OXFORD

# Establishing a framework for the ethical and legal use of web scrapers by cybercrime and cybersecurity researchers: learnings from a systematic review of Australian research

Katie Logos[*,†] , Russell Brewer[*,†], Colette Langos[‡], and Bryce Westlake[**]

## ABSTRACT

The Internet has become an increasingly attractive location for collecting data about cyber threats, driven by the abundance of quality data available and accessible online. As such, researchers and practitioners have turned to automated data collection technologies (ADCT), including 'web crawlers' and 'web scrapers', to study these threats. The rapid proliferation of ADCT has meant directions for their ethical and legal operation have been slow to adapt, with no clear guidelines regulating their use for research. This article identifies the relevant ethical and legal frameworks guiding the deployment of ADCT in Australia for cybersecurity research. This is accomplished through a systematic review of research within this context, coupled with ethical and jurisprudential analysis. We argue that the use of ADCT can be both ethical and legal, but only where mitigating measures are implemented. We provide a series of practical directions to guide researchers and practitioners when navigating this novel terrain.

**KEYWORDS:** Web scrape; Web crawl; Cybercrime; Cybersecurity; Ethics

[*] Lecturer, Cyber Security Cooperative Research Centre, University of Adelaide, North Terrace, Adelaide, South Australia 5005, Australia. Telephone: +61-8-8313-5633; E-mail: katie.logos@adelaide.edu.au.
[†] Associate Professor, School of Social Sciences, University of Adelaide, North Terrace, Adelaide, South Australia 5005, Australia. Telephone: +61-8-8313-5964; E-mail: russell.brewer@adelaide.edu.au.
[‡] Senior Lecturer, Adelaide Law School, University of Adelaide, North Terrace, Adelaide, South Australia 5005, Australia. Telephone: +61-8-8313-9166; E-mail: colette.langos@adelaide.edu.au.
[**] Associate Professor, Department of Justice Studies, San Jose State University, 1 Washington Sq, San Jose, CA 95192, USA. Telephone: +1-408-924-2743; E-mail: bryce.westlake@sjsu.edu.

## INTRODUCTION

Over the past several decades, the Internet has become an increasingly attractive location for collecting data that can be used by researchers and practitioners to understand and assist in combating cyber threats and enhancing security (eg the proliferation of child sexual abuse material,[1] online illicit trade[2]). Two driving forces behind this are the abundance of quality data available (eg personal information, communications, media) and the ease by which such data can be accessed. As the volume of data available online increases, those looking to collect data as a means of understanding or combating said threats (eg researchers, defence, intelligence and law enforcement agencies) have turned to automated data collection technologies ('ADCT'). These technologies are typically described as 'web crawlers' and 'web scrapers'. From a technical perspective, web crawling is the process whereby a tool is given a set of criteria and sent to 'crawl' websites looking for relevant information to capture, whereas web scraping is when a tool 'scrapes' or extracts the data from a website. While these terms are often used interchangeably, for the purpose of this paper, ADCT will refer to both of these functions.

There has been a rapid rise in use of ADCT[3] which has meant that frameworks for their ethical and legal operation have been slow to develop and adapt, with no clear guidelines regulating their use.[4] As a result, there is often confusion, ambiguity and false assumptions regarding what data collection activities are legal, when and how informed consent needs to be obtained, how to protect privacy and what potential harms exist (and should be mitigated). In particular, there are inconsistencies in the guidance provided to researchers regarding the ethical and legal frameworks for this type of work—often through institution-specific guides that are not responsive to technological advances.[5] Additionally, such guidance often does not consider the unique or novel context of the research being carried out,[6] and may be vulnerable to bias based on the expertise and knowledge of members of the institutions' ethics committee or legal department.[7] Accordingly, scholars and practitioners alike have highlighted the necessity for the development of consistent guidelines to govern online automated data collection, and ensure research in this space is adequately regulated and not unduly restricted.[8]

It is imperative that the relevant ethical and legal frameworks are well understood to inform a best-practice approach to research or data collection activities involving ADCT. Some early pioneering work has begun to flesh out legal frameworks for the use of ADCT in specific contexts, such as for 'open-source' data for Defence.[9] Elsewhere, select criminological studies[10]

---

[1] See Janis Dalins and others, 'Laying Foundations for Effective Machine Learning in Law Enforcement. Majura – A Labelling Schema for Child Exploitation Materials' (2018) 26 Digit Investig 40; Janis Dalins, Campbell Wilson and Mark Carman, 'Criminal Motivation on the Dark Web: A Categorisation Model for Law Enforcement' (2018) 24 Digit Investig 62.

