

Figure 11 provides an overview of third party requests blocked by filter lists or detected as tracking requests according to BehaviorTrack. Out of all 4,216,454 third party requests in the *full dataset*, 2,558,921 (60.7%) requests were blocked by EL&EP, 2,757,903 (65.4%) were blocked by Disconnect, and 2,724,020 (64.6%) were detected as performing tracking by BehaviorTrack.

Requests blocked only by filter lists: Figure 11 shows that EL&EP block 661,523 (15.69% out of 4,216,454) requests that were not detected as performing tracking by BehaviorTrack. These requests originate from 2,121 unique third party domain. Disconnect blocks 720,977 (17.10%) requests not detected by BehaviorTrack. These requests originate from 1,754 distinct third party domains.

These requests are missed by BehaviorTrack because they do not contain any identifier cookie. Such requests may contain other non user-specific cookies (identical across two machines, see Sec. 2.2), however we assume that such cookies are not used for tracking. EL&EP and Disconnect block these requests most likely because they are known for providing analytics or advertising services, or because they perform other types of track-

Filter list(s)	# missed requests	% of 4.2M third-party requests	% of 2.7M tracking requests	# domains of missed requests
EL&EP	826,622	19.60%	25.22%	5,136
Disconnect	687,094	16.30%	30.34%	6,189

Table 7. Overview of third-party requests missed by the filter lists and detected as tracking by BehaviorTrack.

ing through scripts such as fingerprinting which is out of the scope of our study.

5.1 Tracking missed by the filter lists

Table 7 gives an overview of third-party requests missed by EL&EP and Disconnect filter lists and detected by BehaviorTrack as performing tracking. The number of third party domains involved in tracking detected only by BehaviorTrack (e.g., 6,189 for Disconnect) is significantly higher than those only detected by filter lists (e.g., 1,754 for Disconnect reported earlier in this section). BehaviorTrack detects all kind of trackers including the less popular ones that are under the bar of detection of filter list. Because less popular trackers are less prevalent, they generate fewer requests and therefore remain unnoticed by filter lists. This is the reason why we detect a large fraction of domains responsible for tracking.

By further analyzing the requests only detected by BehaviorTrack and missed by EL&EP, we found that 118,314 requests (14.31% of the requests detected only by BehaviorTrack) are false positives. The identifier cookies used in these requests are set by requests blocked by EL&EP. As a result, by simulating the blocking behavior of EL&EP, these cookies should not be included in the analysis. Similarly, we found that 46,285 requests (6.73% of the requests detected only by BehaviorTrack) missed by Disconnect are false positives. We exclude these requests from the following analysis and we further analyze the remaining 708,308 requests missed by EL&EP and the 640,809 missed by Disconnect.

5.1.1 Tracking enabled by useful content

We analyzed the type of content provided by the remaining tracking requests. Table 8 presents the top content types used for tracking and not blocked by the filter lists. We refer to images with dimensions larger than 50×50 pixels as Big images. These kind of images, text, font and even stylesheet are used for tracking. The use of these types of contents is essential for the proper func-

Content type	Missed by EL&EP	Missed by Disconnect
script	33.38%	35.27 %
big images	20.62%	21.73 %
text/html	13.77%	14.73 %
font	8.79%	0.09 %
invisible images	6.68%	12.21 %
stylesheet	6.17%	3.05 %
application/json	4.00%	4.83 %
others	6.59%	8.12%

Table 8. Top content type detected by BehaviorTrack and not by filter lists on the 708,308 requests missed by EL&EP and the 640,809 missed by Disconnect

Service category	EL&EP	Disconnect
Content Servers	23.33 %	23.33 %
Social Networking	16.67 %	0.00%
Web Ads/Analytics	13.33 %	23.33 %
Search Engines/Portals	13.33 %	23.33 %
Technology/Internet	13.33 %	10.00 %
Consent frameworks	3.33 %	3.33 %
Travel	3.33 %	3.33 %
Non Viewable/Infrastructure	3.33 %	0.00%
Shopping	3.33 %	3.33 %
Business/Economy	3.33 %	6.67 %
Audio/Video Clips	3.33 %	0.00%
Suspicious	0.00%	3.33 %

Table 9. Categories of the top 30 tracking services detected by BehaviorTrack and missed by the filter lists.

tioning of the website. That makes the blocking of responsible requests by the filter lists impossible. In fact, the lists are explicitly allowing content from some of these trackers to avoid the breakage of the website as it’s the case for `cse.google.com`.

