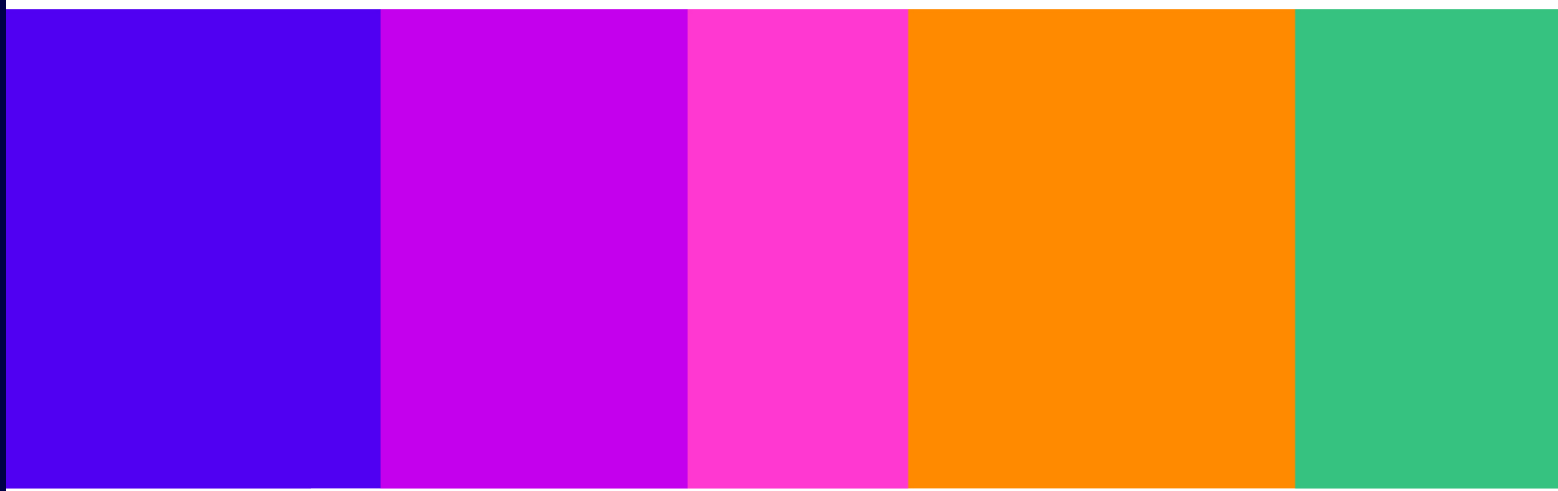


Online content for use in the commission of fraud – accessibility via search services

Ofcom research report

Report

Published 18 September 2023



Contents

Section

1. Overview	3
2. Background	6
3. Methodology.....	7
4. Results.....	12
5. Conclusion.....	23
A1.1: Content assessment result for all terms (counts)	24
A1.2: Content assessment result for all terms (total results / pages assessed).....	25
A2.1: Content assessment result for advertised results only (counts)	26
A2.2: Content assessment result for advertised results only (total results / pages assessed).....	27

1. Overview

This report summarises the findings from research conducted by Ofcom on the existence of content, accessible via general search services¹, relating to offers to supply articles (information or items) for use in the commission of fraud ('the offence').²

The report describes the results from this research.

The research sought to respond to five questions:

- Is content related to offers to supply articles for use in the commission of fraud directly accessible in search results?³
- If so, what is the prevalence of this potentially prohibited content within the first 20 search results (or two pages) delivered by search services?
- How does functionality on search services play a role in surfacing potentially prohibited content related to the offence, if at all (e.g. recommended searches, search result ordering, sponsored ads etc.)?
- How do the answers to the above questions differ between search services tested, if at all (i.e. Google Search, Bing Search)?
- What other observations can be made about the potentially prohibited content accessible from search services that can support the development of Codes of Practice, the Register of Risk and other guidance (and the caveats to these observations given the limitations of this research)?

Please note that some parts of this report containing fraud terminology and examples of content have been redacted for public use. An unredacted version of this report is available upon request.

¹ 'General search services' refers to the category of search services who provide a proprietary database of indexed webpages from which search results are chosen. For the purposes of this research, it is assumed that the findings are also relevant to 'downstream' search services which use the databases of general search services as well as supplementing with their own indexing.

² This refers to the offence of making or supplying articles for use in frauds under section 7 of the Fraud Act 2006, which is a "priority offence" under paragraph 33(c) of Schedule 7 of the Online Safety Bill. References to this offence will henceforth be referred to as 'the offence'.

³ By 'directly accessible', we mean appears in the first two pages of search results, or top 20 individual results, returned by a search service in response to the searched query.

What we found

Content offering to supply articles (information or items) for use in the commission of fraud was easy to find and prevalent on Google and Bing

Search queries used in this research returned large volumes of content within the first 20 search results which we categorised as ‘likely to be prohibited’⁴. For some search terms this was as high as 100% of search results and the vast majority (90%) of the corresponding webpages assessed.

Besides the overall prominence of this type of content, reviewing the type of content that could be accessed and the search terms used to surface it has generated several additional insights:

1. Search service functionality risked directing users towards ‘likely to be prohibited’ content through recommender systems and ranking decisions

In many cases once a keyword was inputted, search engines recommended more targeted search queries in order to provide more relevant results to the searcher. Some of these recommendations provided users with the specific phrases or language needed to surface apparent offers to supply potentially prohibited articles or items more effectively. This occurred most when the original search was based on slang terms⁵. This feature could, therefore, introduce additional risk of accessing prohibited content.

2. Content categorised as ‘likely to be prohibited’ appeared in advertised search results

Most search queries returned at least some advertised search results, many of which were categorised as ‘likely to be prohibited’. This suggested that it is possible to purchase ads promoting this kind of content on the ad exchange provided by search services. This also influences where search results appear in the results feed and may drive users towards those sites over others.

3. Search services direct users to ‘likely to be prohibited’ content on other user-to-user services which may be within scope of the Online Safety Bill

Search results analysed in this work directed users to a range of user-to-user services that hosted content related to the offence. Some content also directed users to other spaces online which enable encrypted one-to-one communication to continue with an attempted purchase of potentially prohibited articles.

4. Searches using fraud-specific terminology appeared to lead to some sites on the dark web

Several search terms returned links to the dark web within the first 20 results. Users may, therefore, be able to easily access dark web sites via search services. These sites pose risks to users and their

⁴ ‘Likely to be prohibited’ is a term developed specifically for the purposes of this research. A full explanation of this can be found in section 3.3 ‘Assessment of search results and webpages’

⁵ This may also be referred to as ‘community-specific language’ where the community refers to those who are involved in the creation and use of this kind of potentially prohibited content

devices from cyber threats, as well as raise serious concerns about potential exposure to further illegal activity due to the higher rate of illegal and malicious activity in such online spaces⁶.