[2] See Matthew Ball and Roderic Broadhurst, 'Data Capture and Analysis of Darknet Markets' (2021) SSRN <http://dx.doi.org/10.2139/ssrn.3344936> accessed 24 May 2023.

[3] Rajeev V Gundur, Mark Berry and Dean Taodang, 'Using Digital Open Source and Crowdsourced Data in Studies of Deviance and Crime' in Anita Lavorgna and Thomas J Hold (eds), Researching Cybercrimes (Palgrave Macmillan, Cham 2021).

[4] Antonella Capriello and Piercarlo Rossi, 'Spidering Scripts for Opinion Monitoring' in Hakikur Rahman and Isabel Ramos (eds), *Ethical Data Mining Applications for Socioeconomic Development* (IGI Global 2013); James Martin and Nicolas Christin, 'Ethics in Cryptomarket Research' (2016) 35 IJDP 84; Mike Thelwall and David Stuart 'Web Crawling Ethics Revisited: Cost, Privacy, and Denial of Service' (2006) 57 JASIST 1771.

[5] Zachary Gold and Mark Latonero, 'Robots welcome? Ethical and Legal Consideration for Web Crawling and Scraping' (2018) 13 WJLTA 275; Thelwall and Stuart (n 4).

[6] David Décary-Hétu and Judith Aldridge, 'Sifting Through the Net: Monitoring of Online Offenders by Researchers' (2015) 2 EROC 122.

[7] Mark McCann, 'The Smartphones Study: An Analysis of Disciplinary Differences in Research Ethics Committee Responses to Phone App-based Automated Data Collection' (2016) 26(suppl. 1) EJPH.

[8] Sophia Alim, 'Automated Data Extraction from Online Social Network Profiles: Unique Ethical Challenges for Researchers' (2013) 5 IJVCSN 24; Emil Chiauzzi and Paul Wicks, 'Digital trespass: Ethical and Terms-of-Use Violations by Researchers Accessing Data from an Online Patient Community' (2019) 21 JMIR <https://www.jmir.org/2019/2/e11985/> accessed 24 May 2023; Gold and Latonero (n 5); Martin and Christin (n 4); Thelwall and Stuart (n 4).

[9] See Lyria Bennett Moses and others, 'Using 'Open Source' Data and Information for Defence, National Security and Law Enforcement: Legal Report' (2018) Data to Decisions CRC.

[10] Décary-Hétu and Aldridge (n 6).

have acknowledged some of the ethical dilemmas facing specific online environments (eg cryptomarkets, forums), but have not considered other contexts (eg social media, peer-to-peer). While this work has been instrumental in setting the scene, we suggest that taking a holistic view of the broader relevant cyber domains utilizing ADCT can provide a fuller understanding of the suite of challenges for researchers and practitioners carrying out this work. The current paper explores this within an Australian context.

In this paper, we identify, and critically interpret both the relevant (i) ethical requirements and (ii) legal frameworks governing the deployment of ADCT for research purposes in Australia. This is accomplished using a multi-stage methodology. First, we undertake a systematic search of available evidence (peer-reviewed research, publicly available government documents and other grey literature (eg conference proceedings, theses)), as a means of providing an overview of the nature and scope of ADCT deployment regarding cybersecurity by researchers in Australia. This extracted data then serves as a foundation for a critical appraisal of the ethical dimensions of the work using the requirements set forth in the *National Statement for Ethical Conduct in Human Research*[11] (the framework governing Australian data collection practices for research, hereafter, the *National Statement*), as well as a jurisprudential analysis of relevant Australian legislation, statutes and case law. The results of this analysis demonstrate that the use of ADCT for research purposes can be both ethical and legal, but that various mitigating measures must be implemented to ensure this. In making these arguments, we also put forward a series of practical directions that will serve to guide those seeking to navigate this challenging terrain to advance research on cybersecurity, whilst also assuaging public concerns about the implications of such practices.