We categorized the top 30 third party services³ not blocked by the filter lists but detected by BehaviorTrack as performing tracking using Symantec’s WebPulse Site Review [13]. Note that new domains such as `consensu.org` are not categorized properly so we manually added a new category called “Consent frameworks” to our categorization for such services.

Table 9 represents the results of this categorization. Web Ads/Analytics represents 13.33% of the services missed by EL&EP and 23.33% of those missed by Disconnect. However, the remaining services are mainly categorized as content servers, search engines and other functional categories. They are tracking the user but not blocked by the lists not to break the websites.

³ Notice that differently from previous sections, where we analyzed the 2nd-level TLD, such as `google.com`, here we report on full domain names, such as `cse.google.com` that give us more information about the service provided.

5.1.2 Why useful content is tracking the user

Tracking enabled by a first party cookie: A cookie set in the first party context can be considered as a third party cookie in a different context. For example, `site.com` cookie is a first party cookie when the user is visiting `site.com`, but it becomes a third party when the user is visiting a different website that includes content from `site.com`. Whenever a request is sent to domain, say `site.com`, the browser automatically attaches all the cookies that are labeled with `site.com` to this request.

For example, when a user visits `google.com`, a first party identifier cookie is set. Later on, when a user visits `w3school.com`, a request is sent to the service `cse.google.com` (Custom Search Engine by Google). Along with the request, `google`'s identifier cookie is sent to `cse.google.com`. The filter list cannot block such a request, and is incapable of removing the first party tracking cookies from it. In our example, filter lists do not block the requests sent to `cse.google.com` on 329 different websites. In fact blocking `cse.google.com` breaks the functionality of the website. Consequently, an identifier cookie is sent to the `cse.google.com` allowing it to track the user across websites.

By analyzing the requests missed by the lists, we found that this behavior explains a significant amount of missed requests: 44.61% requests (316,008 out of 708,308) missed by EL&EP and 32.00% requests (205,088 out of 640,809) missed by Disconnect contain cookies initially set in a first party context.

Tracking enabled by large scope cookies. A cookie set with a 2^{nd} -level TLD domain, can be accessed by all its subdomains. For example, a third party `sub.tracker.com` sets a cookie in the user browser with `tracker.com` as its domain. The browser sends this cookie to another subdomain of `tracker.com` whenever a request to that subdomain is made. As a result of this practice, the identifier cookie set by a tracking subdomain with 2^{nd} -level TLD domain, is sent to all other subdomains even the ones serving useful content.

Large scope cookies are extremely prevalent among requests missed by the filter lists. By analyzing the requests missed by the lists, we found that 77.08% out of 22,606 third-party cookies used in the requests missed by EL&EP and 75.41% out of 24,934 cookies used in requests missed by Disconnect were set with a 2^{nd} -level TLD domain (such as `tracker.com`).

Tracking behavior	Prevalence
Basic tracking	83.90%
Basic tracking initiated by a tracker	13.50%
First to third party cookie syncing	1.42%
Analytics	1.00%
Third to third party cookie syncing	0.09%
Third party cookie forwarding	0.08%

Table 10. Distribution of tracking behaviors in the 379,245 requests missed by EL&EP and Disconnect.

5.2 Panorama of missed trackers

To compare effectiveness of all lists EL&EP and Disconnect combined, we compare requests blocked by these filter lists with requests detected by BehaviorTrack as tracking according to classification from Figure 3. These results are based on the dataset of 4,216,454 third-party requests collected from 84,658 pages of 8,744 domains.

Overall, 379,245 requests originating from 9,342 services (full third-party domains) detected by BehaviorTrack are not blocked by EL&EP and Disconnect. Yet these requests are performing at least one type of tracking, they represent 9.00% of all 4,2M third-party requests and appear in 68.70% of websites.

We have detected that the 379,245 requests detected by BehaviorTrack perform at least one of the tracking behaviors presented in Figure 3. Table 10 shows the distribution of tracking behaviors detected by BehaviorTrack. We notice that the most privacy-violating behavior that includes setting, sending or syncing third-party cookies is represented by the basic tracking that is present in (83.90%) of missed requests.

Table 11 presents top 15 domains detected as trackers and not blocked by the filter lists. We present the domain's category, owners and country of registration that we extracted using whois library [44] and complement it with manual search. We also manually analyzed all the cookies associated to tracking, and report examples of the cookie's name and expiration date.

Out of the 15 presented domains, 7 are tracking the user using first party cookies, These domains use persistent cookies. The search engine Baidu uses a cookie that expires within 68 years while Qualtrics, an experience management company, uses a cookie that expires in 100 years.