5. Slang, coded language and more detailed search keywords or queries led to a higher proportion of content categorised as 'likely to be prohibited'

Slang terms were particularly effective at surfacing content that was categorised as 'likely to be prohibited' compared to when general, non-specialised words or wording was used. This indicates that while such content was generally accessible to a user, it was especially so if a user was familiar with relevant terminology.

⁶ Use of the dark web for criminal activity is relatively well documented (for example, see <https://www.openaccessgovernment.org/exposing-the-criminal-underground-of-the-dark-web/139277/>), and various collated statistics suggest the volume of content assumed to be criminal in nature on the dark web is significant (for example, see <https://blog.gitnux.com/dark-web-crime-statistics>).

2. Background

Requirement for this research

Ofcom is due to gain new responsibilities under the Online Safety Bill⁷, expected to receive Royal Assent in 2023. The Bill sets out risk assessment and safety duties with which in-scope services will need to comply, and which Ofcom will be required to enforce. Our [Roadmap to Regulation](#) explains the different elements of this new regime and the timeline for its implementation.

The Online Safety Bill also sets out a number of ‘priority’ offences. All services will need to conduct an ‘illegal content risk assessment’ which must assess, amongst other things, the risk of individuals encountering illegal content⁸ on a service, the risk of harm presented by illegal content and how the operations and functionalities of a service may reduce or increase these risks. They will also need to put in place proportionate measures to effectively mitigate and manage the risks of harm from illegal content.

This research explores one area of these priority offences, fraud. Specifically, this research was concerned with the extent to which articles (information or items) for use in the commission of fraud could be accessed via search services⁹. Other offences and areas of online harm are explored in other work.

Ofcom has a statutory duty to promote media literacy under Section 11 of the Communications Act 2003. Ofcom has a duty to ‘take such steps, and to enter into such arrangements, as appear to them calculated to bring about, or to encourage others to bring about, a better public understanding of the nature and characteristics of material published by means of the electronic media’ – (Section 11 (1)(a)). Under Section 14 (6)(a) of the Act we have a duty to make arrangements for the carrying out of research into the matters mentioned in Section 11 (1).

In an online environment where external regulatory oversight of individual pieces of content diminishes, the need for a media-literate public increases. Consumers and citizens need to be aware of the risks and opportunities offered across an array of online and mobile service activities.

⁷ Please note that all references to the Online Safety Bill in this report refer to the latest version at the time of writing, published on 19 July 2023.

⁸ “Illegal content” is a new legal concept created by the Online Safety Bill and refers to content that amounts to a “relevant offence”, as defined under the Bill (clause 59). This includes “priority offences” as set out in Schedule 7 of the Bill.

⁹ Under the Online Safety Bill, a search service is an “internet service that is, or includes a Search engine” (clause 3(4)); and a “search engine” “includes a service or functionality which enables a person to search some websites or databases (as well as a service or functionality which enables a person to search (in principle) all websites or databases)” and “does not include a service which enables a person to search just one website or database” clause (230 (1)(a)&(b)).

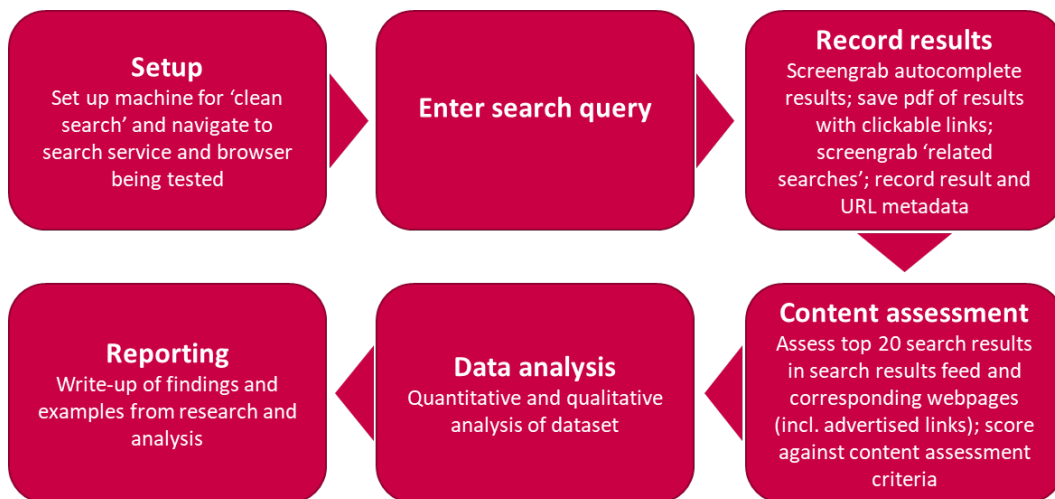
3. Methodology

3.1 Research process

The research was a manual process requiring researchers to search for pre-determined search queries and record and assess the results across Google Search and Microsoft’s Bing Search. To ensure a high quality and consistent approach, and mitigate technical issues, the method for data collection and assessment of the content was developed with input from specialists in Ofcom’s research, technology and legal teams.

Method

The diagram below outlines the main steps researchers took from setting up browsers for the search to recording data and insights from the content being assessed.



Viewing potentially malicious websites

To mitigate against the potential risk of malicious sites, researchers used a site called urlscan.io to view on-site content without opening the webpage itself and creating a cyber security risk.

This tool provided:

- a 'live' screenshot of the webpage as it appears on the web at that moment in time
- IP address and location (although these are unlikely to be the actual location of the prohibited sites)
- Document object model (DOM) providing the script for the site, allowing researchers to look at all text, features and links contained in the webpage safely

Methodological considerations & limitations

Setting up a controlled environment for searches

Researchers followed steps to limit personalisation on the browsers used to conduct the searches – Google Chrome for Google search, and Microsoft Edge for Bing search. This involved browsing in Incognito (Chrome) and InPrivate (Edge) modes; switching off any data personalisation settings; rejecting cookies; and closing the browser and clearing cookies/caches between each search query.

These steps were considered proportionate measures within the scope of the project and suitable to minimise potential search result personalisation, but could not guarantee search results were entirely unaffected by the researchers' searching activity/history within the project.

Testing was limited to Google Search and Bing Search

Research was limited to just two services because of resource limitations. Google Search and Bing Search were chosen due to the significant market share they have for search in the UK.¹⁰

Downstream search services who are accessing the Google Search or Bing Search databases are understood to have limited ability to alter the display of results, suggesting that the search results and search result ranking should be similar on these services.