## REVIEW OF ADCT DEPLOYMENT FOR RESEARCH IN AUSTRALIA

This section provides a detailed review of the ADCT landscape in Australia, where such technologies have been used by researchers and practitioners for research purposes concerning cybersecurity. We begin by outlining the systematic search protocol used to identify these uses of ADCT, followed by an overview of the key characteristics of this research which will be used to inform the ethical and jurisprudential analysis.

### Systematic literature search protocol

To collect all relevant evidence involving Australian uses of ADCT for research, a systematic search was conducted which involved an exhaustive search of predetermined databases using pre-defined search terms, with each source of literature subject to strict inclusion and exclusion criteria (see Figure 1 for a flow chart of the search process). Databases included EBSCO, Google Scholar, HeinOnline, JSTOR, ProQuest and Web of Science. A series of keywords and boolean operators were used to execute the search (see Figure 1).[12] The first series of terms aimed to capture the various references to the deployment of ADCT. The second series aimed to narrow the search to the deployment of such technology within a cybersecurity research context. These terms were devised by members of the research team, and were based on known contexts in

---

[11] National Health and Medical Research Council (NHMRC), 'National Statement for Ethical Conduct in Human Research' (2018) <https://www.nhmrc.gov.au/about-us/publications/national-statement-ethical-conduct-human-research-2007-updated-2018> accessed 24 May 2023.
[12] Given the Google Scholar database does not recognize truncation, these terms were slightly adapted for this platform as follows: ('automated collection technology' OR 'data scrape' OR 'data crawl' OR 'web scrape' OR 'web crawl') AND (cyber OR crypto OR crime OR security OR terror OR 'law enforcement' OR illicit OR deviance OR regulate OR drug OR 'dark web' OR 'dark web') AND (Australia)

which these technologies have been previously deployed for cybersecurity research.[13] Finally, the third series was used to narrow the search to Australian affiliated authors. Various filters were also applied within each database search for consistency, including filtering for language (English only), year (1 January 2000 to 17 November 2021) and the source type (peer-reviewed journal articles, pre-print papers publicly available, book chapters, books, government and other official reports, conference papers and proceedings).

Searches were conducted on 17 November 2021, initially yielding 2135 sources across all databases. Three phases of inclusion criteria were then applied (see Figure 1). In Phase 1, duplicate articles were excluded and author affiliation was manually reviewed by members of the research team and included only if they were (i) Australian authors, or (ii) international authors who were affiliated with an Australian institution at the time of publication/production of the source. The latter criteria were important to include given that those authors could be subject
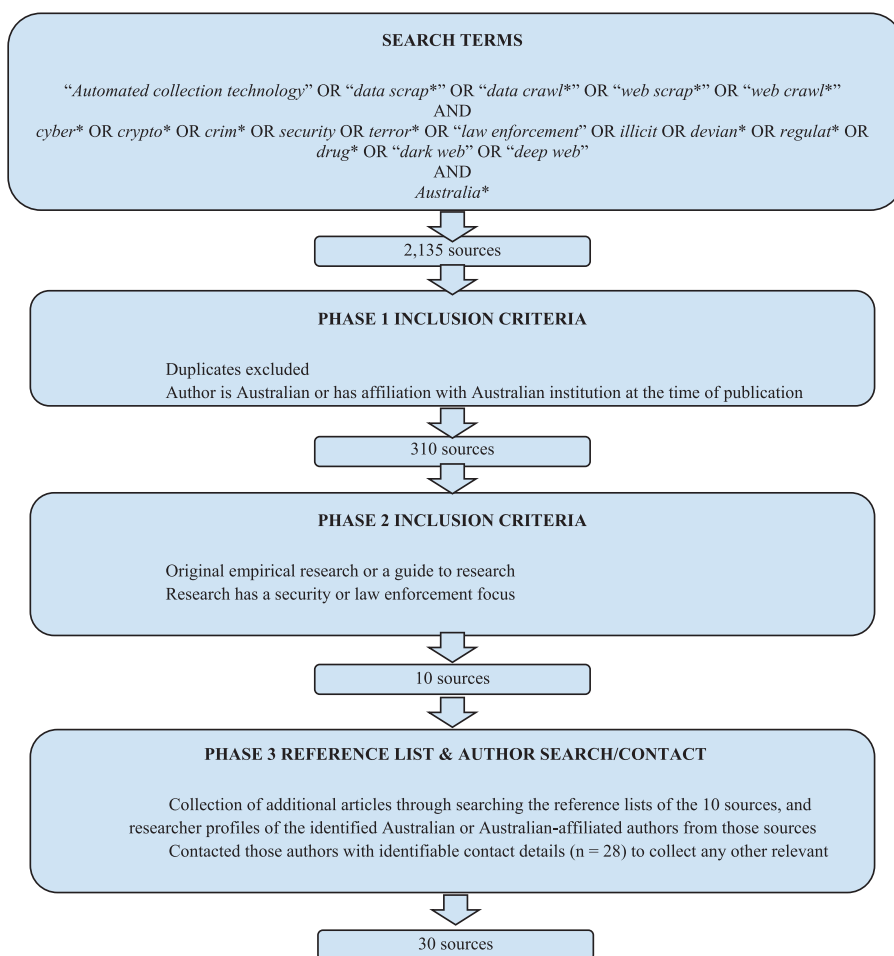