We find that the content from `code.jquery.com`, `s3.amazonaws.com` and `cse.google.com` are explicitly allowed by the filter lists on a list of predefined first-party websites to avoid the breakage of these websites. We identified `static.quantcast.mgr.consensu.org` by IAB Europe that rightfully should not be blocked because they provide useful functionality for GDPR com-

Tracking enabled by a first party cookie						
Full domain	Prevalence of tracking in first-parties	Cookie name	Cookie expiration	Category	Company	Country
code.jquery.com	756 (8.65 %)	__cfduid	1 years	Technology/Internet	jQuery Foundation	US
s3.amazonaws.com	412 (4.71 %)	s_fid	5 years	Content Servers	Amazon	US
ampcid.google.com	282 (3.23 %)	NID	6 months	Search Engines	Google LLC	US
cse.google.com	307 (3.51 %)	NID	1 year	Search Engines	Google LLC	US
use.fontawesome.com	221 (2.53 %)	__stripe_mid	1 years	Technology/Internet	WhoisGuard Protected	—
siteintercept.qualtrics.com	99 (1.13 %)	t_uid	100 years	Business/Economy	Qualtrics, LLC	US
push.zhanzhang.baidu.com	98 (1.12 %)	BAIDUID	68 years	Search Engines	Beijing Baidu Netcom Science Technology Co., Ltd.	CN
Tracking enabled in a third party context						
assets.adobedtm.com	427 (4.88 %)	_gd_visitor	20 years	Technology/Internet	Adobe Inc.	US
yastatic.net	303 (3.47 %)	cto_lwid	1 year	Technology/Internet	Yandex N.V.	RU
s.sspqns.com	278 (3.18 %)	tuuid	6 months	Web Ads/Analytics	HI-MEDIA	FR
tags.tiqcdn.com	276 (3.16 %)	utag_main	1 year	Content Servers	Tealium Inc	US
cdnjs.cloudflare.com	206 (2.36 %)	__cfduid	1 year	Content Servers	Cloudflare	US
static.quantcast.mgr.consensu.org	157 (1.80 %)	_cmpQc3pChkKey	Session	Consent frameworks	IAB Europe	BE
a.twiago.com	133 (1.52 %)	deuxesse_uxid	1 month	Office/Business Applications	REDACTED FOR PRIVACY	—
g.alicdn.com	121 (1.38 %)	_uab_collina	10 years	Content Servers	Alibaba Cloud Computing Ltd.	CN

Table 11. Top 15 domains missed by EL&EP and Disconnect but detected by BehaviorTrack to perform tracking.

pliance. We detect that the cookie values seemed to be unique identifiers, but are set without expiration date, which means such cookies will get deleted when the user closes her browser. Nevertheless, it is known that users rarely close browsers, and more importantly, it is unclear why a consent framework system sets identifier cookies even before the user clicks on the consent button (remember that we did not emit any user behavior, like clicking on buttons or links during our crawls).

We identified tag managers – these tools are designed to help Web developers to manage marketing and tracking tags on their websites and should not be blocked not to break the functionality of the website. We detected that two such managers, tags.tiqcdn.com by Tealium and assets.adobedtm.com by Adobe track users cross-sites, but both of them have an explicit exception in EasyList.

6 Discussion

Domains can at the same time provide useful features to the website, and also track users. We have shown that domains serving functional content such as Content Servers or Search Engines may track the user. However, these domains are not blocked because the web-

site functionality will break otherwise. The functional content doesn't have to set the identifier cookie to perform tracking. It could simply be receiving cookies set by other services from the same domain. In fact, a service can choose to set a cookie with the 2nd-level domain as host, which makes the cookie accessible by all subdomains. Even if the tracking is not intentional, and the domain is not using the identifier cookie it receives to create user's profile currently, these cookie leakage is still a privacy concern that could be exploited by the service anytime. Hence, browser's cookie policy should be revised in order to limit the scope of the cookies.

First party cookies can be exploited to perform cross site tracking. Other browsers should get inspire from Safari and prohibit cookies from being used in a third party context. In its Intelligent Tracking Prevention 2.0 introduced in 2018, Safari allowed cookies to be used in a third party context only in the initial 24 hours. Again the 24 hours rule may be exploited by trackers to keep tracking the user. For example, whenever the user go through google's website the cookie can be updated, hence, google can keep using the cookie from the third party context for a longer period of time.

