

Anonymous Submission

Who Knows I Like Jelly Beans? An Investigation Into Search Privacy

Abstract: Internal site search is an integral part of how users navigate modern sites, from restaurant reservations to house hunting to searching for medical solutions. Search terms on these sites may contain sensitive information such as location, medical information, or sexual preferences; when further coupled with a user's IP address or a browser's user agent string, this information can become very specific, and in some cases possibly identifying.

In this paper, we measure the various ways that websites send search terms to third parties when a user submits a search query. We developed a reliable methodology for identifying and interacting with search components, which we implemented on top of an instrumented headless browser. We used this crawler to visit the Tranco top one million websites and analyzed search term leakage across three vectors: URL query parameters, payloads, and the Referer HTTP header. Our crawler found that 512,701 of the top 1 million sites had internal site search, where it successfully interacted with 92.1% of these and captured over 68 million web requests. We found that 81.3% of websites containing internal site search sent our search terms to third parties in some form. We then compared our results to the expected results based on a natural language analysis of the privacy policies of those websites (where available) and found that about 87% of those privacy policies do not mention search terms explicitly. However, about 75\% of these privacy policies seem to mention the sharing of some information with third-parties in a generic manner, which means that the onus is on the reader to fathom whether such generic wording includes search terms. We then present our countermeasure: a browser extension to warn users about imminent search term leakage to third parties, available for all major browsers. We conclude this paper by making recommendations on clarifying the privacy implications of internal site search to end users.

Keywords: privacy policy, web tracking, measurement DOI Editor to enter DOI

1 Introduction

Internal site search is integral to how users discover and interact with a wide range of web content including shopping, travel planning, medical information, and social network features. Internal search is increasingly being used by businesses to drive revenue, since 30% of e-commerce website visitors use internal site search which contributes to 14% of all revenue [6].

Users may use these search boxes to type in highly personal terms expressing racial identity, sexual or religious preferences, and medical conditions. Prior work has shown how easy it is to de-anonymize users based on their search terms [46, 70]. In 2012, the US retailer Target inferred that a young woman was pregnant from her purchasing habits and search history, and started showing her advertising for baby products, before she was aware of her pregnancy herself [4]. Alarmingly, Guha et al. found that advertisers collect and use sensitive data such as sexual preference for ad personalization: for example, nursing college degree ads appeared to target gay users exclusively [43]. In 2006, after AOL published a large amount of anonymized search queries, several participants were completely de-anonymized based on the uniqueness and specificity of their search terms [3]. Therefore, the secrecy of search terms is paramount to any reasonable definition of privacy on the web.

A parallel trend shows the increasing proliferation of third-party scripts and fingerprinting techniques (including persistent cookies, tracking pixels, and other technologies - collectively called *trackers*) on the web, mostly for advertising purposes [73]. A 2014 crawl of the top 1 million websites¹ showed ubiquitous user tracking: 62.9% of sites spawned third-party cookies, 83.0% loaded third-party JavaScript libraries (potentially for browser fingerprinting), and 88% of all sites made a third-party network request in some form [54]. A follow-up study in 2019 found the percentage of sites which track users increased to 90%, despite the introduction of privacy-protecting legislation in the European Union at that time [66].



In this paper, we explore the question of search privacy by visiting the Tranco top 1 million websites [52] and performing an automated internal site search. We developed a crawler on top of a headless browser which can handle dynamically generated websites, identify search inputs, and capture outgoing requests. We then analyzed the content of outgoing requests to determine whether our search terms appeared in the data or metadata. In particular, we examine search term leakage along three vectors:

- 1. **Leakage via** *Referer*². Presence of search terms in the HTTP *Referer* request header sent to third parties.
- 2. **Leakage via URL**. Presence of search terms in third-party URLs and query parameters.
- 3. **Leakage via payload**. Presence of search terms in HTTP request payloads sent to third parties.

We discover that 81.3% of visited websites send search terms to third parties in some form, representing a potential privacy threat to users. Additionally, despite many websites having privacy policies, we found the topic of search term processing was either not explicitly covered by the privacy policies (about 87% of the time) or simply ignored by those websites (we were unable to find a privacy policy for about 50% of sites).

The contributions of this paper are as follows:

- We present the first large-scale measurement study on the leakage of search terms via internal web search, visiting one million websites, intercepting over 68 million network requests, and discovering 512,701 unique domains with an internal search functionality.
- 2. We build a robust methodology to identify search inputs across a range of languages during a large-scale crawl. We use this methodology to collect search input and search-related element selectors for each domain. We then open-source this labeled dataset of discovered search inputs per domain to the community for further study³.
- 3. We perform a sophisticated analysis of both naturallanguage and P3P-specified [58] privacy policies to determine if these privacy policies mention processing search terms or sharing some information with third parties.
- 4. We provide a framework to characterize the various vectors for leaking search terms to third par-

ties. Using this framework, we found that 81.3% of sites send search terms in some form to third parties. By breaking this down according to our vectors, we find that 75.8% of websites send data via the *Referer* header (72.5% directly and 10.6% indirectly); 71% via URL query parameters; and 21.2% via the payload. Additionally, 87.4% of websites had the potential to send sensitive query terms to a new third-party if a link on the page was clicked by the user. Finally, this framework allowed us to identify those third parties that are actively involved in providing search results.

 We develop a browser extension to inform users about potential privacy leaks via internal search, based on our findings, which we release ⁴.

2 Background

2.1 Referer HTTP Header

The Referer HTTP request header is automatically set by the client (browser) to indicate to a website how a user found that website. For example, when clicking on a link to a New York Times article from the Hacker News aggregator website, the New York Times web server would see these HTTP headers by default:

GET /2020/10/15/technology/ignore-phone-companies-5g.html

Host: www.nytimes.com

User-Agent: Mozilla/5.0 XXX Firefox/80.0
Referer: https://news.ycombinator.com/

The Referer header has been a feature in HTTP since at least HTTP 1.0 [1]. This first RFC stipulates that the Referer header may be useful in debugging broken links, cache optimization, and logging. It also contains a warning about the privacy implications, and suggests that users should have the option to control this value. This recommendation is also made in the RFC for HTTP 1.1 [2]. To our knowledge, no major browser allows users to easily configure this value as suggested.

² original spelling

³ https://jellybeans-paper.com/

f 4 under web store review, which should be complete by camera-ready deadline

2.2 Referrer-Policy HTTP Header

The Referrer-Policy response header is set by the website owner, and specifies how the Referer header will be set upon requests originating from that site. The Referrer-Policy header has several acceptable values, which are shown in Table 1^5 .

The referrer policy can be set either in the HTTP header, in a <meta> tag in the header section of the document, or as a noreferrer attribute for a particular link on the site. It may also be set implicitly via inheritance.

As can be seen from this table, a significant amount of attention has been paid to not sending *Referer* values over insecure (non-HTTPS) channels. It was observed by the early writers of the specification that the *Referer* values may contain sensitive information that man-in-the-middle attackers may want to steal. Our concern is however not with man-in-the-middle actors but with third-parties with resources on a first-party web page.

2.3 Leaking Search Terms via the Referer

The Referer header can therefore accidentally (or intentionally) leak search terms to third parties. This attack has been previously documented by Libert [53], Krishnamurthy [50], Malandrino [59], and several others.

A real-world attack scenario is shown in Figure 1⁶. A user performs a search for "pancreatic cancer" on the WebMD website. The WebMD server returns a webpage with the results, which then loads a variety of third-party JavaScript tracking scripts, including one from Google Ad Services. When the request is made to fetch that script from the Google server, the default behaviour (and the one observed in the wild) is sending the URL together with the query string, which contains the sensitive search term. The raw request header values can be seen below⁷:

GET /pagead/conversion_async.js HTTP/2

Host: www.googleadservices.com

User-Agent: Mozilla/5.0 XXX Firefox/80.0

Referer: https://www.webmd.com/search/search_results/

default.aspx?query=pancreatic%20cancer

In this paper, we use the term **search term leakage** to refer to the transmission of search terms to third

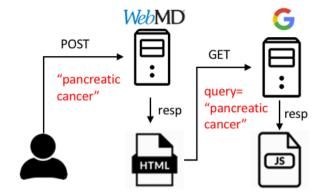


Fig. 1. After searching for "pancreatic cancer" on WebMD, when WebMD loads the Google Ad Service JavaScript file, "pancreatic cancer" appears in the *Referer* header.

parties. We note this transmission might be either intentional or accidental with respect to the website owner.

2.4 What is a Third Party?

What constitutes a first party and a third party can be difficult to define. In this study, we rely on effective Top-Level Domains as defined in the Public Suffix List [23], specifically the eTLD+1, so we consider domains that do not share the same eTLD+1 third parties.