**SEARCH TERMS**

"*Automated collection technology*" OR "*data scrap*" OR "*data crawl*"" OR "*web scrap*"" OR "*web crawl*""
AND
*cyber*\* OR *crypto*\* OR *crim*\* OR *security* OR *terror*\* OR "*law enforcement*" OR *illicit* OR *devian*\* OR *regulat*\* OR *drug*\* OR "*dark web*" OR "*deep web*"
AND
*Australia*\*

2,135 sources

**PHASE 1 INCLUSION CRITERIA**

Duplicates excluded
Author is Australian or has affiliation with Australian institution at the time of publication

310 sources

**PHASE 2 INCLUSION CRITERIA**

Original empirical research or a guide to research
Research has a security or law enforcement focus

10 sources

**PHASE 3 REFERENCE LIST & AUTHOR SEARCH/CONTACT**

Collection of additional articles through searching the reference lists of the 10 sources, and researcher profiles of the identified Australian or Australian-affiliated authors from those sources
Contacted those authors with identifiable contact details (n = 28) to collect any other relevant

30 sources

**Figure 1.** Database search keywords, inclusion criteria and results

---

[13] See Russell Brewer and others, 'The Ethics of Web Crawling and Web Scraping in Cybercrime Research: Navigating Issues of Consent, Privacy, and Other Potential Harms Associated with Automated Data Collection' in Anita Lavorgna and Thomas J Hold (eds), *Researching Cybercrimes* (Palgrave Macmillan, Cham 2021).

to Australian legislation based on their institutional affiliation. Phase 1 reduced the number of sources to 310. In Phase 2, sources were excluded if (i) automated collection technologies were not used within the research, or data collected by such technologies were not analysed, (ii) it was not original empirical research (eg the source only provided a systematic review or meta-analysis of existing evidence), or (iii) the research did not have a cybersecurity focus (eg health-related). This second phase narrowed the sources to 10. Phase 3 was then carried out to identify further relevant sources through members of the research team manually examining the reference lists, and author profiles of each Australian or Australian-affiliated author, from the sources included after Phase 2. This led to the identification of a further 18 sources. Phase 3 also involved contacting all Australian authors with identifiable contact details ($n =$ 28) from those sources to enquire about any additional sources they had produced that may be relevant. This yielded responses from 16 authors, which led to the identification of two further relevant sources. Additional sources gathered during Phase 3 adhered to the same inclusion criteria from the searches and Phases 1 and 2. This three-phase process resulted in a total of 30 sources of research that utilized ADCT for security purposes by Australian-affiliated authors (see Supplementary Material 1 for information pertaining to each source). This included 28 pieces of empirical research. Additionally, our search uncovered two sources providing descriptive guides for researchers on the use of ADCT—one focusing on the dark web and the other on wildlife trade. While not included in our ethical and jurisprudential analysis, these two sources did provide important considerations for ethical and legal use and are thus referenced throughout the recommendations.

### Key characteristics of the use of ADCT in Australian research

Within the 28 empirical sources, researchers and practitioners used a myriad of web crawlers and scrapers—including custom-built tools specifically designed for the website from which they were examining,[14] or automated technologies deployed for data collection by researchers outside of Australia.[15] Not all sources involved researchers generating unique datasets through the use of ADCT, with five sources accessing existing datasets collected as part of a previous project by the authors, or other researchers.[16] However, regardless of how the data has come to be in their possession, researchers must always be mindful of how to use such data. Therefore, the discussion and guidance regarding the ethical and legal considerations presented in this paper are pertinent for all researchers who *collect*, or *use*, data collected via ADCT.