Sample size and data collection timings

Data was collected over a three-week period in January and February 2023. The findings represent a snapshot of search results and webpages returned by the chosen search queries at this time. Using the same method during a different period may provide different results.

Due to the resource-intensive approach for data collection and content assessment, the dataset was limited to results from 11 search queries across Google Search and Bing Search, providing a total of 448 search results and corresponding webpages, including advertised links.¹¹ This represents only a tiny fraction of the content accessible via search services, and focused only on the main search results feed (i.e. image and video search were not included in this work). While the size of this dataset allowed for robust quantitative and qualitative analysis, it should be treated as a targeted snapshot of content surfaced by the two search services tested. It is not a representative sample of all content that could be returned by search services as a result of any search query related to the making or supply of articles for use in the commission of fraud.

Data integrity

Steps were taken to ensure consistency during the content assessment process. This included the use of a framework for assessing content and recording data, and a quality assurance process to review and challenge how content had been categorised, with all results subsequently updated based on these decisions.

A pilot phase was conducted to test and adjust the data collection process and content assessment guidance before fieldwork commenced. Data in this report is from the full fieldwork stage only.

¹⁰ Estimates vary but the combined market share of Google and Bing is assumed to be in the region of 93%. For example: <https://www.impressiondigital.com/blog/bing-differ-google/#bing-vs-google-market-share-in-2022>

¹¹ Search queries across Google and Bing did not provide equal numbers of standard and advertised results on every search results page. The total results and corresponding webpages, including advertised links, assessed per search query (per search service) varied between 17 and 24 results.

Accessing URLs and cyber security risks

Many of the websites that needed to be accessed were considered to be a cyber security risk (e.g. links to dark web sites). The use of urlscan.io allowed researchers to review content on webpages without accessing them directly. However, this did come with limitations. In certain cases, content in the webpage screenshot was harder to review, and, in a small number of instances, the pages themselves were not accessible with the tool (e.g. because the site had already been blacklisted by urlscan.io). In these instances, the webpage content was categorised as ‘unknown’.

3.2 Search queries

The research focused on a small selection of search queries due to the manual nature of the data collection and analysis process. As such, the search queries represent only a fraction of the potential terms that could be used to surface content that contains an offer to supply articles to be used in the commission of fraud. They were not intended to act as a representative sample of all relevant search queries, but a collection of terms that might surface a variety of content if it was present online, indexed by the search engines and not filtered out or downranked by the search service.

The individual queries and the reason for their inclusion can be found in the table below. They:

- Cover a range of articles and items relating to the commission of fraud (i.e. data as well as physical objects)
- Include some terms with “buy” in the search query
- Intend to represent motivated searches that someone might make – i.e. the primary focus was on whether someone looking for the likely prohibited content could find it, not on whether a user might encounter such content inadvertently

Please note that the exact search queries have been removed throughout this report, but the table below describes the kind of information they contained.

Search queries tested:

Search query	Description
Query 1	Slang terminology for contact details of people who are considered to be more vulnerable to scams
Query 2	Name for a piece of equipment used to steal payment card information
Query 3	Slang terminology referring to an online shop/store purporting to sell full credit card details (often alongside other things)
Query 4	Query combining multiple versions/abbreviations for the same thing. Terminology on sites claiming to supply stolen credit card information often combines various terms for the same thing, assumed to improve search engine optimisation
Query 5	Slang terminology referring to complete payment card and personal details
Query 6	Slang terminology for batches of stolen credit card information
Query 7	Search query referring to buying pin numbers
Query 8	Query referring to specific type of credit card information

Search query	Description
Query 9	Slang terminology for 'how to' guides on committing scams
Query 10	Slang terminology for batches of stolen credit card information (similar to query 6)
Query 11	Slang terminology referring to specific part of stolen credit card information

3.3 Assessment of search results and webpages

The scope of this research was on content that was directly accessible from the search results within 'one click', i.e. content that was displayed on the search results feed or on the landing page of any of the webpages linked to in the search results.

Assessing individual search results and the corresponding linked webpages to determine whether the content met certain criteria to be labelled as 'likely to be prohibited'¹² was central to this research.

The basis of this assessment was whether the content contained an apparent offer to sell or supply articles or items that are prohibited based on their use in the commission of fraud, as well as other indicators:

1. Whether there is an apparent offer to supply the relevant article or item (i.e. the article or item mentioned in the search query)
2. Whether there was an apparent route to purchase for a UK-based user
3. Whether there was a means of contacting the individual or group purporting to be able to supply the article(s) or item(s) in question
4. Whether the content suggested or encouraged making contact via other online services

Of these, the first two were the most important. If there was an apparent offer to supply and no obvious reason why a UK user would be prohibited from securing the article or item in question, then the content could be considered 'likely to be prohibited'. Criteria 3 and 4 were captured for additional detail.

The four categories used to label search results and corresponding webpage content were:

¹² As stated above, "illegal content" is a new legal concept created by the Online Safety Bill and refers to content that amounts to a "relevant offence", as defined under the Bill (clause 59). There are "relevant offences" relating to the offer to supply articles for use in the commission of fraud (i.e. under section 7 of the Fraud Act 2006 and section 49(3) of the Criminal Justice and Licensing (Scotland) Act 2010) - as set out in Schedule 7 of the Bill under "Priority Offences". For some of these items and articles, it can be straightforward to determine whether or not the online marketing of them is potentially illegal content. For others, it is much less clear because whether content could be considered 'illegal' will depend on offline circumstances too. For the purposes of this research, we have looked for content featuring these items where certain indicative factors (as set out at 3.3 above) are present to suggest that a person is offering to supply the items. We refer to this content as "likely to be prohibited" and the relevant items as "likely prohibited items".

Search query	Assessment criteria	Example
Unrelated	Unrelated to the supply or offers to supply articles for use in the commission of fraud <i>OR</i> Purely descriptive/explanatory	News articles about fraud Blogs or articles about tactics used by scammers Informative blogs or articles about relevant items or articles (e.g. information about how credit cards work)
Mitigating / warnings	Content focused on warning against the purchase of these kinds of tools/ information required to commit fraud	Warnings about the illegality of purchasing certain articles / committing certain acts
Below threshold	Overt offer to supply articles for use in fraud/attempted fraud <i>BUT</i> UK purchasers are excluded <i>AND/OR</i> No apparent route to purchase despite claims to be “for sale” (i.e. not clear how purchase would be achieved)	Result title claims to have articles for “for sale”, but webpage is only a collection of images of products with no payment links or means to contact the supplier Result or webpage explicitly states no UK purchasers
Likely to be prohibited	Overt offers to supply/sell articles for use in fraud/attempted fraud <i>AND</i> Apparent route to purchase (i.e. a way to progress the interaction to a point where articles would change hands)	A site with purchase functionality (e.g. shopping basket, ‘Buy now’ buttons, prices and stock availability) and/or contact information for arranging purchase (e.g. “For all sales email [...] directly”, “For [article] join our Telegram channel”)

Considerations for the assessment process

The content on the webpages and displayed in the search results was the focus of the assessment: it was not the domain/whole website itself that was under scrutiny. A label of ‘Likely to be prohibited’ reflects that the content displayed – text, images, search features – met the criteria, not that the website itself would or ought to be prohibited in its entirety (although in many cases this might be expected). See example in Annex 3.