Previous studies have tried to refine this definition by considering a third-party as "a subdomain of the main domain whose IP address is registered to a different autonomous system (AS) than the IP address of the first-party website." [69] While this approach would catch those tracking entities that do not want to be detected, we accept that in some cases our third-party results might be slightly lower that what they might be in reality. However, we think that this might be balanced by the fact that in other cases, two domains that do not share the same eTLD+1 might actually be owned by the same entity. As mentioned by authors from the W3C Privacy Community Group's First-Party Sets work group, "the website the user is interacting with may be deployed across multiple domain names" [14]. So ideally, when identifying third-party requests, we would like to exclude those domains that may be owned by the entity that also owns the site that is being analyzed (e.g. datahug.com and sap.com are part of the same entity after SAP's acquisition of DataHug but nothing in their respective WHOIS records can link one to the other). While this violates the *Origin* concept, it is closer to a user's mental model of what a website is and its associated trust relationship. However, we may

⁵ These values are taken from the RFC [5]

 $[{]f 6}$ WebMD's privacy policy includes language suggesting that search terms may be transmitted to third parties

⁷ some values were redacted to protect the author's privacy

Value	Description
no-referrer	Referer will be empty
no-referrer-when-downgrade (default	sends full Referer header value if both this and the target page are HTTPS otherwise sends
for most browsers)	nothing
same-origin	same origin requests contain full Referer information while no Referer is sent for cross-origin
	requests
origin	only the origin, not path or query information, are sent
strict-origin	same as above, but only for HTTPS pages. Non-HTTPS pages result in no Referer being sent
origin-when-cross-origin	cross-origin requests behave as if $origin$ was the specified policy while same-origin requests
	behave as if $same-origin$ was specified
strict-origin-when-cross-origin	same as above, but only for HTTPS pages. Non-HTTPS pages result in no Referer being sent
(default for Safari and Chrome v85+)	
unsafe-url	Full Referer is sent in all contexts. This is generally advised against

Table 1. Referrer-Policy header valid values, derived from the W3C specification [5]. no-referrer-when-downgrade is the default policy for Firefox, Chrome versions 84 and below, and Edge. strict-origin-when-cross-origin is the default policy for Safari and Chrome versions 85 and higher [27].

not be able to reliably identify such relationships, since domain sets which claim to provide this information are not completely reliable (e.g. WHOIS) or exhaustive (e.g. Disconnect's Tracker Protection lists [30]). So for the purpose of this study, we rely solely on eTLD+1 and define **search term leakage** as any transmission of the search string between different eTLD+1 entities.

2.5 Privacy Policies

To determine whether website operators explicitly inform users about sharing their search queries with third-parties, we examine the privacy policies of websites (where possible). We investigate both machine-readable privacy policies (provided by the Platform for Privacy Preferences (P3P) framework [37]) as well as privacy policies authored using natural language.

While P3P is officially deprecated and no longer supported by browsers, it is a machine-readable format which could help us easily determine whether sharing users' search terms with third parties is covered by a website's privacy policy. According to the P3P W3C Recommendation [58], the interactive data category covers "data actively generated from or reflecting explicit interactions with a service provider through its site, such as queries to a search engine." This definition precisely corresponds to what we want to measure. While P3P policies are extremely convenient to parse, since it is a deprecated standard we expect these to be infrequent.

Since P3P policies are rarely encountered in the wild, natural-language privacy policies are increasingly common. While some jurisdictions do not mandate the

posting of privacy policies, "many individual [US] states do have this type of legislation. Additionally, the countries within the European Union (EU) have also enacted similar legislation that requires websites that collect personal or identifiable data to post privacy policies regarding its use." [16] A recent study focusing on a set of 6,579 popular websites found that 84.5% of these had privacy policies [38].

3 Data Collection

In this section, we describe in detail the crawling environment used to simulate user searches over the TRANCO top 1M domains⁸ [31]. In order to complete this task, we instrumented a headless Chromium-based browser to visit each domain and detect the presence of all search inputs contained on the landing webpage. For each detected search input, the crawler simulates a real user search by typing the keyword "JELLYBEANS" into the search input box one character at a time. While doing so, the crawler simultaneously intercepts all inbound and outbound network traffic, request redirection chains, payloads, and HTML content. We can then search for our dummy search string in the intercepted outbound network request data and metadata to detect search string sharing with third-party entities. To facilitate replication, we have open-sourced our crawler ⁹.

 $^{8\,}$ We used the list created on April 27, 2020. Available at https://tranco-list.eu/list/NYJW.

 $[{]f 9}$ the code for the crawler can be found here: https://jellybeans-paper.com/

3.1 Web Crawler

Considering the trade-offs from the crawler categories minutely discussed in [32], we developed a User Layer based crawling system built on top of the Puppeteer headless browser [24], written in TypeScript, which we call SlicedPie. This architecture was chosen to maximize the overall accuracy, coverage, and completeness levels of the crawling process while allowing for some moderate extensibility. SlicedPie's design facilitated simulating real user's interactions with the browser and enabled (i) driving Chromium as a full-fledged browser capable of handling dynamically generated content; (ii) launching on-demand JavaScript execution tasks; and still (iii) benefiting from lower level access to the Chrome Development Tools (CDP) [9] API (through CDP sessions [25] instantiation).

In terms of computational power, the underlying system was configured with 32 CPU cores running 64 threads and 800 GB RAM running Ubuntu 18.04LTS and Linux kernel version 4.15.0-72-generic. This setup was used to handle more than 100 concurrent worker-thread [20] instances. Each worker-thread was initially assigned to process a single input domain URL, and during its execution process it could run multiple parallelized browser instances.

3.2 Crawling Stages

The crawling process encompasses four interdependent crawling stages. Each stage contains additional intermediate steps where data validation and information extraction tasks are performed in order to reduce error rates, improve performance, and maximize results coverage. Crawls are scheduled and subsequently parsed by a crawler manager (main thread), and the collected data is primarily stored in a NoSQL database, with records larger than 16MB being offloaded to disk.

3.2.1 Preflight Request (Stage I)

In this stage we perform a *preflight* request using Node fetch [11] to confirm the URL's availability. Common system errors such as DNS failures are captured at this stage and therefore prevent further stages from executing. Additionally, we implement a mechanism to handle timeouts to circumvent the fetch API's inability to perform as expected [12].

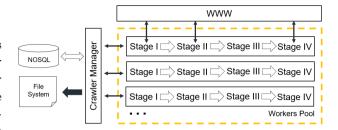


Fig. 2. Crawler Architecture. Stage I (Preflight Request), Stage II (Extracting HTML), Stage III (Detecting Search Functionality), Stage IV (Simulating User Search)

3.2.2 Extracting HTML (Stage II)

During this stage the HTML text is extracted and analysed in order to:

- Detect additional languages. The HTML language attribute [18] is extracted from the HTML text body in an attempt to identify additional languages present in the page structure. Each additional language found (apart from English) is used in *Stage III* to generate additional selectors.
- Detect embedded search functionalities. Websites may contain embedded search capabilities[28]. These share a common structure [17] such that they can be detected using a regular expression and parsed into a search URL like http://example.com/?s={search-term-string}. Thus, as an additional step of this stage we try to identify and extract the search endpoint URL from the HTML string. Using this method, we discover that 8.5% of crawled sites rely solely on this search functionality.
- Collect all hyperlinks found in the page, including those pointing to privacy policies.
 The identification of privacy policy links is described in Section 3.5.

3.2.3 Detecting Search Functionality (Stage III)

In order to detect search functionalities, we first rely on a list of query selectors [26] to identify all elements on a Web page that contain the keyword *search* as part of their attributes. While English is the default language used, the keyword varies according to any additional language detected during stage I. In order to use quality translations for the term *search*, we leveraged translations from Mozilla's Firefox Translation Project [13].

The list of elements is subsequently filtered and bucketed into two unique query selector lists of search inputs and search-related elements. Search inputs represent all elements into which it is possible to enter text. In the search-related inputs category fall all other HTML elements (for example: div, 1i, span, etc.) that cannot be used to perform a search but indicate via their attributes that they are somehow related to search functionality.

3.2.4 Simulating User Search (Stage IV)

All information gathered during previous stages is used here to define a strategy to simulate a search. There are two base approaches to do so:

- 1. **Search Inputs.** If the element is of type input, the approach is to visit the page, locate the element, simulate the user typing the keyword "JELLY-BEANS" followed by the return key.
- 2. **Embedded Search URL.** If an embedded search URL has been identified and extracted during *Stage II*, it is used to conduct a search. In this case the keyword is added as an input parameter while the page is visited using a URL in the format of http://example.com?s=JELLYBEANS.

3.3 Challenges and Edge Cases

The SlicedPie crawler can deal with a number of challenging situations that arise in the wild including dynamically-generated element attributes, interstitials, hidden search inputs, and search results displayed in another tab. The mechanics of our solutions for these are detailed in the Appendix.