With regards to Internet locations where ADCT were deployed, five broad categories were identified (Table 1). The majority (18) deployed ADCT on the dark web, with 13 targeting marketplaces specifically (eg sale of illicit goods and services) and five targeting the dark web more generally. Social media websites, which included blogs and forums, was the next most common (6), followed by E-commerce (4) and pornography websites (1).

The type of data collected across these locations also varied considerably (Table 2), and could be categorized in three ways: sale listings (21), digital communications (8) and graphic media files (2). Within 'sale listings', research focused largely on cryptomarkets and classifieds,

---

[14] See, for example, Ball and Broadhurst (n 2); Jack Foye and others, 'Illicit Firearms and Other Weapons on Darknet Markets' (2021) 622 Trends and Issues in Crime and Criminal Justice <https://doi.org/10.52922/ti78009> accessed 24 May 2023; Tong Chen and others, 'A Hidden Astroturfing Detection Approach Base on Emotion Analysis' in Gang Li and others (eds) *Lecture Notes in Computer Science Volume 10412* (Springer International Publishing 2017).

[15] See, for example, Julian Broséus and others, 'Forensic Drug Intelligence and the Rise of Cryptomarkets. Part I: Studying the Australian Virtual Market' (2017) 279 FSI 288; James Martin and others, 'Effect of Restricting the Legal Supply of Prescription Opioids on Buying through Online Illicit Marketplaces: Interrupted Time Series Analysis' (2018) 361 BMJ <https://www.bmj.com/content/361/bmj.k2270> accessed 24 May 2023.

[16] See, for example, Broséus and others (n 15)

**Table 1:** Number of Sources by Location of Data Collection

| Location of data collection | No. of sources |
|---|---|
| Dark web marketplaces[a] | 13 |
| Dark web (general)[a] | 5 |
| Social media websites (incl. forums, blogs) | 6 |
| E-commerce websites | 4 |
| Pornography websites | 1 |

[a]  One source deployed technologies across multiple locations

**Table 2:** Number of Sources by Data Type

| Data type | Content | No. of sources |
|---|---|---|
| Sale listings | Illicit drugs[a] | 11 |
| | Illicit wildlife | 3 |
| | Illicit firearms | 1 |
| | Other illicit goods/services | 4 |
| | COVID-19 illicit products | 1 |
| | Clothing | 1 |
| Digital communications | Social media posts and profiles | 3 |
| | Forum posts[a] | 4 |
| | Blog posts | 1 |
| Graphic media files | Child sexual abuse material | 1 |
| | Adult pornography | 1 |

[a]  One source deployed technologies to collect multiple data types

with the objective of that work being to better understand how these markets (and their users) operate. The nature of the data extracted from this work was diverse and included information about the goods/services being sold, transactions (eg sale price, date, location) and the users of the platform (eg vendor and buyer attributes).

Elsewhere, eight sources used ADCT to extract digital communications (ie text-based data) between users across social media (3), web forums (4) and blogs (1). These data have been used to produce insights into the discourse surrounding illicit drug use, malware attacks, the spread of misinformation and astroturfing, and public perceptions of image-based sexual abuse. Text-based data have also been used to extract features to assist in the identification of anomalous communications on social media that may indicate the potential for illegal behaviour. Finally, two sources collected graphic media files depicting pornography (1) and child sexual abuse material (CSAM) (1). This research sought to generate knowledge about the spread of CSAM online, and develop new law enforcement capabilities (ie training machine learning models to classify sexual imagery).

## ANALYTICAL FRAMEWORK

The sources (and contexts) identified through the systematic literature search were subject to an ethical and jurisprudential review. This was carried out to determine the relevant themes and

practical measures implemented by those deploying ADCT, and to identify the application of relevant ethical dimensions and laws that may guide appropriate deployment. The frameworks guiding this review are described in turn.