The threshold applied was relatively high: content had to meet the strict ‘offer to supply’ and ‘route to purchase’ criteria to be labelled ‘likely to be prohibited’. Unless both of these criteria were clearly met, content would instead be labelled only as ‘below threshold’. For example, a webpage that contained reviews of other sites that overtly offer to supply credit card details would not necessarily contain its own overt offer to supply, and therefore could not be labelled ‘likely to be prohibited’.

Search results and the webpages they linked to were assessed separately: in many cases the result and the webpage received the same label (e.g. both considered ‘likely to be prohibited’), but there were instances where the webpage content was different, missing or not making such overt claims as the text in the search result.

4. Results

The main questions this research sought to respond to were:

1. Is content related to offers to supply articles for use in the commission of fraud accessible by interacting with search results?
2. If so, what is the prevalence of this potentially prohibited content within the first 20 search results (or two pages) delivered by search services?¹³

Both these questions can be better understood by looking at how content was categorised within search results and the linked webpages.

Tables showing the full breakdown for search results and corresponding webpages can be found in the annexes. The following sections focus on specific findings from quantitative and qualitative analysis of the sample.

In the following sections, the analysis focused on the following aspects of the dataset:

Term	Description
Search result(s)	The title and short description provided in the search engine results page on Google or Bing
Webpage(s), URLs	The webpage a user is taken to when clicking on a search result

Please note, other terms may be used to qualify where analysis refers to a specific subset of the data. E.g. “analysable” webpages refers to those webpages researchers were able to view and assess, and excludes the small number that were inaccessible due to being blocked by urlscan.io and not considered safe enough to click through to view in the browser.

4.1 ‘Likely to be prohibited’ content was prevalent within the first 20 search results

Of the search terms tested, 10 out of 11 returned at least some content within the search results feed and the linked webpages that was categorised as ‘likely to be prohibited’.

Among those 10 terms that returned this kind of content, as many as 100% of the search results and 90% of the linked webpages in the first two results pages contained content that was categorised as ‘likely to be prohibited’ (see Figure 1 below).

Across all search terms, the dataset contains 380 unique URLs. Of these:

- 191 search results (50%) appeared to contain (based on the criteria outlined above at 3.3) an overt offer to supply likely prohibited items

¹³ The vast majority of traffic on major search engines falls within the first 10 or 20 results, or the first two results pages. For example, research from 2020 found that click-through rates fell to 2.5% on just the tenth Google search result (<https://www.sistrix.com/blog/why-almost-everything-you-knew-about-google-ctr-is-no-longer-valid/>), while older research (2013) found that the first page of Google search results captured 92% of all search traffic (<https://www.searchenginejournal.com/the-value-of-google-result-positioning/65176/#close>).

- 111 (29%) analysable webpages appeared to contain an overt offer to supply likely prohibited items
- 104 (27%) analysable webpages appeared to contain both an offer to supply likely prohibited items *and* an apparent route to purchase them

As can be seen from the totals above, search results were more likely to meet the threshold for being considered ‘likely to be prohibited’ than the pages they linked to. This was largely due to discrepancies in the content (i.e. text) displayed in search results compared to what was analysable to researchers once webpages were accessed.

Why were search results more likely to contain content that was ‘likely to be prohibited’ than the webpages they linked to?

This discrepancy is counterintuitive in one sense: because webpages host significantly more content than the search results, they should present far more opportunity to contain content that meets the threshold for being considered ‘likely to be prohibitive’.

However, search results are not necessarily only showing a snapshot of content that exists on the source webpage. The ability to add meta descriptions and add search engine optimisation (SEO) on webpages provides opportunities for actors to present content within the search results that meets the requirement of the search query – e.g. a search for offers to supply prohibited articles – but can be different from the webpage content once accessed.

Common examples included instances where:

- Text containing offers to supply had been applied to an otherwise legitimate site and remained visible in the search result, if not on the webpage itself
- Search results led to webpages that contained no content, only a login page acting as a gateway to the page(s) beyond which potentially contain the source content driving what appears in the search result

Figure 1: Search result and webpage content categorised as ‘Likely to be prohibited’ by search term and search service

Search term	Service	Total results / pages assessed (incl. ads)	Categorised as: ‘Likely to be prohibited’			
			Search results content		Linked webpage content	
			Count	%	Count	%
Query 1	Google	19	0	0%	0	0%
	Bing	19	0	0%	0	0%
Query 2	Google	21	11	52%	7	33%
	Bing	20	10	50%	12	60%
Query 3	Google	24	21	88%	0	0%
	Bing	20	15	75%	1	5%
Query 4	Google	20	7	35%	2	10%
	Bing	17	6	35%	1	6%
Query 5	Google	23	20	87%	4	17%
	Bing	21	17	81%	11	52%
	Google	21	17	81%	10	48%

Search term	Service	Total results / pages assessed (incl. ads)	Categorised as: 'Likely to be prohibited'			
			Search results content		Linked webpage content	
			Count	%	Count	%
Query 6	Bing	22	19	86%	15	68%
Query 7	Google	20	2	10%	2	10%
	Bing	20	8	40%	8	40%
Query 8	Google	20	20	100%	18	90%
	Bing	20	16	80%	15	75%
Query 9	Google	19	4	21%	3	16%
	Bing	20	11	55%	7	35%
Query 10	Google	20	9	45%	1	5%
	Bing	20	18	90%	13	65%
Query 11	Google	21	3	14%	3	14%
	Bing	21	5	24%	4	19%

As can be seen in Figure 1 above, results were not consistent across all search terms. This was expected, as the terms were intended to cover a range of items/articles and vary in their specificity.

Because the number of search terms was small, only limited conclusions can be drawn about the composition of search queries and the affect this has on search results. However, from the data collected, the use of very specific or slang terms did appear to be an important factor. For instance, the terms which used very specific slang (e.g. Query 5, Query 6) seemed to return very high levels of 'likely to be prohibited' content. Whereas terms which contained more general language (e.g. Query 4, Query 7), returned more legitimate content which focused on the common terms within those search queries.

This could be considered the expected outcome where search engines are working as intended; delivering targeted and accurate results based on the input and the search engine optimisation efforts of the content creators.