We have noticed that some domains are not performing tracking themselves. but they redirect or include other domains that perform tracking. We refer to such domains as *tracking initiators* Tracking initiators

do not necessarily include trackers intentionally, whenever a third party service is included, such service has the capability to track the user. However, such behavior has strong impact on the publisher’s liability (the publisher is only liable for the inclusion of tracking initiator) and on the users’ privacy (included trackers are capable of recreating user’s browsing history). In the full dataset, we have identified tracking initiators on 11.24% of websites. Such initiators redirect or include trackers that perform at least one of the cross-site tracking behaviors shown in Figure 3. In our dataset, the most popular tracking initiator is `translate.googleapis.com` that redirects the requests to `gstatic.com`.

To sum up, tracking detection is a complex task and preventing this tracking is even more complex. Blocking domains using filter lists such as Disconnect or EL&EP is not efficient. These lists are widely used in research literature, in browsers and ad-blocking extensions for tracker detection, but they lead to both; false positives, and false negatives. Faced with the choice between protecting one’s privacy or keeping the functionality, we clearly need a more fine-grained tracking approach.

7 Related Work

In this section we first overview previous works on measuring invisible pixels. We then examine state-of-the-art techniques to detect online tracking: behavior-based techniques and methods leveraging the filter lists.

7.1 Invisible pixels, known as web bugs

Invisible pixels are extensively studied starting from 2001 [5, 17, 32, 34, 40]. Invisible pixels, called “web bugs” in previous works, were primarily used to set and send third-party cookies that would come along the request or response when the browser fetches such image. In 2003, Martin et al. [32] studied web bugs on two different lists of websites – the 84 most popular websites and 289 random websites – and found that 58% of the popular websites and 36% of the random websites contain at least one web bug. To help enhancing privacy, in 2002, Alsaïd and Martin [5] deployed a tool (*Bugnosis*) to detect the web bugs. The main goal of the tool was to raise awareness among the public. At that time, the tool helped people to understand the privacy impact and the threat of these web bugs. In fact, the tool was used by more than 100,000 users. However, *Bugnosis*

was only generating warning messages without actively blocking the bugs. Moreover, it was only supported by Internet Explorer 5 that is deprecated today.

Dobias[17] showed that web bugs lead to new privacy threats, such as fingerprinting. Ruohonen and Leppänen [40], studied the presence of invisible pixels in Alexa’s top 500 websites. They showed that invisible pixels are still widely used. Differently from our work, where we detect all effectively delivered images from the response headers, the authors analyze the source code of landing HTML page and extract images from the `` tag. Compare to their work, we crawled the first 500 Alexa’s top websites with all the links, creating a dataset similar to the original work [40]. We have detected 2,698,177 images from which 746,974 are invisible pixels (27.68%). Ruohonen and Leppänen instead extracted 30,572 images from a similar dataset, from which they detected 324 invisible pixels (1.05%). This comparison shows that the detection by the `` tag indeed misses a significant number of invisible images. The significant number of studies on invisible pixels shows that it is a well known problem. The goal of our study is different: *we aim to use invisible pixels that are still widely present on the Web to detect different tracking behaviors and collaborations*.

7.2 Detection of online tracking

Detection of trackers by analysing behavior: Roesner et al. [39] and Lerner et al. [29] were the first to analyze trackers based on their behavior. They have proposed a classification of tracking behaviors that make a distinction between analytics and cross-domain tracking. With respect to this previous work, we propose a more *fine-grained classification of tracking behaviors* that includes not only previously known behaviors, but also specific categories of cookie sharing and syncing (see more details in Section 3). Yu et al.[45], identify trackers by detecting unsafe data without taking into account the behavior of the third party domain and the communications between trackers.

Previous studies [1, 10, 23, 36] were measuring cookie syncing on thousands of websites. Olejnik et al. [36] considers as identifiers sufficiently long cookies. If the value of these cookies is shared between domains, such sharing is classified as cookie syncing. Additionally, they studied the particular case of `doubleclick` to detect cookie matching based on the URL patterns.

In our study, we show that domains are using more complex techniques to store and share the identifier

cookie. We build our technique based on the studies of Acar et al. [1], and Englehardt and Narayanan [23], where they show that companies may use various techniques to store the identifier in the cookie. However, they only check for the identifiers that are stored and shared in a clear text. We further detect synchronization of encoded cookies and even encrypted ones in the case of doubleclick.net. Bashir et al. [10], used ads to detect cookie syncing. The authors filter out all images with dimensions lower than 50×50 pixels which excludes invisible pixels that are the basis of our work. We show that cookie syncing is used with invisible images as well.