### Identifying and examining relevant ethical dimensions

There are myriad ethical challenges individuals potentially face when using ADCT for research. In Australia, the requirements for appropriate ethical conduct in the collection, use or dissemination of data pertaining to human subjects for research purposes is governed by institutional review boards, taking guidance from the *National Statement*. The *National Statement* consists of a series of guidelines in accordance with the *National Health and Medical Research Council Act 1992* (Cth), which are guided by several core themes that must be attended to where research involves humans (or secondary data that has been generated by humans). More specifically, these themes call attention to, and highlight researcher consideration of, (i) the risks and benefits of the research, in conjunction with (ii) participants' consent. Accordingly, the forthcoming analysis will review the research identified through the systematic search, and adopt an analytical framework derived from these themes and resultant guidelines. In particular, using the procedures outlined in Chapter 2.1, risks (ie the potential for harm, discomfort or inconvenience to participants or others)[17] present in the literature will be identified, their probability and severity assessed against the potential benefits of the research and risk management/mitigation strategies examined. In addition, we also review this research to identify and interpret conditions that might necessitate participant consent (as per the Guidelines set out in Chapter 2.2), as well as those conditions where a waiver of consent (as per the Guidelines set out in Chapter 2.3) has been, or could appropriately be applied. Each of these will be discussed in turn, and in conjunction with relevant legal frameworks, which are briefly introduced immediately below.

### Identifying and examining relevant legal frameworks

In light of the nature and types of data collected within the research identified by the systematic search, relevant legislation and pertinent case law was identified for jurisprudential analysis. Beyond international law, namely Article 32 the *Budapest Convention on Cybercrime*[18] (for which Australia is a signatory), there is no specific law which expressly governs the use of ADCT in Australia. Further, there is very scant judicial guidance.

A comprehensive review of federal and state legislation registers revealed four primary Australian legal frameworks that have direct relevance to ADCT. These included (i) the *Data and Availability and Transparency Act 2022* ('the DATA Scheme'), (ii) the *Privacy Act 1988* (Cth) ('Privacy Act') and associated State- and Territory-based privacy frameworks, (iii) the *Copyright Act 1968* (Cth) ('the Copyright Act'), and (iv) Commonwealth, State, and Territory-based criminal codes. Case law pertaining to each of these legal frameworks was identified through case law search engines, including Westlaw Au, HeinOnline and Lexis Advance. Below we provide an overview of these legal frameworks, and to whom (and what) they apply.

### *DATA Scheme*

The *Data and Availability and Transparency Act 2022* ('the DATA Scheme') establishes a regulatory framework under which federal bodies are authorized to collect and share public sector

---

[17]  NHMRC (n 11), [12]

[18]  This describes the rights for any party to, without the authorization of another party, 'access publicly available (open-source) stored computer data, regardless of where the data is geographically located' *Budapest Convention on Cybercrime,* opened for signature 23 November 2001 ETS No. 185 (entered into force 1 July 2004) art 32.

data[19] with particular users. *The DATA Scheme* applies to three parties, namely (i) 'data custodians' who are Commonwealth government bodies controlling public sector data, (ii) 'accredited users' who are state and territory bodies and Australian universities[20] and (iii) 'accredited data service providers' who are state and territory bodies and Australian universities granted accreditation to provide complex data integration, de-identification and secure data access services to support data sharing.[21] The *DATA Scheme* stipulates general privacy principles which prohibit the sharing of biometric data unless by express consent of the individual to whom the data relates; prohibit an accredited entity with whom personal data is shared from storing, accessing, providing access to the data or the output of the project outside Australia; and prohibit reidentification of de-identified human subject data.[22] It also expressly bars the sharing of data in certain situations including, but not limited to,[23] when the data is operational data originally held or received by the AFP;[24] where the sharing of data contravenes copyright or intellectual property rights to which the data is subject;[25] if sharing of data is inconsistent with any of Australia's binding international agreements;[26] or a copy of the data being shared is being held as evidence before a court.[27] The enforcement scheme is set out in Chapter 5 of the legislation.

The recently enacted 2022 legislation broadens the legal basis for lawful data sharing, which is likely to foster increased levels of engagement between entities, such as government agencies and Australian universities, particularly in the pursuit of research and development. Each project will need to be assessed against the requirements of the legislation to ensure compliance. Given the infancy of *the DATA Scheme*, there is no case law which provides judicial guidance on the interpretation of the legislation.