A much larger sample of search terms and results would be needed to quantify the impact of search query wording changes, such as using "buy" as a prefix. A more detailed breakdown of common keywords found on webpages which contained content that was 'likely to be prohibited' can be found in section 4.10.

There were very few instances of false positives where a 'likely to be prohibited' search result led to an 'unrelated' webpage upon closer inspection. Across all search terms, on Google or Bing, there were only 15 instances of this. In most cases this was a 404 error, where the source content that appeared to be populating the search result had been removed from the website.

The assumption is that these slang terms are not used in other contexts, therefore there is very limited crossover with legitimate content. For example, some search queries do not have another meaning and do not appear to be used outside of the context of fraud.

4.2 ‘Likely to be prohibited’ content regularly appeared in the top three search results

Among those search terms that returned content classified as ‘likely to be prohibited’, it was common for this to appear in the first three search results, on both search services tested. The table below shows how the first three results for each search term were categorised – only two terms (Query 1 and Query 9) contained no potentially prohibited results within the first three.

Figure 2: Category of first three search results, per search term, per search service

Search term	Service	Search result #		
		1	2	3
Query 1	Google	Unrelated	Unrelated	Unrelated
	Bing	Unrelated	Unrelated	Unrelated
Query 2	Google	Prohibited	Prohibited	Prohibited
	Bing	Below threshold	Below threshold	Prohibited
Query 3	Google	Prohibited	Below threshold	Prohibited
	Bing	Below threshold	Prohibited	Prohibited
Query 4	Google	Unrelated	Grey area	Unrelated
	Bing	Unrelated	Mitigating	Mitigating
Query 5	Google	Prohibited	Unrelated	Prohibited
	Bing	Prohibited	Prohibited	Prohibited
Query 6	Google	Prohibited	Prohibited	Prohibited
	Bing	Prohibited	Prohibited	Prohibited
Query 7	Google	Prohibited	Prohibited	Unrelated
	Bing	Unrelated	Unrelated	Prohibited
Query 8	Google	Prohibited	Prohibited	Prohibited
	Bing	Prohibited	Prohibited	Prohibited
Query 9	Google	Below threshold	Unrelated	Unrelated
	Bing	Unrelated	Unrelated	Unrelated
Query 10	Google	Prohibited	Below threshold	Unrelated
	Bing	Prohibited	Prohibited	Prohibited
Query 11	Google	Prohibited	Unrelated	Unrelated
	Bing	Unrelated	Unrelated	Unrelated

4.3 Search services surfaced different volumes and types of ‘likely to be prohibited’ content

In the study, Bing appeared to link to a higher proportion of webpages that were categorised as ‘likely to be prohibited’ compared to Google, for eight out of eleven search terms. And in total, 40% of webpages in the sample (87 of 220) returned by Bing in response to the queries were categorised as ‘likely to be prohibited’, compared to 22% (50 of 228) of those returned by Google – a difference

which is statistically significant¹⁴. Given the search queries were the same this suggests there may be a material difference in which webpages the two services have indexed and/or how they deal with the search queries tested.

This difference is not necessarily surprising, given that each service has its own proprietary algorithms to assess signals from, and calculate rankings for, the billions of web pages stored in their index. However, it is not possible from this work to say which aspect of these processes caused the differences observed between Bing and Google in this sample. This analysis is also limited to the search queries tested and the periods in which the study was undertaken, and it is not possible to determine whether these differences would be consistent across all potential search queries for fraud-facilitating content at any given time.

Figure 3: Search result and webpage content categorised as ‘Likely to be prohibited’ by search term and search service

Search term	Service	Total results / pages assessed (incl. ads)	Categorised as: ‘Likely to be prohibited’				Which service had more ‘likely to be prohibited’ webpage content?
			Search results content		Linked webpage content		
			Count	%	Count	%	
Query 1	Google	19	0	0%	0	0%	=
	Bing	19	0	0%	0	0%	
Query 2	Google	21	11	52%	7	33%	Bing
	Bing	20	10	50%	12	60%	
Query 3	Google	24	21	88%	0	0%	Bing
	Bing	20	15	75%	1	5%	
Query 4	Google	20	7	35%	2	10%	Google
	Bing	17	6	35%	1	6%	
Query 5	Google	23	20	87%	4	17%	Bing
	Bing	21	17	81%	11	52%	
Query 6	Google	21	17	81%	10	48%	Bing
	Bing	22	19	86%	15	68%	
Query 7	Google	20	2	10%	2	10%	Bing
	Bing	20	8	40%	8	40%	
Query 8	Google	20	20	100%	18	90%	Google
	Bing	20	16	80%	15	75%	
Query 9	Google	19	4	21%	3	16%	Bing
	Bing	20	11	55%	7	35%	
Query 10	Google	20	9	45%	1	5%	Bing
	Bing	20	18	90%	13	65%	
Query 11	Google	21	3	14%	3	14%	Bing
	Bing	21	5	24%	4	19%	

¹⁴ At the 95% confidence level.

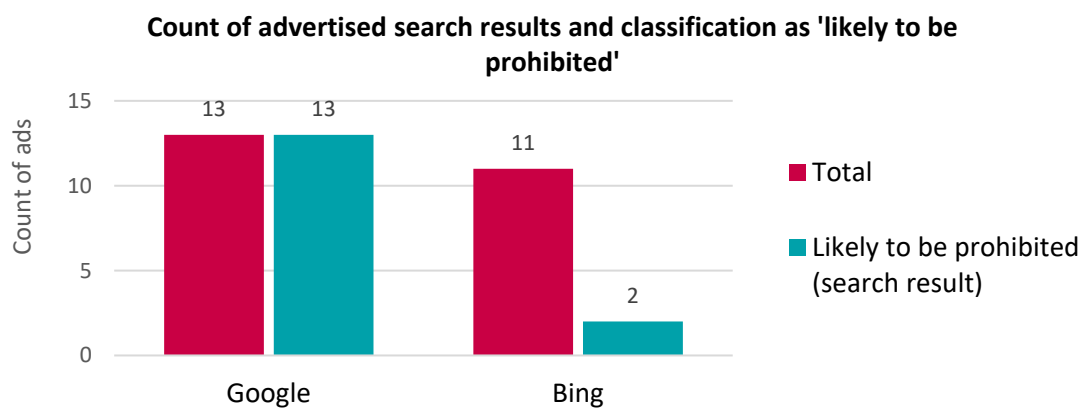
4.4 Search services were advertising a small number of ‘likely to be prohibited’ content

During content assessment of the results feed, researchers recorded a small number of advertisements. These results and their corresponding webpages were assessed in the same manner as all other results but flagged as adverts in the dataset. These advertisements are text-based and displayed among search results on a Google or Bing results page.