3.4 Crawl Summary

Table 2 outlines the results of our crawl. In summary, we tried to crawl 1 million sites but 15.7% of those generated errors preventing us from determining whether the site might contain a search component. For 32.9% of the remaining sites, we were unable to detect any search component. This left us with 512,701 sites, where our search simulations succeeded 92.1% of the time.

3.5 Collecting Privacy Policies

To determine whether users are notified about the leakage of their search terms to third parties, we also col-

Error Code	Description	Domains
PREFLIGHT ERR	Preflight request failures due to DNS, invalid HTTP Status Code (4xx or 5xx), and timeouts.	157,357
NO INPUTS FOUND	No elements suggesting the presence of search functionality on the page were detected.	329,942
INVALID INPUTS	Search was attempted but discovered elements were not interactive (could not input text and simulate user interaction).	40,368
SUCCESS	At least one interactive element was found and used to complete the search simulation.	472,333
Total		1,000,000

Table 2. Crawl Summary.

lected the machine-readable and natural-language privacy policies (where possible) using our instrumented crawler SlicedPie. Machine-readable policies are easily collected using the P3P HTTP header (where present). To collect natural-language policies, we followed the following two-step approach:

- 1. Privacy Policy Link Identification
 - Check page links in reverse order
 - Use English terms
 - Use detected language terms if required
- 2. Privacy Policy Document Retrieval
 - Perform simple GET request
 - Use headless browser if simple GET request fails to retrieve enough content
 - Perform PDF to text conversion if required

Our first step follows the approach proposed in [55]. However, our keywords are not limited to English. While English is the default language used, the keywords vary according to any additional language detected during the HTML extraction presented in Figure 2, once again leveraging translations from Mozilla's Firefox Translation Project [13].

4 Results

In Section 3, we successfully executed a search for our dummy query string, "JELLYBEANS", on 472,333 websites and captured the outgoing network requests with our crawler. In this section, we analyze the collected data and categorize the observed privacy leakages using

the ontology presented in Figure 3, inspired by Starov et al. [69].

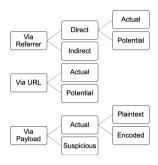


Fig. 3. Types of privacy leakages analyzed

While Starov et al. categorized privacy leakages as intentional or accidental, we think it is quite difficult to be confident whether a leakage is one or the other. Instead, we analyze leakages using several criteria. First, we focus on the mechanism by which the search query was leaked, be it via the *Referer HTTP* request header, via a *URL query parameter* or via the actual *payload* of the HTTP request. We further classify *Referer*-based leakages depending on whether our dummy query string appears in a *Referer* header that was either directly sent by the analyzed site or indirectly by one of the site's third-parties.

Finally, we classify leakages depending on whether they actually occurred during our crawl, or whether they could potentially occur if the crawl's scope was expanded. Since we analyze the document containing the search results, we inspect all links present on that page to determine whether further leakages might occur if these links were followed.

Also, we break down the results in the *payload* category using additional criteria. Since we observed some large request payloads in our dataset, we suspect some of them may contain our search query, but that is not necessarily easy to confirm since some of these payloads are encoded.

All of our results are based on the number of sites where our crawler was able to execute a search without encountering any error. As presented in Table 2, this number is 472,333. It should be noted that two or more sites may redirect to the same domain. While the final number of distinct response domains is 447,428, we report all of our results based on a total of 472,333.

4.1 Leakage via Referer

Our dummy query was leaked via the *Referer* header by 75.8% of the sites, be it directly or indirectly. The example in Listing 1 shows both scenarios when conducting a search on the *www.autoguide.com* site. The first part of the example shows a request to a third party domain (*ib.adnxs.com*) with our search query present in the *Referer* header. Subsequently, our query string appears in the *Referer* header of a request to *sync.placelocal.com* even though this third party domain was never contacted directly by *www.autoguide.com*. This flow is visualized in Figure 4.

```
{'url': 'https://ib.adnxs.com/ut/v3/prebid',
     'resourceType': 'xhr', 'headers':
    {'sec-fetch-mode': 'cors', 'referer':
     'https://www.autoguide.com/search.html?q=JEL
    LYBEANS', 'payload':
     '...'}}
{'url': 'https://sync.placelocal.com/openadid=864
    616440172932000; redirect=https%3A%2F%2Fsync.
    placelocal.com%2Fsyncdatapartnersimg%3F0bypa
    sssync%3D1%26blob%3Dcdc4c4f81e1770430aea8361
    Ofd967e8%253Aef9dc7a0939e092d667bc27fbf9dc57
    060a97d39dfb6f590dd1cb76e2dbfa57e6fc69c7ec6c
    fd67ef623c6836288b87e', 'resourceType':
     'image'. 'headers': {'referer':
     'https://tag.placelocal.com/ad/iframe?...auc
    tionid=1789060303328245535&refurl=https%253A
    %252F%252Fwww.autoguide.com%252Fsearch.html%
    253Fq%253DJELLYBEANS...'}}
```

Listing 1. Examples (shortened) of referrer-based leakages

4.1.1 Direct leakage

72.5% of sites directly send our query string to at least one third party domain via the *Referer* header. This is due to the fact that most sites do not set any pagelevel *Referrer-Policy*, be it using a response header or an HTML directive (using a meta element). The most frequently observed pairs of policy values are shown in Table 3.

It is interesting to note that the most frequently used values are ineffective at preserving users' privacy. no-referrer-when-downgrade is often never relevant, as most of today's traffic is relying on HTTPS and downgrades are uncommon. always is not a valid policy

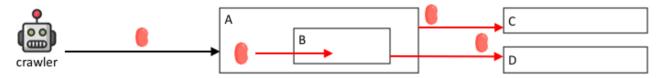


Fig. 4. Crawler searches for "JELLYBEANS" on site A. A then passes this search query to A's child iframe B. A then makes a network request to C and "JELLYBEANS" is leaked to C via the *Referer* header. B makes a request to D, where B inserts "JELLYBEANS" into the URL. Thus in this example "JELLYBEANS" is leaked to three domains (B, C, D) via two different methods (*Referer* header and URL).

Rank	Header Policy	Meta Policy	Frequency
1	None	None	93.6%
2	no-referrer-when-downgrade	None	3.2%
3	(empty string)	None	1.4%
4	None	always	0.6%
5	None	unsafe-url	0.4%
6	unsafe-url	None	0.2%

Table 3. Most frequent page-level *Referrer-Policies*, as specified in the HTTP response header and the HTML meta element.

and unsafe-url is discouraged except in specific usecases, as shown in Table 1.

4.1.2 Indirect leakage

As shown in Listing 1, leaks via *Referer* can also occur indirectly when a third-party domain is contacted by another third-party domain. This situation was observed for 10.6% of the sites, with the most prevalent domains (and their associated entities based on the Disconnect list[30]) shown in Table 4.

This table shows that many of these domains are classified as tracking entities according to the Disconnect list. However, some domains, such as those associated with Google, may be more ambiguous as they might deliver actual search results via a custom search engine (cse.google.com) [15]. This possibility is explored in Section 4.3.

4.1.3 Potential leakage

To conclude this analysis of *Referer*-based leakages, we are interested in examining those links present on search results pages that would leak the search query to a third-party domain, if those links were clicked by the user (if the query would not previously have been leaked to any of these third-party domains). To determine whether a link would leak or not, we analyzed the page-level

Domain	Entity	Frequency
google.com	Google	3.3%
fonts.googleapis.com	None	2.6%
google-analytics.com	Google	1.8%
moatads.com	Moat	1.7%
gstatic.com	Google	1.6%
2mdn.net	Google	1.6%
quantserve.com	Quantcast	1.5%
quantcount.com	None	1.5%
doubleclick.net	Google	1.1%
googletagmanager.com	None	1.0%
www.googleapis.com	None	1.0%
facebook.net	Facebook	0.8%
adsafeprotected.com	Integral Ad Science	0.7%
doubleverify.com	DoubleVerify	0.7%
yahoo.com	Yahoo!	0.7%
iasds01.com	Integral Ad Science	0.7%
flashtalking.com	Flashtalking	0.6%
googlesyndication.com	Google	0.5%
googletagservices.com	Google	0.5%
ajax.googleapis.com	None	0.5%

Table 4. Most frequent domains present in indirect referrers.

referrer policies (as shown in Table 1 but also element-level policies that may be present. We found that links with the potential to leak the search query occurred for 87.4% of sites.

The results in Table 5 show that the most prevalent domains in links on search results pages are linked with social media services (e.g. Facebook, Twitter, Instagram, Youtube, Linkedin and Pinterest). A generic link shortening service is also present in that list (bit.ly).