Detection of trackers with filter lists: To detect domains related to tracking or advertisement, most of the previous studies [9–11, 22, 23, 26–28, 38] rely on filter lists, such as EasyList [19] and EasyPrivacy [20] (EL&EP). EL&EP became *de facto* approach to detect trackers in security, privacy, and web measurement community. Only from the latest three years venues, we identified *nine papers that rely on EL&EP* to detect or block third-party tracking and advertising (see Table 12 in Appendix for details).

Englehardt and Narayanan [23] seminal work on measuring trackers on 1 million top Alexa websites relies on EL&EP as a ground truth to detect requests to trackers and ad-related domains. This work has set up a methodology to detect trackers that then have been used in the number of follow-up studies.

Three papers by Bashir et al. [9–11] customize EL&EP to detect 2^{nd} -level domains of tracking and ad companies: to eliminate false positives, a domain is considered only if it appears more than 10% of the time in the dataset. Lauinger et al. [28] use EL&EP to identify advertising and tracking content in order to detect what content has included outdated and vulnerable JavaScript libraries in Web applications. Razaghpanah et al. [38] use EasyList as an input to their classifier to identify advertising and tracking domains in Web and mobile ecosystems. Ikram et al. [26] analysed how many tracking JavaScript libraries are blocked by EL&EP lists based on 95 websites. They have found that EasyList is able to block 44% of tracking scripts. Englehardt et al. [22] apply EL&EP lists on third-party leaks caused by invisible images in emails. Iordanou et al. [27] rely on EL&EP as a ground truth for detecting ad- and tracking-related third party requests. Only one work by Papadopoulos et al. [37] use Disconnect list [16] to detect tracking domains.

To the best of our knowledge, *we are the first who compares the behavior-based detection method to filter lists extensively used in the literature.*

Effectiveness of filter lists: Merzdovnik et al [33] studied the effectiveness of the most popular tracking blocking extensions. They evaluate how many third party requests are blocked by each extension. In their evaluation, they don’t distinguish tracking third party requests from non tracking ones, which affect their evaluation. In our work, we detect trackers using a behavior-based detection method and then we evaluate how much of these trackers are blocked.

Das et.al [14], studied the effectiveness of filter lists against tracking scripts that misuse sensors on mobile. They show that filter lists fail to block the scripts that access the sensors. We instead evaluate effectiveness of filter lists against third party requests in web applications that contain identifier cookie

8 Conclusion

Web tracking remains an important problem for privacy of Web users. Even after the General Data Protection Regulation (GDPR) came in force in May 2018, third party companies continue tracking users with various sophisticated techniques based on cookies without their consent. According to our study, 91.92% of websites incorporate at least one type of cookie-based tracking.

In this paper, we define a new classification of Web tracking behaviors, thanks to a large scale study of invisible pixels collected from 84,658 webpages. This invisible images are frequently used in the web: they are present on more than 92.85% of the crawled webpages. We then applied our classification to the full dataset which allowed us to uncover different relationships between domains. The redirection process and the different behaviors that a domain can adopt are an evidence of the complexity of these relationships. We show that even the most popular consumer protection lists fail to detect these complex behaviors. We find out that the browser extensions based on EasyList and EasyPrivacy and Disconnect respectively miss 25.22% and 30.34% of tracking requests we detect. Therefore, these consumer protection lists should not be considered as ground truth to identify trackers, and industry should complement the usage of the lists with detection of trackers based on their behavior.