In total, across the 11 search terms, there were 24 advertised links (see breakdown in Annex 2.1, 2.2). 68% of these advertised links on the search results pages were categorised as ‘likely to be prohibited’, but only 18% of the corresponding webpages could be categorised this way. One of the key reasons for this was that many of the advertised links had login pages that prevented researchers from accessing the webpage. In some cases, these sites also blocked urlscan.io’s attempt to scan their website, meaning the on-site assessment was not possible.

With this sample of search queries, Google presented slightly more advertisements than Bing, and notably more advertised results categorised as ‘likely to be prohibited’ (see Figure 4 below).

Figure 4: Total number of advertised links and categorised as ‘likely to be prohibited’, by search engine



4.5 Recommender systems helped users find ‘likely to be prohibited’ content

Certain functionality on search services can help users to find prohibited content by improving the relevance or specificity of search queries.

Autocomplete suggestions and ‘related searches’ were the main functionalities that facilitated this, providing more detailed or accurate search suggestions for the kind of prohibited articles or items a user might be searching for. This was more common when the original keyword was already specific.

How do these features work?

Autocomplete and search suggestion features exist to aid individuals in completing a search they had in mind and find more relevant ways of expressing their queries. To generate these predictions, search services analyse real searches conducted by users,

and display frequently occurring and trending keywords relevant to the searcher’s initial query, location and search history. As users type new searches, these predictions are updated accordingly. In these fraud-related examples, this functionality is simply working as intended by supporting the user to find the content they are looking for.

4.6 Search services acted as a gateway to illegal content on other services which could be within scope of the Online Safety Bill

Search engines acted as gateways to ‘likely to be prohibited’ content hosted on other services that will fall in scope of the Online Safety Bill, with keywords common to illegal activities creating pathways between online spaces. Searches for fraud-related keywords led to instructive content and seemingly overt offers to supply potentially prohibited items on social media platforms.

Figure 5: breakdown of the number of results in the total sample leading to another platform that is, or is set to be, regulated under Video Service Provider or Online Safety regulation

Service	Count
YouTube	4
Google groups	4
Reddit	3
Pinterest	2
TikTok	3

Some content returned by the search services included instructions guiding users into closed online spaces. This often appeared in the form of an advert or informational post that recommended or required users to contact a user via an encrypted messaging service, phone number or email address to continue with a purchase. Encrypted services allow individuals or groups involved in illegal activities to protect themselves and their customers from scrutiny. This created a dynamic where search results were used to demonstrate an offer to supply prohibited items or articles, but the path to purchase involved routing users away from the search service.

4.7 Search services directed users to the dark web

The search queries tested also returned results that led to dark web sites.¹⁵ The presence of links to dark web sites within the first two pages of search results suggests that users could either find this content, or even be directed towards it, with minimal effort by starting their search on a mainstream search engine.

¹⁵ The “dark web” is a special part of the internet that is made up of lots of untraceable online websites. Specific software (TOR/I2P) must be used to access the websites.

30 search results, representing 27 different websites in the sample, had a direct or indirect connection to specific websites on the dark web. These were either direct dark web links; promises of a direct dark web link; or articles describing certain methods on the dark web.

Some search terms were more likely to return these dark web links and references than others. *Query 8*, for instance, delivered nearly half of all dark web links identified. Within the sample of search terms tested, Bing returned more dark web-related results than Google.

4.8 Content that had to be categorised ‘unknown’ or ‘below threshold’ was common

As highlighted in 3.3 *Assessment of search results and webpages*, there was a strict threshold for labelling content as ‘likely to be prohibited’. As such, the sample contains a relatively large number of search results and corresponding webpage content labelled as ‘below threshold’ or ‘unknown’.

‘Unknown’ refers to webpage content which could not be assessed accurately. This was often associated with webpages hidden behind login pages, or 404 errors where the content showing in the search result was no longer present on the webpage itself.

Figure 6: Search result and webpage content categorised as ‘Unknown’, by search term and search service.

Search term	Service	Total results / pages assessed (incl. ads)	Categorised as: ‘Unknown’			
			Search results content		Linked webpage content	
			Count	%	Count	%
Query 1	Google	19	-	-	0	0%
	Bing	19	-	-	0	0%
Query 2	Google	21	-	-	2	10%
	Bing	20	-	-	1	5%
Query 3	Google	24	-	-	15	63%
	Bing	20	-	-	13	65%
Query 4	Google	20	-	-	5	25%
	Bing	17	-	-	1	6%
Query 5	Google	23	-	-	12	52%
	Bing	21	-	-	5	24%
Query 6	Google	21	-	-	3	14%
	Bing	22	-	-	1	5%
Query 7	Google	20	-	-	0	0%
	Bing	20	-	-	0	0%
Query 8	Google	20	-	-	0	0%
	Bing	20	-	-	1	5%
Query 9	Google	19	-	-	5	26%
	Bing	20	-	-	9	45%
	Google	20	-	-	3	15%

Search term	Service	Total results / pages assessed (incl. ads)	Categorised as: 'Unknown'			
			Search results content		Linked webpage content	
			Count	%	Count	%
Query 10	Bing	20	-	-	5	25%
Query 11	Google	21	-	-	3	14%
	Bing	21	-	-	5	24%

As can be seen in Figure 6 above, for some search terms a notable proportion of the webpages that were assessed were classified as 'unknown', particularly Query 3 and Query 5. This was due to these queries returning several sites that required logins to access.

Whether these login pages themselves could be categorised as 'likely to be prohibited' depended on the content on the landing page.

'Below threshold' predominantly refers to content which contains an apparent offer to supply, but no apparent route to purchase, suggesting that someone would not actually be able to use the content to commit an offence (see 3.3 *Assessment of search results and webpages*). In practice, this included content such as reviews of online stores purporting to sell stolen credit card information which, despite linking to sites that themselves would be 'likely to be prohibited', did not contain a route to purchase directly from the review site/blog.

Content categorised as 'below threshold' was largely associated with the same search terms as 'unknown' content, making up a notable proportion of those search results and webpages.

Figure 7: Search result and webpage content categorised as 'below threshold', by search term and search service.