4.2 Leakage via URL

71% of sites send our search query as part of URL query parameters. As shown in Listing 1, however, a URL query parameter may contain the full search page URL rather than the actual query. Distinguishing these two scenarios seems important because some of our observations (that could be qualified as *leakages*) may be

Domain	Entity	Frequency
twitter.com	Twitter	44.8%
facebook.com	Facebook	38.3%
instagram.com	Facebook	34.2%
youtube.com	Google	29.2%
linkedin.com	LinkedIn	15.8%
pinterest.com	Pinterest	9.0%
google.com	Google	8.8%
wordpress.org	None	5.7%
apple.com	None	4.0%
vk.com	VKontakte	2.0%
t.me	None	1.9%
bit.ly	None	1.8%
vimeo.com	Vimeo	1.5%
flickr.com	Yahoo!	1.5%

Table 5. Most frequent domains from links on search results pages.

absolutely necessary for the search functionality to be delivered to the user. For instance, if a website is relying on a third-party search engine provider (be it Google Custom Search Engine or any other service), our search query may be passed on directly to that third-party using technologies such as AJAX. In order to address this scenario, we try to identify those domains that receive specific query parameters (such as q or query) with our query string (or parts of it) as the value. By observing whether URL parameter values contain substrings of "JELLYBEANS", such as q=J or q=JE (which is the case when services provide auto-complete suggestions), we can reliably identify actual third-party search engine providers. Table 6 shows the most common domains we are able to identify using this strategy. Some of them are associated with Google, but others may not be as well-known. However, manually checking the twenty five most frequent domains confirms that they actually are search services that can be added to websites such as content management systems.

4.3 Leakage via Payload

Finally we are interested in examining whether sites leak search queries with third parties via the HTTP request payload. We find that 21.2% of sites actually do, the most prevalent third parties being listed in Table 7.

This list is quite varied since it contains social media entities (e.g. Facebook, Twitter), advertizers (Federated Media), analytics services (sumo.com, Google Analytics) and actual search service providers (algolia.net) identified using the methodology described in the previ-

Domain	Sites	Frequency
google.com	5063	0.01%
yandex.ru	692	0.0014%
searchanise.com	609	0.0012%
wikipedia.org	453	0.0009%
wiktionary.org	444	0.0009%
ksearchnet.com	284	0.0006%
doofinder.com	260	0.0005%
mybcapps.com	241	0.0005%
nextopiasoftware.com	237	0.0005%
searchspring.net	236	0.0004%
addsearch.com	223	0.0004%
cludo.com	207	0.0004%
searchspring.io	166	0.0003%
swiftype.com	156	0.0003%
resultspage.com	151	0.0003%
algolia.net	18	0.00003%

Table 6. Most frequent domains identified as search results providers after receiving query parts via specific URL query parameter

Entity	_
Littity	Frequency
Facebook	7.1%
Hotjar	3.8%
AppNexus	2.6%
Criteo	1.5%
PubMatic	1.5%
Twitter	1.0%
Federated Media	0.6%
None	0.6%
Google	0.6%
Flashtalking	0.6%
None	0.5%
	Facebook Hotjar AppNexus Criteo PubMatic Twitter Federated Media None Google Flashtalking

Table 7. Most frequent domains receiving query via payload.

ous section. It is also worth highlighting the presence of hotjar.com, which provides session replay capabilities to site owners as reported in [40]. A deeper analysis of the request payloads (as well as the responses) should allow us to identify specific entities (e.g. search providers) as we did in the previous section using specific query parameters.

To conclude this analysis, we report that 0.9% of sites make use of third parties that rely on Base 64-encoded payloads that we were able to decode in order to detect our search query. We also report that some requests contain unusually large payloads (over 10kb per request), most often sent to entities that offer session replay services, like yandex.ru, hotjar.com, fullstory.com, smartlook.cloud, clicktale.net, quantummetric.com, sessioncam.com or inspectlet.com.

4.4 Analysis of Privacy Policies

Now that we have observed how sites leak search queries to third parties, we want to determine whether this is a practice that website operators explicitly mention to their users. In order to measure how often website operators in fact include this type of data leakage in privacy policy documents, we analyze both machine-readable (P3P) and natural-language privacy policies.

4.4.1 Machine readable policies

We find that 4.5% of all sites use the P3P response header to communicate some of their data handling practices. When we exclude sites that do not leak our search query, we still observe 4.4% of sites using this header. Focusing specifically on those sites which mention both interactive data and third parties, however, we find that the number goes down from 4.4% to 1.3%.

4.4.2 Natural language privacy policies

4.4.2.1 Identifying Privacy Policy Documents

Using our approach from Section 3.5, we were able to find privacy policy links on 50.5% of input domains (we only considered those sites leaking our search query to at least one entity). While this number may seem quite low, it is not unexpected given we relied on a fully automatic link identification approach. Previous work has shown that a word-list approach can work for specific site types (such as news) but miss links on other sites, thus requiring a complementary manual identification step [38]. Using such a manual step was not feasible based on the size of our dataset.

Using the identified links, we were able to retrieve some content 96.5% of the time. However, while the keywords used in our first step are generally a good indicator that the linked content is an actual privacy policy, there is no guarantee that this is always the case. For example, some sites rely on intermediate pages to guide users towards a relevant privacy policy if they have more than one (e.g. one per product, or one per platform). So while the link identified might point to a page containing the word *privacy*, there may actually not be any actual privacy policy content in it.

Previous studies have used some heuristics such as keyword spotting to justify their scraping [55], but we think this is too prone to errors. We therefore decided to train a classifier to filter out any content that does not

appear to contain any privacy policy content. To do so, we used the sanitized OPP-115 corpus [76] as training data and relied on a one-class SVM classifier (RBF kernel with gamma set to auto) to perform outlier detection using the Scikit-learn framework [63] whose implementation is based on Libsvm [35]. Using a very small nu value (0.05) to reduce the number of training errors to a minimum, we extracted TF-IDF features (n-grams in the range 1-3) after removing English stopwords from the content. The resulting model detected 4 documents from our training set as outliers, and when applied to all of the 127,996 English policies retrieved detects 37.9% of them as outliers. This might be very conservative but ensures that we are looking for search-related mentions in the appropriate content.

To validate our model, we manually sampled 100 of the 79,515 English documents labelled as privacy policies by our model and found that all of them were indeed privacy policy documents. Our pipeline is almost directly comparable to the one described in [56] which achieved a privacy policy link candidate detection rate of about 43% (ours is 50.5% with duplicates and 44.7% after removing duplicates) and an actual privacy policy document detection rate of about 15% using a CNN (ours is 20.7%). Our results are summarized in Table 8. We also make available all discovered privacy policy links to future researchers ¹⁰.

Description	Count
Input domains	384101
Domains with policy link identified	193815
Link found but scraping error	6703
Link found and scraping success	187112
Uniq. links with document	165139
Uniq. documents detected as English	127996
Uniq. English docs labelled as privacy policies	79515
Uniq. English docs including search-related terms	10914
Uniq. English docs including generic info. sharing	60212

Table 8. Analysis of privacy policies statistics. Only domains with observed search term leakage were used as input domains.

4.4.2.2 Parsing Privacy Policies

To investigate whether websites inform users of the fact that their search queries may be sent to third parties,

¹⁰ https://jellybeans-paper.com

Source	Category	Text	Match type
1582_msn.com.csv	Third Party Sharing and Collection	results page URL (which may contain your search query)	explicit
$1099_enthus ia stnetwork.com.csv$	Third Party Sharing and Collection	the pages you visit, and which links you click, which ads you see and click on, and the categories of search terms you enter	explicit
1703_sports-reference.com.csv	First Party Collection and Use	SRL may record information identifying the visitor or linking the visitor to the search performed .	explicit
105_amazon.com.csv	Third Party Sharing and Collection	search term and search result information from some searches conducted through the Web search features	explicit
1300_bankofamerica.com.csv	Third Party Sharing and Collection	certain information about your activities on our Sites, such as pages visited and search key words entered	explicit
1034_aol.com.csv	Data Retention	Information about your interactions with the websites, apps, and other services you use, the content you view, the search queries you submit,	explicit
1089_freep.com.csv	Third Party Sharing and Collection	We may share technical or aggregate information about your interaction with our Site, such as type of pages viewed and categories of interest, from our Site with these service providers	generic
856_sciencemag.org.csv	Third Party Sharing and Collection	We may share information, such as your IP address, with third parties as might be required for technical purposes	generic
686_military.com.csv	Third Party Sharing and Collection	We may share your information with third parties who help us in the delivery of our own products and services to you	generic
1070_wnep.com.csv	Third Party Sharing and Collection	We may share your information with third-party advertisers and advertising networks or others with which we have a contractual relationship	generic
175_mlb.mlb.com.csv	Third Party Sharing and Collection	A permitted network advertiser may use cookies, web beacons or similar technologies to collect information about your interaction with our Services	generic

Table 9. Examples of relevant OPP-115 corpus fragments.

we must parse the extracted privacy policies and find sections relevant to processing search terms.