References

- [1] Gunes Acar, Christian Eubank, Steven Englehardt, Marc Juárez, Arvind Narayanan, and Claudia Díaz. The web never forgets: Persistent tracking mechanisms in the wild. In *Proceedings of the ACM SIGSAC Conference on Computer and Communications Security*, pages 674–689, 2014.
- [2] Gunes Acar, Marc Juárez, Nick Nikiforakis, Claudia Díaz, Seda F. Gürses, Frank Piessens, and Bart Preneel. Fp-detective: dusting the web for fingerprinters. In *2013 ACM SIGSAC Conference on Computer and Communications Security (CCS'13)*, pages 1129–1140, 2013.
- [3] Adblock list parser. <https://github.com/scrapinghub/adblockparser>.
- [4] Adblock Plus Official website. <https://adblockplus.org/>.
- [5] David Martin Adil Alsaid. Detecting web bugs with bugnosis: Privacy advocacy through education. In *Privacy Enhancing Technologies*, 2002.
- [6] Alexa official website. <https://www.alexa.com/>.
- [7] Google-analytics service. <https://developers.google.com/analytics/devguides/collection/protocol/v1/devguide>.
- [8] Mika D Ayenson, Dietrich James Wambach, Ashkan Soltani, Nathan Good, and Chris Jay Hoofnagle. Flash cookies and privacy ii: Now with html5 and etag respawning. Technical report, Available at SSRN: <https://ssrn.com/abstract=1898390> or <http://dx.doi.org/10.2139/ssrn.1898390>, 2011.
- [9] Muhammad Ahmad Bashir, Sajjad Arshad, Engin Kirda, William K. Robertson, and Christo Wilson. How tracking companies circumvented ad blockers using websockets. In *Internet Measurement Conference 2018*, pages 471–477, 2018.
- [10] Muhammad Ahmad Bashir, Sajjad Arshad, William Robertson, and Christo Wilson. Tracing Information Flows Between Ad Exchanges Using Retargeted Ads. In *Proceedings of the 25th USENIX Security Symposium*, 2016.
- [11] Muhammad Ahmad Bashir and Christo Wilson. Diffusion of User Tracking Data in the Online Advertising Ecosystem. In *Proceedings on Privacy Enhancing Technologies (PETS 2018)*, 2018.
- [12] Yinzhi Cao, Song Li, and Erik Wijmans. (cross-)browser fingerprinting via os and hardware level features. In *24th Annual Network and Distributed System Security Symposium, NDSS 2017, San Diego, California, USA, 26 February - 1 March, 2017*, 2017.
- [13] Symantec categorization service. <http://sitereview.bluecoat.com/#/>.
- [14] Anupam Das, Gunes Acar, Nikita Borisov, and Amogh Pradeep. The web's sixth sense: A study of scripts accessing smartphone sensors. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS 2018, Toronto, ON, Canada, October 15-19, 2018*, 2018.
- [15] Disconnect Official website. <https://disconnect.me/>.
- [16] Disconnect List. <https://disconnect.me/trackerprotection/blocked>.
- [17] Jaromir Dobias. Privacy effects of web bugs amplified by web 2.0. In *Privacy and Identity Management for Life - 6th IFIP WG 9.2, 9.6/11.7, 11.4, 11.6/PrimeLife International Summer School, Helsingborg, Sweden, August 2-6, 2010, Revised Selected Papers*, pages 244–257, 2010.
- [18] Doubleclick cookie syncing. <https://developers.google.com/ad-exchange/rtb/cookie-guide>.
- [19] EasyList filter lists. <https://easylist.to/>.
- [20] EasyPrivacy filter lists. <https://easylist.to/easylist/easyprivacy.txt>.
- [21] Peter Eckersley. How Unique is Your Web Browser? In *Proceedings of the 10th International Conference on Privacy Enhancing Technologies, PETS'10*, pages 1–18. Springer-Verlag, 2010.
- [22] Steven Englehardt, Jeffrey Han, and Arvind Narayanan. I never signed up for this! privacy implications of email tracking. In *Privacy Enhancing Technologies*, 2018.
- [23] Steven Englehardt and Arvind Narayanan. Online tracking: A 1-million-site measurement and analysis. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security ACM CCS*, pages 1388–1401, 2016.
- [24] Steven Englehardt, Dillon Reisman, Christian Eubank, Peter Zimmerman, Jonathan Mayer, Arvind Narayanan, and Edward W. Felten. Cookies that give you away: The surveillance implications of web tracking. In *Proceedings of WWW 2015*, pages 289–299, 2015.
- [25] The new Firefox. Fast for good. <https://www.mozilla.org/en-US/firefox/new/>.
- [26] Muhammad Ikram, Hassan Jameel Asghar, Mohamed Ali Kaafar, Anirban Mahanti, and Balachandar Krishnamurthy. Towards seamless tracking-free web: Improved detection of trackers via one-class learning. In *Privacy Enhancing Technologies*, 2017.
- [27] Costas Iordanou, Georgios Smaragdakis, Ingmar Poese, and Nikolaos Laoutaris. Tracing cross border web tracking. In *ACM Internet Measurement Conference (IMC)*, 2018.
- [28] Tobias Lauinger, Abdelberi Chaabane, Sajjad Arshad, William Robertson, Christo Wilson, and Engin Kirda. Thou shalt not depend on me: Analysing the use of outdated javascript libraries on the web. In *Network and Distributed System Security Symposium, NDSS*, 2017.
- [29] Adam Lerner, Anna Kornfeld Simpson, Tadayoshi Kohno, and Franziska Roesner. Internet jones and the raiders of the lost trackers: An archaeological study of web tracking from 1996 to 2016. In *25th USENIX Security Symposium (USENIX Security 16)*. USENIX Association, 2016.
- [30] Timothy Libert, Lucas Graves, and Rasmus Kleis Nielsen. Changes in third-party content on european news websites after gdpr., 2018. https://timlibert.me/pdf/Libert_et_al-2018-Changes_in_Third-Party_Content_on_EU_News_After_GDPR.pdf.
- [31] Timothy Libert and Rasmus Kleis Nielsen. Third-party web content on eu news sites: Potential challenges and paths to privacy improvement, 2018. https://timlibert.me/pdf/Libert_Nielsen-2018-Third_Party_Content_EU_News_GDPR.pdf.
- [32] David Martin, Hailin Wu, and Adil Alsaid. Hidden surveillance by web sites: Web bugs in contemporary use. 2003.
- [33] Georg Merzdovnik, Markus Huber, Damjan Buhov, Nick Nikiforakis, Sebastian Neuner, Martin Schmiedecker, and Edgar Weippl. Block me if you can: A large-scale study of tracker-blocking tools. In *2nd IEEE European Symposium on Security and Privacy*, Paris, France, 2017. To appear.