Search term	Service	Total results / pages assessed (incl. ads)	Categorised as: 'Below threshold'			
			Search results content		Linked webpage content	
			Count	%	Count	%
Query 1	Google	19	0	0%	0	0%
	Bing	19	0	0%	0	0%
Query 2	Google	21	1	5%	2	10%
	Bing	20	5	25%	1	5%
Query 3	Google	24	2	8%	9	38%
	Bing	20	3	15%	4	20%
Query 4	Google	20	7	35%	3	15%
	Bing	17	2	12%	5	29%
Query 5	Google	23	1	4%	4	17%
	Bing	21	0	0%	2	10%
Query 6	Google	21	1	5%	4	19%
	Bing	22	0	0%	3	14%
Query 7	Google	20	0	0%	0	0%
	Bing	20	1	5%	1	5%
	Google	20	0	0%	0	0%

Search term	Service	Total results / pages assessed (incl. ads)	Categorised as: 'Below threshold'			
			Search results content		Linked webpage content	
			Count	%	Count	%
Query 8	Bing	20	0	0%	0	0%
Query 9	Google	19	1	5%	1	5%
	Bing	20	4	20%	1	5%
Query 10	Google	20	6	30%	7	35%
	Bing	20	1	5%	2	10%
Query 11	Google	21	0	0%	0	0%
	Bing	21	2	10%	0	0%

4.9 Content 'unrelated' to the offence was predominant only around a small number of search terms

Of the results and webpages assessed, less than half (147 unique URLs, 39% of the total) contained content that was categorised as 'unrelated' to offers to supply likely prohibited items.

Importantly, as can be seen in Figure 8 below, although most terms returned at least some 'unrelated' content, this was predominantly associated with a small number of search terms including "Query 1" and "Query 9". In these instances, the unrelated content was largely comprised of news articles on these subjects.

Figure 8: Search result and webpage content categorised as 'Unrelated', by search term and search service.

Search term	Service	Total results / pages assessed (incl. ads)	Categorised as: 'Unrelated'			
			Search results content		Linked webpage content	
			Count	%	Count	%
Query 1	Google	19	19	100%	19	100%
	Bing	19	19	100%	19	100%
Query 2	Google	21	9	43%	10	48%
	Bing	20	5	25%	6	30%
Query 3	Google	24	1	4%	0	0%
	Bing	20	2	10%	2	10%
Query 4	Google	20	5	25%	11	55%
	Bing	17	3	18%	3	18%
Query 5	Google	23	2	9%	3	13%
	Bing	21	1	5%	1	5%
Query 6	Google	21	3	14%	4	19%
	Bing	22	3	14%	3	14%
	Google	20	18	90%	18	90%

Search term	Service	Total results / pages assessed (incl. ads)	Categorised as: 'Unrelated'			
			Search results content		Linked webpage content	
			Count	%	Count	%
Query 7	Bing	20	11	55%	11	55%
Query 8	Google	20	0	0%	2	10%
	Bing	20	4	20%	4	20%
Query 9	Google	19	14	74%	14	74%
	Bing	20	5	25%	3	15%
Query 10	Google	20	4	20%	8	40%
	Bing	20	0	0%	0	0%
Query 11	Google	21	18	86%	16	76%
	Bing	21	14	67%	12	57%

4.10 'Likely to be prohibited' content referred to common keywords and offence-specific terminology

Websites with 'likely to be prohibited' content often contained paragraphs of keywords related to the search term that surfaced them. These blocks of text were assumed to be present in order to improve the position of the webpage for related search queries – an example of keyword stuffing¹⁶. Where possible, researchers collected the paragraphs of keywords which provided a collection of terms related to each initial search query. Analysis was conducted on this collection of terms to better understand what common language and terminology was used.

¹⁶ Keyword stuffing is where the same, or similar, keywords are used over and over again on a webpage, sometimes visibly in paragraphs of text where the keywords are just separated by a comma, in other cases keywords might be hidden by putting text on a background of the same colour or adding extra words to alt text functions etc.

5. Conclusion

This research sought to provide greater clarity on the extent to which items and/or information for use in the commission of fraud could be accessed via search services. A systematic process of entering search queries and assessing the search results and corresponding webpages returned from Google Search and Bing Search demonstrated that, within the context of this research, content that is likely to be prohibited could be accessible to any user inputting relevant search queries.

In fact, we found that content categorised as 'likely to be prohibited' was extremely common – comprising up to 100% of the first 20 search results associated with some search queries. Among the general prominence of this kind of content, 'autocomplete' and 'related searches' functions recommended search queries that returned similar results, and websites claiming to sell illicit articles and items appeared in the advertised search results.

The ability of researchers to surface this content, and the patterns identified with respect to the functions of the search services tested, raise important questions about how the prominence of this kind of content may be addressed in the long term. It is hoped the insights within this report contribute to constructive conversations about this.

A1.1: Content assessment result for all terms (counts)

Search query	Service	Total results / pages assessed (incl. ads)	Results feed content categorisation				Webpage content categorisation					Webpage content features			
			Unrelated	Mitigating / warnings	Below threshold	Likely to be prohibited	Unrelated	Mitigating / warnings	Below threshold	Likely to be prohibited	Unknown	Offer to supply	Apparent route to purchase	Contact provided for purchase	Links to other online service (for sale or info)
Query 1	Google	19	19	0	0	0	19	0	0	0	0	0	0	0	0
	Bing	19	19	0	0	0	19	0	0	0	0	0	0	0	0
Query 2	Google	21	9	0	1	11	10	0	2	7	2	8	7	7	4
	Bing	20	5	0	5	10	6	0	1	12	1	13	12	10	7
Query 3	Google	24	1	0	2	21	0	0	9	0	15	7	4	7	8
	Bing	20	2	0	3	15	2	0	4	1	13	6	2	6	2
Query 4	Google	20	5	1	7	7	11	1	3	2	5	2	1	5	4
	Bing	17	3	6	2	6	3	7	5	1	1	4	3	6	3
Query 5	Google	23	2	0	1	20	3	0	4	4	12	2	2	3	1
	Bing	21	1	3	0	17	1	2	2	11	5	10	5	3	6
Query 6	Google	21	3	0	1	17	4	0	4	10	3	10	4	6	6
	Bing	22	3	0	0	19	3	0	3	15	1	15	4	6	6
Query 7	Google	20	18	0	0	2	18	0	0	2	0	2	2	2	1
	Bing	20	11	0	1	8	11	0	1	8	0	8	8	1	0
Query 8	Google	20	0	0	0	20	2	0	0	18	0	14	15	14	3
	Bing	20	4	0	0	16	4	0	0	15	1	15	15	2	3
Query 9	Google	19	14	0	1	4	14	0	1	3	5	4	3	2	1
	Bing	20	5	0	4	11	3	0	1	7	9	8	5	5	1
Query 10	Google	20	4	1	6	9	8	1	7	1	3	0	0	5	5
	Bing	20	0	1	1	18	0	0	2	13	5	8	5	4	2
Query 11	Google	21	18	0	0	3	16	0	0	3	3	2	2	1	2
	Bing	21	14	0	2	5	12	0	0	4	5	3	3	1	4

A1.2: Content assessment result for all terms (% - count / 'total results / pages assessed')