Since search queries are instances of personal information (rather than personally identifiable information), we would expect privacy policies to mention them in sections pertaining to personal data collection or handling. In order to verify this hypothesis, we manually check the annotations from the OPP-115 corpus to find occurrences of phrases containing the token search. While we find such mentions in the corpus (as shown in Table 9), we notice that search queries collection is explicitly mentioned in segments labelled as "First Party Collection/Use" by the OPP corpus annotators more frequently than segments labelled as "Third Party Sharing/Collection" (24 vs. 11). A closer inspection of these "Third Party Sharing/Collection" segments reveals that some vague or generic phrases are used to inform the reader that some of their information (without necessarily specifying which) may be shared with third parties. Examples of such generic phrasing is also shown in Table 9.

Based on this discovery, we decided to manually review all 672 segments (originating from 112 distinct privacy policies of the OPP-115 corpus) to identify those that seem to mention (either explicitly or implicitly) the collection by a third party of information (such as search queries) or the sharing of such information with a third party when the user interacts with the site. The results of this manual review are as follows: 116 segments seem to mention such a practice (labelled as "yes") while 462 do not (labelled as "no") and for 94 segments, we are unable to decide (labelled as "unsure").

We further examined these 116 "yes" segments to extract the sentences (or parts of sentences) that refer to the sharing of information with a 3rd party. Once these parts were identified, we defined patterns using morphological and syntactic information with a view to create detection rules of such mentions in newly collected privacy policies. We use rules instead of a machine learning approach due to the small amount of labelled data (116 sentence parts). These rules were implemented using Spacy's entity ruler [45] so that a sen-

tence was deemed to mention information sharing with a third party if it contained the following elements: a third party (which may also be described using words such as vendor or partner, a verb describing the act of sharing or collecting, such as disclose or gather, and a reference to information or data. In order to avoid some false positives in specific contexts (such as information sharing with law enforcement entities or information transfer in the event of a corporate merger), we penalized sentences containing specific terms. These rules showed an accuracy of 89% when run on the annotated segments from the OPP-115 corpus. We then used these rules on the 79,515 privacy policies we collected and found that 60,212 of them (about 75%) contained at least one mention of information sharing with a third party. An example of detection of a generic mention is: "These third-party ad servers or ad networks can use technology which displays the advertisements that appear on our site directly to your browser, enabling them to collect information." 11 10,914 of these policies (about 13%) did mention search terms explicitly, an example being: "We and our service providers and third-party advertising partners also use these technologies to track your activities on our Site and other websites, including the websites and web pages that you visited, the ads or content that you clicked on, any items you may have purchased and the search terms you used, in order to deliver tailored advertising to you."12 These results must be considered as upper bounds since our detection rules can still be prone to false positives.

To validate our detection results, we read 12 of the 100 validation documents manually and confirmed our automated approach did not find any mention of users' search terms in this sample. However, 2 of these 12 documents contained some mention of user interaction with some services and 8 other documents in the sample did mention users' usage of the site/services, which could encompass the use of the search functionality.

These results suggest that privacy policies tend to be worded in such a way that search query handling by a third party is very often described in a very generic manner. This means that users who are looking for an answer to the question are my search terms shared with or collected by a third party? will have to carefully read the entire privacy policy (or specific sections of it) to find such information (instead of simply searching for a keyword like search using a browser shortcut). In prac-

tice, this means that most users will not have the time to perform such a careful review and will be unaware that their search terms won't be kept private by the site they are interacting with. This is something we decided to address by designing a novel browser extension, as described in Section 5.4 on countermeasures.

5 Discussion

5.1 First Party Categorization

We used a proprietary tool to group the domains according to their content category to determine if privacy leakage varies across website categories. These results can be seen in Figure 5. By and large, the worst offenders were "Personal Sites" (92.2% Any), which is consistent with the findings of [33]. The other categories of bad actors were "Restaurants/Dining/Food" (88.1%) and "Shopping" (86.7%), which might be consistent with sites that are most likely to have the most advertising. Interestingly, some of the most well-behaved sites were sites categorized as "Piracy/Copyright Concerns" (77.6%), "Suspicious" (77.5%), and "Pornography" (77.7%). This may be because these sites exist outside the usual ad ecosystem that powers the conventional web, as documented by [71].

Finally, the "Search Engines/Portals" category ironically had the lowest value for Direct Referrer leaks (55.3%), possibly because these tend to be large sites whose central focus is search, and are no doubt aware of the sensitivity of search terms and are careful against this data leaking into the hands of competitors.

5.2 Expert Opinions on Privacy Policies

To determine the accuracy of our analysis above, we spoke to several legal professionals to determine how they might craft privacy policies and what language might be included to cover search terms. The legal professionals with whom we consulted told us that there is no standard way to write a privacy policy (though in practice one of a handful of templates are typically used). They further stated that search terms are usually not explicitly covered by privacy policies, and are often grouped under more general terms such as "site interaction". Finally, these legal professionals disagreed on whether the GDPR and CCPA regulations specifically covered search terms. One prominent privacy legal

¹¹ http://www.miicharacters.com/privacy.php

¹² https://www.audi-forums.com/help/privacy

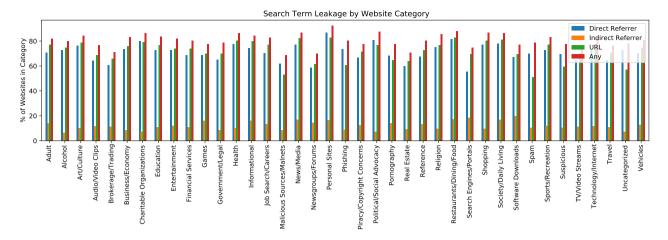


Fig. 5. Leakage breakdown by category and leakage type

expert told us he did not view search terms coupled with an IP address as personal information, showing a disconnect between the sensitivity of search terms as demonstrated by computer science research, and the interpretation of current law by the legal community.

When asked, "How might an average user determine if search terms are covered by a site's privacy policy?", our experts were pessimistic. They re-iterated that privacy policies are generally written in broad terms, and agreed that they tend to be quite long (in 2008, an average privacy policy took 10 minutes to read and this number is likely larger today as privacy policies have only grown in scope [61]). They said that users might look for specific phrases, similar to the ones extracted in Section 4.4.2.2.

5.3 Findings

As mentioned in Section 4, we successfully executed a search for our dummy query string, "JELLYBEANS", on 472,333 websites. 81.3% of these leak search terms to third parties in some form, according to our definition of leakage in Section 2.4. 75.8% of websites leaked data via the *Referer* header (72.5% directly and 10.6% indirectly); 71% via URL query parameters; 21.2% via the payload. Additionally 87.4% of websites had the potential to leak sensitive query terms to a new third-party if a link on the page was clicked.

The high incidence of leakage via the *Referer* header is likely related to the fact that over 93.6% of sites did not explicitly specify a *Referer-Policy*, and when they did it was either at the same level or worse than the default. After we completed our crawl, the

Google Chrome team announced that starting with Chrome version 85, the default *Referer-Policy* will be strict-origin-when-cross-origin [19]. This change should have a significant positive change on the privacy of users on the web, especially if adopted by the other major browsers ¹³.

Most sites leaked search terms through more than one vector. This is extremely concerning, as often times these search terms refer to sensitive topics such as medical conditions, sexual preferences, racial identity, or financial situation. Even for seemingly innocuous cases such as shopping, products can reveal sensitive information about users - consider products such as pregnancy tests, wedding bands, or anti-anxiety medication.

5.4 Countermeasures

First, some browsers such as Chrome v85+ and Safari have changed their default Referrer-Policy to the more secure "strict-origin-when-cross-origin" value. We would like other popular browsers (such as Firefox, QQ Browser and Edge) to adopt this default, and have standards bodies such as W3C change the guidance for browsers to make this the default value. We would also like to see tracking protection tools (such as those built into Firefox) flag sites that downgrade the Referrer-Policy using a server-side header or HTML, and prevent that behavior when tracking protection has been turned on by the user. We would also like website maintainers to acknowledge privacy leakage through this vector, and

¹³ Chrome is estimated to control 66% of all browser market share [8]

sites to be held accountable for the fact that this privacy leakage is not present in their privacy policies using the GDPR privacy framework or other legal standards. Finally, we would hope for a technology similar to Content Security Policies (CSPs) that deals with sharing of user data to third-parties, which would be enforced by the browser and could be inspected by automated tools.

In the meantime, we developed a browser extension which displays our findings for each crawled site, as shown in Figure 6. This extension also links to the privacy policy where possible.