- [34] Steve Nichols. Big brother is watching: An update on web bugs. In *SANS Institute*, 2001.
- [35] Nick Nikiforakis, Alexandros Kapravelos, Wouter Joosen, Christopher Kruegel, Frank Piessens, and Giovanni Vigna. Cookieless monster: Exploring the ecosystem of web-based device fingerprinting. In *IEEE Symposium on Security and Privacy, SP 2013*, pages 541–555, 2013.
- [36] Lukasz Olejnik, Minh-Dung Tran, and Claude Castelluccia. Selling off user privacy at auction. In *21st Annual Network and Distributed System Security Symposium, NDSS 2014, San Diego, California, USA, February 23-26, 2014*, 2014.
- [37] Panagiotis Papadopoulos, Pablo Rodríguez Rodríguez, Nicolas Kourtellis, and Nikolaos Laoutaris. If you are not paying for it, you are the product: how much do advertisers pay to reach you? In *Internet Measurement Conference, IMC*, pages 142–156, 2017.
- [38] Abbas Razaghpahan, Rishab Nithyanand, Narseo Vallina-Rodriguez, Srikanth Sundaresan, Mark Allman, Christian Kreibich, and Phillipa Gill. Apps, trackers, privacy, and regulators: A global study of the mobile tracking ecosystem. In *Network and Distributed System Security Symposium, NDSS*, 2018.
- [39] Franziska Roesner, Tadayoshi Kohno, and David Wetherall. Detecting and defending against third-party tracking on the web. In *Proceedings of the 9th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2012*, pages 155–168, 2012.
- [40] Jukka Ruohonen and Ville Leppänen. Invisible pixels are dead, long live invisible pixels! In *Workshop on Privacy in the Electronic Society, WPES@CCS*, pages 28–32, 2018.
- [41] Same Origin Policy. https://www.w3.org/Security/wiki/Same_Origin_Policy.
- [42] Ashkan Soltani, Shannon Canty, Quentin Mayo, Lauren Thomas, and Chris Jay Hoofnagle. Flash cookies and privacy. In *AAAI Spring Symposium: Intelligent Information Privacy Management*, 2010.
- [43] uBlock Origin - An efficient blocker for Chromium and Firefox. Fast and lean. <https://github.com/gorhill/uBlock>.
- [44] Whois library. <https://pypi.org/project/whois/>.
- [45] Zhonghao Yu, Sam Macbeth, Konark Modi, and Josep M. Pujol. Tracking the trackers. In *International Conference on World Wide Web, WWW*, pages 121–132, 2016.

9 Appendix

9.1 Detecting identifier sharing

GA sharing: Google-analytics serves invisible pixels on 69.89% of crawled domains as we show in Figure 12. By analyzing our data, we detect that the cookie set by google-analytics script is of the form GAX.Y.Z.C, while the *identifier cookies* sent in the parameter value to google-analytics is actually Z.C. This case is not detected by the previous cookie syncing detection techniques for two reasons. First, "." is not considered as

a delimiter. Second, even if it was considered as a delimiter, it would create a set of values {GAX, Y, Z, C} which are still different than the real value Z.C used as an identifier by google-analytics.

Base64 sharing: When a domain wants to share its *identifier cookie* with doubleclick.net, it should encode it in base64 before sending [18]. For example, when adnxs.com sends a request to doubleclick.net, it includes a random string into a URL parameter. This string is the base64 encoding of the value of the cookie set by adnxs.com in the user's browser.