Search query	Service	Total results / pages assessed (incl. ads)	Results feed content categorisation				Webpage content categorisation					Webpage content features			
			Unrelated	Mitigating / warnings	Below threshold	Likely to be prohibited	Unrelated	Mitigating / warnings	Below threshold	Likely to be prohibited	Unknown	Offer to supply	Apparent route to purchase	Contact provided for purchase	Links to other online service (for sale or info)
Query 1	Google	19	100%	0%	0%	0%	100%	0%	0%	0%	0%	0%	0%	0%	0%
	Bing	19	100%	0%	0%	0%	100%	0%	0%	0%	0%	0%	0%	0%	0%
Query 2	Google	21	43%	0%	5%	52%	48%	0%	10%	33%	10%	38%	33%	33%	19%
	Bing	20	25%	0%	25%	50%	30%	0%	5%	60%	5%	65%	60%	50%	35%
Query 3	Google	24	4%	0%	8%	88%	0%	0%	38%	0%	63%	29%	17%	29%	33%
	Bing	20	10%	0%	15%	75%	10%	0%	20%	5%	65%	30%	10%	30%	10%
Query 4	Google	20	25%	5%	35%	35%	55%	5%	15%	10%	25%	10%	5%	25%	20%
	Bing	17	18%	35%	12%	35%	18%	41%	29%	6%	6%	24%	18%	35%	18%
Query 5	Google	23	9%	0%	4%	87%	13%	0%	17%	17%	52%	9%	9%	13%	4%
	Bing	21	5%	14%	0%	81%	5%	10%	10%	52%	24%	48%	24%	14%	29%
Query 6	Google	21	14%	0%	5%	81%	19%	0%	19%	48%	14%	48%	19%	29%	29%
	Bing	22	14%	0%	0%	86%	14%	0%	14%	68%	5%	68%	18%	27%	27%
Query 7	Google	20	90%	0%	0%	10%	90%	0%	0%	10%	0%	10%	10%	10%	5%
	Bing	20	55%	0%	5%	40%	55%	0%	5%	40%	0%	40%	40%	5%	0%
Query 8	Google	20	0%	0%	0%	100%	10%	0%	0%	90%	0%	70%	75%	70%	15%
	Bing	20	20%	0%	0%	80%	20%	0%	0%	75%	5%	75%	75%	10%	15%
Query 9	Google	19	74%	0%	5%	21%	74%	0%	5%	16%	26%	21%	16%	11%	5%
	Bing	20	25%	0%	20%	55%	15%	0%	5%	35%	45%	40%	25%	25%	5%
Query 10	Google	20	20%	5%	30%	45%	40%	5%	35%	5%	15%	0%	0%	25%	25%
	Bing	20	0%	5%	5%	90%	0%	0%	10%	65%	25%	40%	25%	20%	10%
Query 11	Google	21	86%	0%	0%	14%	76%	0%	0%	14%	14%	10%	10%	5%	10%
	Bing	21	14	0	2	5	12	0	0	4	5	3	3	1	4

A2.1: Content assessment result for advertised results only (counts)

Search query	Service	Total results / pages assessed (incl. ads)	Results feed content categorisation				Webpage content categorisation					Webpage content features			
			Unrelated	Mitigating / warnings	Below threshold	Likely to be prohibited	Unrelated	Mitigating / warnings	Below threshold	Likely to be prohibited	Unknown	Offer to supply	Apparent route to purchase	Contact provided for purchase	Links to other online service (for sale or info)
Query 1	Google	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	Bing	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Query 2	Google	1	0	0	0	1	1	0	0	0	0	0	0	0	0
	Bing	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Query 3	Google	4	0	0	0	4	0	0	0	0	4	0	0	0	0
	Bing	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Query 4	Google	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	Bing	2	2	0	0	0	2	0	0	0	0	0	0	0	0
Query 5	Google	3	0	0	0	3	0	0	0	0	3	0	0	0	0
	Bing	1	1	0	0	0	1	0	0	0	0	0	0	0	0
Query 6	Google	1	0	0	0	1	0	0	0	0	1	0	0	0	0
	Bing	3	3	0	0	0	3	0	0	0	0	0	0	0	0
Query 7	Google	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	Bing	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Query 8	Google	3	0	0	0	3	0	0	0	3	0	3	3	3	0
	Bing	3	1	0	0	2	1	0	0	1	1	1	1	0	0
Query 9	Google	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	Bing	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Query 10	Google	1	0	0	0	1	0	0	1	0	0	0	0	1	1
	Bing	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Query 11	Google	0	0	0	0	0	0	0	0	0	0	0	0	0	0
	Bing	2	2	0	0	0	2	0	0	0	0	0	0	0	0

A2.2: Content assessment result for advertised results only (% - count / 'total results / pages assessed')

Search query	Service	Total results / pages assessed (incl. ads)	Results feed content categorisation				Webpage content categorisation					Webpage content features			
			Unrelated	Mitigating / warnings	Below threshold	Likely to be prohibited	Unrelated	Mitigating / warnings	Below threshold	Likely to be prohibited	Unknown	Offer to supply	Apparent route to purchase	Contact provided for purchase	Links to other online service (for sale or info)
Query 1	Google	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
	Bing	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
Query 2	Google	1	0%	0%	0%	100%	100%	0%	0%	0%	0%	0%	0%	0%	0%
	Bing	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
Query 3	Google	4	0%	0%	0%	100%	0%	0%	0%	0%	100%	0%	0%	0%	0%
	Bing	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
Query 4	Google	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
	Bing	2	100%	0%	0%	0%	100%	0%	0%	0%	0%	0%	0%	0%	0%
Query 5	Google	3	0%	0%	0%	100%	0%	0%	0%	0%	100%	0%	0%	0%	0%
	Bing	1	100%	0%	0%	0%	100%	0%	0%	0%	0%	0%	0%	0%	0%
Query 6	Google	1	0%	0%	0%	100%	0%	0%	0%	0%	100%	0%	0%	0%	0%
	Bing	3	100%	0%	0%	0%	100%	0%	0%	0%	0%	0%	0%	0%	0%
Query 7	Google	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
	Bing	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
Query 8	Google	3	0%	0%	0%	100%	0%	0%	0%	100%	0%	100%	100%	100%	0%
	Bing	3	33%	0%	0%	67%	33%	0%	0%	33%	33%	33%	33%	0%	0%
Query 9	Google	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
	Bing	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
Query 10	Google	1	0%	0%	0%	100%	0%	0%	100%	0%	0%	0%	0%	100%	100%
	Bing	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
Query 11	Google	0	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
	Bing	2	100%	0%	0%	0%	100%	0%	0%	0%	0%	0%	0%	0%	0%