Fig. 6. JellyBeans browser plugin warning a user when they attempt to use internal search.

5.5 Limitations

Crawling Limitations. During the crawl process, we failed to visit 15.7% of domains, which is roughly in line with the expected results in [32]. However, we made few efforts to circumvent IP-blocking or solve CAPTCHAs during our crawl. Some other teams use Tor to evade IP-blocking, while others employ anti-crawling-detection stealth toolkits. We do not believe such approaches are consistent with the spirit of the web, and allowed for such failures. However, it may be the case that these 15.7% of websites differ in some key characteristics from the ones that we successfully visited.

We made every effort to locate the correct internal search inputs on every website regardless of language. However, our approach was not perfect. For example, we only visited the landing page of each webpage, and did not navigate the site to search for the internal search feature on other pages. We anticipate future work will improve on our approach and be able to find more search inputs. For this reason, we will release the labeled corpus of detected search input selectors for each domain under open source for other researchers to use.

Analysis Limitations. During our analysis, we found some cases of large payloads sent to third party sites, which appeared to be encrypted or encoded, and which we suspected contained the "JELLYBEANS" search string. We utilized Ciphey to try to decode all large payloads, but sometimes were not successful [10].

Therefore, the payload leakage may in fact be higher than reported in this paper.

In addition to privacy policies, some sites also have a cookie policy document or a cookie banner. While we looked for mentions of search term handling and information sharing in privacy policies, we did not analyze these cookie-related documents. However, we do not believe that this would substantially change our results.

6 Related Work

Crawling for Privacy Leaks. The use of web crawlers to measure privacy leakages is ubiquitous in the literature: a recent survey paper by Ahmad et. al. [32] found that 350 papers published at top-tier venues from 2015-2020 relied on data gathered from web crawlers. These web crawlers were most often implemented as headless browsers on top of the PhantomJS [22]. Selenium [29], Mozilla OpenWPM [21], or Puppeteer [24] frameworks. Such crawlers were often deployed in largescale scans to discover privacy leakages around the web. often using either the Tranco [31] or the Alexa [7] top domain datasets as inputs. The resulting studies measured privacy leakages found in email tracking [41], link tracking [65], third-party cookies [66], the ad ecosystem [42], HTTP headers [51], contact forms [69], and registration forms [36]. Many other crawls characterized privacy leakages (often by looking at domains of third-party requests) more broadly on the web [33, 50, 54, 55, 68], or in specific sectors such as adult websites [60, 71], health websites [53], online collaboration services [47], and the financial sector [64]. Ahmad et al. noted that challenges exist in crawling the web in an automated way: only approximately 77-93% of the top 500 domains can be successfully reached by a crawler ¹⁴, including blocking due to regional restrictions, IP-based blocking, and encountering CAPTCHAs.

Privacy of the Referer HTTP header. The Referer HTTP header has been known to be a potential source of privacy leakage since at least 2011 [39, 50]. The Referrer-Policy HTTP header was specifically created for website owners to control the leakage of the referrer [39]. The issue of potentially sensitive queries being leaked in referrer values was showcased by Krishnamurthy et al. using the example of searching for pancreatic cancer on a health information website, and hav-

 $[{]f 14}$ depending on the crawler used

ing that search string transmitted to third parties [50]. Similar examples were shown in work done by Libert in 2014 [53]. Lavrenovs et al. demonstrated in 2018 that the Referrer-Policy HTTP header was only explicitly set in 0.05% of HTTP responses and 0.33% of HTTPS responses [51], though this work looked at HTTP headers only, while the policy may also be set inside HTML. In 2017, Dolnak characterized the Referrer-Policy HTTP header status of 7 million websites and suggested 56% of those websites might leak sensitive queries via the Referer HTTP header[39].

Search Privacy. Search queries, especially over a longer time period and associated with a specific IP address, are well-known to lead to significant loss of user privacy. Jones et al. showed that search queries can be used to efficiently de-anonymize users [46]. White et al. showed that it was possible to use search queries to diagnose neurodegenerative disorders [75]. Libert showed that health terms in medical search engines were leaked in 2014 [53]. However, Libert started his analysis by searching for likely medical terms using popular search engines, and then analyzed the resulting 80,142 pages for privacy leakage. Note that this set of pages corresponds to a much smaller set of domains, as Libert includes multiple pages per domain in his analysis. In contrast to this approach, we focus on performing an internal site search directly on each domain, which requires a more complex crawler to correctly locate the search input fields. We also scan a much larger set of domains: 1 million in all, and do not restrict ourselves to the health and medical field.

Privacy Policies. Privacy policies provide another way to measure the intended privacy practices of a company [34], but previous work has shown that their usefulness can be limited due to their complexity [62]. The adoption of privacy policies by companies has been measured in numerous studies [38, 57, 67, 77]. Some of these studies have focused on specific platforms (such as mobile apps [78] or specific categories of web sites, such as adult sites [72] [60]. For instance, Maris et al. used the webXray tool to crawl and analyze 22,484 adult websites for privacy leakages. [60] They developed a tool called policyXray to locate privacy policy links on the pages, looking for links with text such as "privacy" and "privacy policy". The privacy policy text was then extracted using Google Chrome's Readability. js library, and its reading difficulty and time to read was assessed. Maris et al. also extracted any explicitly disclosed third parties in the privacy policy. In this work, Maris et al. did not try to perform a more in-depth linguistic analysis of the extracted privacy policies. Finally, privacy policies

have also been recently analyzed automatically using deep learning so that they can be queried by users using natural language questions [44].

Users' Perception of Search Privacy. Many studies have been conducted to analyze users' perception of privacy. For instance, standardized privacy policy information formats were designed to determine whether users would find them more useful than traditional policies [49]. However, very few have specifically analyzed the users' mental models of online search activity. A recent study focusing on a tracking visualization tool [74] did find that a majority of users did not want to have their search activity tracked, while a previous study found that lay people had simpler mental models than technical models that omitted concepts such as Internet levels and entities [48] (suggesting that a very large number of users does not realize that their search queries are shared with third parties).

7 Conclusions

In this paper, we showed how terms entered into the internal site search feature of the Tranco top 1 million websites are leaked to third parties, substantially compromising user privacy. We developed a crawler called SlicedPie, built on top of Puppeteer, which can interact with modern dynamic websites, identify search inputs, and capture outgoing network requests. Our crawler was able to locate an internal site search feature on 512,701 websites, and was able to successfully execute a search for the dummy term "JELLYBEANS" 92.1% of the time.

We analyzed privacy leakage across three vectors: (i) Referer HTTP header; (ii) URL including query parameters; (iii) request payload. We detected 81.3% of these websites leaking the "JELLYBEANS" string to third-party domains through at least one of these three vectors. We released a browser extension to inform users about potential privacy leaks via internal search, based on our findings.

We also developed a novel technique to analyze both P3P and natural language privacy policies and determine whether they mention sharing search terms with third parties. We extracted privacy policies from about 50% of the websites leaking search queries, and found that only 13% of those privacy policies explicitly mentioned search terms. In most privacy policies (about 75%), the sharing of information (which may include search terms) is also described using generic wording.

References

- [1] 1996. RFC 1945 Hypertext Transfer Protocol HTTP/1.0. https://tools.ietf.org/html/rfc1945#section-10.13. (1996). Accessed: 2019-12-13.
- 1999. RFC 2616 Hypertext Transfer Protocol HTTP/1.1. https://tools.ietf.org/html/rfc2616#section-14.36. (1999). Accessed: 2019-12-13
- 2006. A Face Is Exposed for AOL Searcher No. 4417749. https://www.nytimes.com/2006/08/09/technology/09aol. html. (2006).
- [4] 2012. How Target Figured Out A Teen Girl Was Pregnant Before Her Father Did. https://www.forbes.com/sites/ kashmirhill/2012/02/16/how-target-figured-out-a-teengirl-was-pregnant-before-her-father-did/#459563216668.
- [5] 2019. Referrer Policy Editor's Draft, 4 December 2019. https://w3c.github.io/webappsec-referrer-policy/#referrerpolicies. (2019). Accessed: 2019-12-13.
- 2020. 24 best practice tips for ecommerce site search. https: //econsultancy.com/24-best-practice-tips-for-ecommercesite-search/#i.e6bb0iahterern. (2020).
- 2020. Alexa. https://www.alexa.com/topsites. (2020).
- 2020. Browser Market Share Worldwide. https://gs. statcounter.com/browser-market-share. (2020).
- [9] 2020. Chrome DevTools Protocol. https://chromedevtools. github.io/devtools-protocol/. (2020).
- [10] 2020. Ciphey. https://github.com/Ciphey/Ciphey. (2020).
- [11] 2020. The fetch API. https://developer.mozilla.org/en-US/docs/Web/API/Fetch_API. (2020).
- [12] 2020. The fetch API. https://github.com/whatwg/fetch/ issues/951. (2020).
- [13] 2020. Firefox Translation Project on Pontoon. https:// pontoon.mozilla.org/projects/firefox. (2020).
- [14] 2020. The First-Party Sets. https://github.com/privacycg/ first-party-sets. (2020).
- [15] 2020. Google's Programmable Search Engine. https: //developers.google.com/custom-search/docs/overview. (2020). Accessed: 2020-11-11.
- [16] 2020. Guidelines for Online Privacy Policies. http://www. privacyalliance.org/resources/ppguidelines. (2020).
- [17] 2020. Implementing Google Search Box. https:// developers.google.com/custom-search/docs/tutorial/ implementingsearchbox. (2020).
- [18] 2020. The language global attribute. https://developer. mozilla.org/en-US/docs/Web/HTML/Global_attributes/
- [19] 2020. A new default Referrer-Policy for Chrome: strictorigin-when-cross-origin. https://developers.google. com/web/updates/2020/07/referrer-policy-new-chromedefault#:~:text=referrer%20.,-Up%20until%20recently& text = Chrome % 20 plans % 20 to % 20 switch % 20 its, % 2D cross %2Dorigin%20by%20default.. (2020).
- [20] 2020. Node Worker Threads. https://nodejs.org/docs/ latest-v12.x/api/worker_threads.html. (2020).
- [21] 2020. OpenWPM. https://github.com/mozilla/OpenWPM. (2020).
- [22] 2020. PhantomJS. https://github.com/ariya/phantomjs. (2020).