Encrypted sharing: When doubleclick.net wants to share its *identifier cookie* with some other domain, it encrypts the cookie before sending, which makes it impossible to detect. Instead we rely on the semantic set by doubleclick to share this identifier that we extract from its documentation [18].

Assume that doubleclick.net is willing to share an identifier cookie with adnxs.com. To do so, Doubleclick requires that the content of adnxs.com includes an image tag, pointing to a RL that contains doubleclick.net as destination and a parameter *google_nid*. The value of *google_nid* will tell Doubleclick that adnxs.com was the initiator of this request. Upon receiving such request, doubleclick.net sends a redirection response pointing to a URL that contains adnxs.com as destination with encrypted doubleclick.net's cookies in the parameters. When the browser receives this response, it redirects to adnxs.com, who now receives encrypted doubleclick.net's cookie.

We detect such behavior by detecting requests to doubleclick.net with *google_nid* parameter and analysing the following redirection. If we notice that the redirection is set to a concrete domain, for example adnxs.com, we conclude that doubleclick.net has shared its cookie with this domain.

9.2 Figures

Table 12 summarizes the usage of EL&EP lists in the previous works that we describe in Section 7.

Figure 12 represents the Top 20 domains involved in invisible pixels inclusion in the 8,744 domains.

Table 5 represents the Top 10 third parties forwarding cookies.

Table 12. Usage of EL&EP lists in security, privacy and web measurement community (venues form 2016-2018). “Detection” describes how EL&EP was used to detect trackers: whether the filterlists were applied only on all requests, on requests and follow-up requests that would be blocked, or whether filterlists were further customised before being applied to the dataset. “Dependency” describes whether the results of the paper rely on EL&EP or authors use these lists to only verify their results.

Paper	Venue	EasyList	EasyPrivacy	Detection	Dependency
Englehardt and Narayanan [23]	ACM CCS 2016	✓	✓	Req.	Rely
Bashir et al. [10]	USENIX Security 2016	✓		Custom.	Rely
Lauinger et al. [28]	NDSS 2017	✓	✓	Req.+Follow	Rely
Razaghpanah et al. [38]	NDSS 2018	✓		Custom.	Rely
Ikram et al. [26]	PETs 2017	✓		Req.+Follow	Verif.
Englehardt et al.[22]	PETs 2018	✓	✓	Req.+Follow	Verif.
Bashir and Wilson [11]	PETs 2018	✓	✓	Custom.	Rely+Verif.
Bashir et al.[9]	IMC 2018	✓	✓	Custom.	Rely
Iordanou et al.[27]	IMC 2018	✓	✓	Req.+Follow	Rely

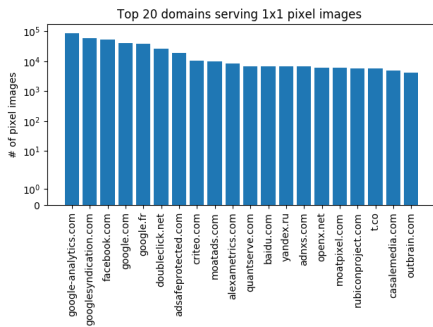


Fig. 12. Top 20 domains responsible for serving invisible pixels

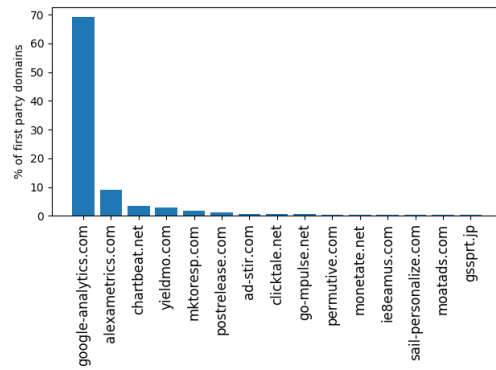


Fig. 14. Analytics: Top 15 receivers in the 8,744 domains.

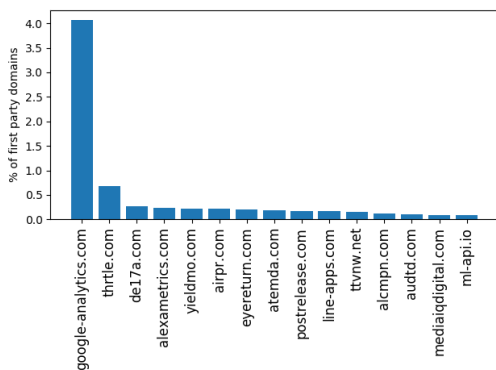


Fig. 13. Third party cookie forwarding : Top 15 receivers in 8,744 domains.