- [23] 2020. The Public Suffix List. https://publicsuffix.org/. (2020).
- [24] 2020. Puppeteer. https://github.com/GoogleChrome/ puppeteer. (2020).
- [25] 2020. Puppeteer Chrome DevTools Protocol Session. https: //devdocs.io/puppeteer/index#class-cdpsession/. (2020).
- [26] 2020. Query selector. https://developer.mozilla.org/en-US/docs/Web/API/Document/guerySelector/. (2020).
- [27] 2020. Referer and Referrer-Policy best practices. https: //web.dev/referrer-best-practices/. (2020).
- [28] 2020. Search Action. https://schema.org/SearchAction.
- [29] 2020. Selenium. https://github.com/SeleniumHQ/selenium. (2020).
- [30] 2020. The Tracker Protection lists. https://github.com/ disconnectme/disconnect-tracking-protection. (2020). Accessed: 2021-02-28.
- [31] 2020. TRANCO list. https://tranco-list.eu/list/NYJW/ 1000000. (apr 2020).
- [32] Syed Suleman Ahmad, Muhammad Daniyal Dar, Muhammad Fareed Zaffar, Narseo Vallina-Rodriguez, and Rishab Nithyanand. 2020. Apophanies or Epiphanies? How Crawlers Impact Our Understanding of the Web. In Proceedings of The Web Conference 2020 (WWW '20). Association for Computing Machinery, New York, NY, USA, 271-280. https://doi.org/10.1145/3366423.3380113
- [33] Amirhossein Aleyasen, Oleksii Starov, Alyssa Phung Au, Allan Schiffman, and Jeff Shrager. 2015. On the Privacy Practices of Just Plain Sites. In Proceedings of the 14th ACM Workshop on Privacy in the Electronic Society (WPES '15). Association for Computing Machinery, New York, NY, USA, 1-10. https://doi.org/10.1145/2808138.2808140
- [34] Ryan Amos, Gunes Acar, Elena Lucherini, Mihir Kshirsagar, Arvind Narayanan, and Jonathan Mayer. 2020. Privacy Policies over Time: Curation and Analysis of a Million-Document Dataset. (2020). arXiv:cs.CY/2008.09159
- [35] Chih-Chung Chang and Chih-Jen Lin. 2011. LIBSVM: A library for support vector machines. ACM Transactions on Intelligent Systems and Technology 2 (2011), 27:1-27:27. Issue 3. Software available at http://www.csie.ntu.edu.tw/ ~cilin/libsvm.
- [36] Manolis Chatzimpyrros, Konstantinos Solomos, and Sotiris Ioannidis. 2020. You Shall Not Register! Detecting Privacy Leaks Across Registration Forms. In Computer Security, Apostolos P. Fournaris, Manos Athanatos, Konstantinos Lampropoulos, Sotiris Ioannidis, George Hatzivasilis, Ernesto Damiani, Habtamu Abie, Silvio Ranise, Luca Verderame, Alberto Siena, and Joaquin Garcia-Alfaro (Eds.). Springer International Publishing, Cham, 91-104.
- [37] Lorrie Faith Cranor. 2003. P3P: Making privacy policies more useful. IEEE Security & Privacy 1, 6 (2003), 50-55.
- Martin Degeling, Christine Utz, Christopher Lentzsch, Henry Hosseini, Florian Schaub, and Thorsten Holz. 2019. We Value Your Privacy ... Now Take Some Cookies - Measuring the GDPR's Impact on Web Privacy. Inform. Spektrum 42, 5 (2019), 345-346. https://doi.org/10.1007/s00287-019-
- [39] Ivan Dolnák. 2017. Implementation of referrer policy in order to control HTTP Referer header privacy. In 2017 15th International Conference on Emerging eLearning Technolo-

- gies and Applications (ICETA). IEEE, 1-4.
- [40] Steven Englehardt. 2017. No boundaries: Exfiltration of personal data by session-replay scripts. https://freedomto-tinker.com/2017/11/15/no-boundaries-exfiltration-ofpersonal-data-by-session-replay-scripts/. (2017).
- [41] Steven Englehardt, Jeffrey Han, and Arvind Narayanan. 01 Jan. 2018. I never signed up for this! Privacy implications of email tracking. Proceedings on Privacy Enhancing Technologies 2018, 1 (01 Jan. 2018), 109 – 126. https://doi.org/10.1515/popets-2018-0006
- [42] Kiran Garimella, Orestis Kostakis, and Michael Mathioudakis. 2017. Ad-blocking: A study on performance, privacy and counter-measures. In *Proceedings of the 2017 ACM on Web Science Conference*. 259–262.
- [43] Saikat Guha, Bin Cheng, and Paul Francis. 2010. Challenges in measuring online advertising systems. In Proceedings of the 10th ACM SIGCOMM conference on Internet measurement. 81–87.
- [44] Hamza Harkous, Kassem Fawaz, Rémi Lebret, Florian Schaub, Kang G. Shin, and Karl Aberer. 2018. Polisis: Automated Analysis and Presentation of Privacy Policies Using Deep Learning. (2018). arXiv:cs.CL/1802.02561
- [45] Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength Natural Language Processing in Python. (2020). https://doi.org/10. 5281/zenodo.1212303
- [46] Rosie Jones, Ravi Kumar, Bo Pang, and Andrew Tomkins. 2007. "I know what you did last summer" query logs and user privacy. In Proceedings of the sixteenth ACM conference on Conference on information and knowledge management. 909–914.
- [47] Beliz Kaleli, Manuel Egele, and Gianluca Stringhini. 2019. On the Perils of Leaking Referrers in Online Collaboration Services. In International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment. Springer, 67–85
- [48] Ruogu Kang, Laura Dabbish, Nathaniel Fruchter, and Sara Kiesler. 2015. "My Data Just Goes Everywhere": User Mental Models of the Internet and Implications for Privacy and Security. In Proceedings of the Eleventh USENIX Conference on Usable Privacy and Security (SOUPS '15). USENIX Association, USA, 39–52.
- [49] Patrick Gage Kelley, Lucian Cesca, Joanna Bresee, and Lorrie Faith Cranor. 2010. Standardizing Privacy Notices: An Online Study of the Nutrition Label Approach. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '10). Association for Computing Machinery, New York, NY, USA, 1573–1582. https://doi.org/10.1145/1753326.1753561
- [50] Balachander Krishnamurthy, Konstantin Naryshkin, and Craig Wills. 2011. Privacy leakage vs. protection measures: the growing disconnect. In *Proceedings of the Web*, Vol. 2. 1–10
- [51] Arturs Lavrenovs and F Jesús Rubio Melón. 2018. Http security headers analysis of top one million websites. In 2018 10th International Conference on Cyber Conflict (CyCon). IEEE, 345–370.
- [52] Victor Le Pochat, Tom Van Goethem, Samaneh Tajalizadehkhoob, Maciej Korczyński, and Wouter Joosen. [n. d.]. TRANCO: A Research-Oriented Top Sites Ranking

- Hardened Against Manipulation. https://doi.org/10.14722/ndss.2019.23386
- [53] Tim Libert. 2014. Privacy implications of health information seeking on the web. arXiv preprint arXiv:1404.1951 (2014).
- [54] Timothy Libert. 2015. Exposing the hidden web: An analysis of third-party HTTP requests on 1 million websites. arXiv preprint arXiv:1511.00619 (2015).
- [55] Timothy Libert. 2018. An automated approach to auditing disclosure of third-party data collection in website privacy policies. In *Proceedings of the 2018 World Wide Web Con*ference. International World Wide Web Conferences Steering Committee, 207–216.
- [56] Thomas Linden, Rishabh Khandelwal, Hamza Harkous, and Kassem Fawaz. 01 Jan. 2020. The Privacy Policy Landscape After the GDPR. Proceedings on Privacy Enhancing Technologies 2020, 1 (01 Jan. 2020), 47 – 64. https://doi.org/10.2478/popets-2020-0004
- [57] Chang Liu and Kirk P. Arnett. 2002. Raising a Red Flag on Global WWW Privacy Policies. *Journal of Computer Information Systems* 43, 1 (2002), 117–127. https://doi. org/10.1080/08874417.2002.11647076
- [58] Massimo Marchiori Martin Presler-Marshall Joseph Reagle Lorrie Cranor, Marc Langheinrich. 2002. The platform for privacy preferences 1.0 (P3P1.0) specification. http://www. w3.org/TR/P3P/. (2002).
- [59] Delfina Malandrino, Andrea Petta, Vittorio Scarano, Luigi Serra, Raffaele Spinelli, and Balachander Krishnamurthy. 2013. Privacy awareness about information leakage: Who knows what about me?. In Proceedings of the 12th ACM workshop on Workshop on privacy in the electronic society. ACM, 279–284.
- [60] Elena Maris, Timothy Libert, and Jennifer R Henrichsen. 2020. Tracking sex: The implications of widespread sexual data leakage and tracking on porn websites. *New Media & Society* 22, 11 (2020), 2018–2038.
- [61] Aleecia M McDonald and Lorrie Faith Cranor. 2008. The cost of reading privacy policies. *Isjlp* 4 (2008), 543.
- [62] Aleecia M. McDonald, Robert W. Reeder, Patrick Gage Kelley, and Lorrie Faith Cranor. 2009. A Comparative Study of Online Privacy Policies and Formats. In *Privacy Enhancing Technologies*, Ian Goldberg and Mikhail J. Atallah (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 37–55.
- [63] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research* 12 (2011), 2825–2830.
- [64] Agus Purwanto and Andi Wahju Rahardjo Emanuel. [n. d.]. The State of Website Security Response Headers in Indonesia Banking. ([n. d.]).
- [65] Iskander Sanchez-Rola, Davide Balzarotti, Christopher Kruegel, Giovanni Vigna, and Igor Santos. 2020. Dirty Clicks: A Study of the Usability and Security Implications of Click-related Behaviors on the Web. In *Proceedings of The Web Conference 2020*. 395–406.
- [66] Iskander Sanchez-Rola, Matteo Dell'Amico, Platon Kotzias, Davide Balzarotti, Leyla Bilge, Pierre-Antoine Vervier, and Igor Santos. 2019. Can I Opt Out Yet? GDPR and the Global Illusion of Cookie Control. In *Proceedings of the 2019*

- ACM Asia Conference on Computer and Communications Security. 340–351.
- [67] Kanthashree Mysore Sathyendra, Florian Schaub, Shomir Wilson, and Norman Sadeh. 2016. Automatic extraction of opt-out choices from privacy policies. In FS-16-01 (AAAI Fall Symposium - Technical Report). AI Access Foundation, United States, 270–275. 2016 AAAI Fall Symposium; Conference date: 17-11-2016 Through 19-11-2016.
- [68] Sebastian Schelter and Jérôme Kunegis. 2016. On the ubiquity of web tracking: Insights from a billion-page web crawl. arXiv preprint arXiv:1607.07403 (2016).
- [69] Oleksii Starov, Phillipa Gill, and Nick Nikiforakis. 01 Jan. 2016. Are You Sure You Want to Contact Us? Quantifying the Leakage of PII via Website Contact Forms. *Proceedings* on *Privacy Enhancing Technologies* 2016, 1 (01 Jan. 2016), 20 – 33. https://doi.org/10.1515/popets-2015-0028
- [70] Bill Tancer. 2008. Click: What millions of people are doing online and why it matters. Hachette Books.
- [71] Pelayo Vallina, Álvaro Feal, Julien Gamba, Narseo Vallina-Rodriguez, and Antonio Fernández Anta. 2019. Tales from the porn: A comprehensive privacy analysis of the web porn ecosystem. In *Proceedings of the Internet Measurement Conference*. 245–258.
- [72] Pelayo Vallina, Álvaro Feal, Julien Gamba, Narseo Vallina-Rodriguez, and Antonio Fernández Anta. 2019. Tales from the Porn: A Comprehensive Privacy Analysis of the Web Porn Ecosystem. In Proceedings of the Internet Measurement Conference (IMC '19). Association for Computing Machinery, New York, NY, USA, 245–258. https://doi.org/10.1145/3355369.3355583
- [73] Ben Weinshel, Miranda Wei, Mainack Mondal, Euirim Choi, Shawn Shan, Claire Dolin, Michelle L Mazurek, and Blase Ur. 2019. Oh, the Places You've Been! User Reactions to Longitudinal Transparency About Third-Party Web Tracking and Inferencing. In Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security. 149–166.
- [74] Ben Weinshel, Miranda Wei, Mainack Mondal, Euirim Choi, Shawn Shan, Claire Dolin, Michelle L. Mazurek, and Blase Ur. 2019. Oh, the Places You've Been! User Reactions to Longitudinal Transparency About Third-Party Web Tracking and Inferencing (CCS '19). Association for Computing Machinery, New York, NY, USA, 149–166. https://doi.org/10.1145/3319535.3363200
- [75] Ryen W White, P Murali Doraiswamy, and Eric Horvitz. 2018. Detecting neurodegenerative disorders from web search signals. NPJ digital medicine 1, 1 (2018), 1–4.
- [76] Shomir Wilson, Florian Schaub, Aswarth Abhilash Dara, Frederick Liu, Sushain Cherivirala, Pedro Giovanni Leon, Mads Schaarup Andersen, Sebastian Zimmeck, Kanthashree Mysore Sathyendra, N. Cameron Russell, Thomas B. Norton, Eduard Hovy, Joel Reidenberg, and Norman Sadeh. 2016. The Creation and Analysis of a Website Privacy Policy Corpus. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, Berlin, Germany, 1330–1340. https://doi.org/10.18653/v1/P16-1126
- [77] Razieh Nokhbeh Zaeem and K. Suzanne Barber. 2017. A study of web privacy policies across industries. Journal of

- Information Privacy and Security 13, 4 (2017), 169–185. https://doi.org/10.1080/15536548.2017.1394064
- [78] Sebastian Zimmeck, Peter Story, Daniel Smullen, Abhilasha Ravichander, Ziqi Wang, Joel Reidenberg, N. Cameron Russell, and Norman Sadeh. 01 Jul. 2019. MAPS: Scaling Privacy Compliance Analysis to a Million Apps. *Proceedings on Privacy Enhancing Technologies* 2019, 3 (01 Jul. 2019), 66 – 86. https://doi.org/10.2478/popets-2019-0037

A Crawler Challenges and Edge Cases

A.1 Interacting with the Right Inputs

As part of Stage III results, inputs not related to search are mistakenly matched because some of their attributes erroneously trigger our input detection query selectors. To anticipate and minimize the impact of triggering non-search functionalities, we interact with each input individually (one per headless browser visit) by following these steps: (i) detect the presence of a single input element via provided query selector; (ii) focus on the input element; (iii) type "JELLYBEANS"; (iv) wait for 500ms; and (v) type the return key. By doing so, we avoid interacting with non-search-related functionalities. For example, login forms require multiple fields to be filled out before triggering a useful action when the enter key is hit.

A.2 Dynamically Generated Attributes

Some websites dynamically generate random attribute names for their elements for each page visit. We detect these cases when *Stage IV* yields exceptions indicating that none of the elements collected during *Stage III* are found. To address the problem, we merge both stages in a single fallback task.

A.3 Hidden Search Inputs

Other websites do not explicitly display the search functionality to visitors, in some cases input elements are not even detectable using JavaScript. To work around this problem, we rely on the elements previously categorized as search-related (captured in *Stage III*). We can then infer from its attributes values that a search box is likely to be made available upon click. We validate the hypothesis by instrumenting the crawler to simulate a

click and waiting for navigation until any elements of type input are displayed or a timeout period is met.

A.4 Search Results Displayed Within **Another Tab**

Some websites deliberately display search results within a new browser tab. This situation requires the crawler to specifically listen and intercept all new tab creation events and capture any links associated with the ongoing search simulation. Whenever this happens a new visit to the captured URL is independently performed.