### *Privacy Act*

At the federal level, the *Privacy Act 1988* (Cth) ('Privacy Act') regulates the use of personal information,[28] including sensitive information,[29] with the aim of protecting the privacy of individuals. Embedded within the *Privacy Act* are 13 Australian Privacy Principles (APPs) which apply to Australian federal government agencies and private sector organizations (APP entities).[30] The Australian Federal Police (AFP) are specifically named as an APP agency. Websites based outside of Australia, such as many of the large social media sites (eg Facebook), may be APP entities by virtue of an 'Australian link' which includes among other factors incorporation in Australia, carrying on business in Australia, and collecting or holding data in Australia.[31] The

---

[19]  Data Availability and Transparency Act 2022 (Cth), s 9—'data lawfully collected, created or held by or on behalf of a Commonwealth body, and includes ADSP-enhanced data'.

[20]  ibid, s 11(4)—'Accredited users' are defined, noting that 'excluded entities' are listed in 11(3) and include the AFP.

[21]  ibid, s 11; See also <https://www.datacommissioner.gov.au/the-data-scheme> accessed 24 May 2023.

[22]  Data Availability and Transparency Act 2022 (Cth), s 16A(1)-(3).

[23]  ibid, s 17.

[24]  ibid, s 17(2)(ii).

[25]  ibid, s 17(3)(a)(i).

[26]  ibid, s 17(5)(a)(i).

[27]  ibid, s 17(6)(a).

[28]  Privacy Act 1988 (Cth), s 6(1), states that 'personal information' is information about an identified individual, or an individual who is reasonably identifiable: (i) whether the information or opinion is true or not; and (ii) whether the information or opinion is recorded in a material form or not. Note, personal information that has been de-identified will no longer be personal information.

[29]  ibid, s 6(1), states that 'sensitive information' is information or an opinion about an individual's: racial or ethnic origin; racial or ethnic origin; political opinions; membership of a political association; religious beliefs or affiliations; philosophical beliefs; membership of a professional or trade association; membership of a trade union; sexual orientation or practices; criminal record; that is also personal information or health information about an individual or genetic information about an individual that is not otherwise health information or biometric information that is to be used for the purpose of automated biometric verification or biometric identification or biometric templates.

[30]  ibid, s 6C, states that 'organisations' include an individual, a partnership, a body corporate, a trust and any other incorporated association.

[31]  ibid, s 5B(1)(a).

**Table 3:** Applicable Legal Framework Governing Privacy

| Level | Framework | Applicable to |
|---|---|---|
| Federal | *Privacy Act 1988* (Cth) | Australian federal government agencies<br>Private sector agencies (eg Australian Federal Police, websites based outside Australia with an 'Australian link', private Australian universities, Australian National University), any researcher/practitioner conducting joint research with one of the above |
| State | | |
| SA | South Australian Cabinet Administrative Instruction | Public sector agencies<br>Not public universities |
| WA | *State Records Act 2000* (WA) | Public sector agencies<br>Not public universities |
| VIC | *Privacy and Data Protection Act 2014* (Vic) | Public sector agencies<br>Not public universities |
| QLD | *Information Privacy Act 2009* (Qld) | Public sector agencies<br>Unclear re public universities |
| NSW | *Privacy and Personal Information Protection Act 1998* (NSW) and *Health Records and Information Privacy Act 2002* (NSW) | Public sector agencies<br>Unclear re public universities |
| NT | *Information Act 2002* (NT) | Public sector agencies<br>Unclear re public universities |
| ACT | *Information Privacy Act 2014* (ACT) | Public sector agencies<br>Unclear re public universities |
| TAS | *Personal Information and Protection Act 2004* (Tas) | Public sector agencies<br>Public universities |

2022 Federal Court decision of *Facebook Inc v Australian Information Commissioner* ('*Facebook*' case)[32] provided specific guidance on this point:

> [A]n Australian link will only be present where an organisation has collected or held personal information in Australia, and it is that information which is alleged to have been misused or mishandled in contravention of the Act. If, for example, Facebook Inc collected personal information from users in Australia and then separately collected and misused personal information from users in the United States, such conduct would be beyond the scope of the Privacy Act.[33]

Notably, the APPs apply to private Australian universities and the Australian National University (ANU)[34] and, thus, must be complied with by researchers of these institutions. The legal requirements imposed under the federal privacy framework *do not* apply to public

---

[32] [2022] FCAFC 9.

[33] ibid 22.

[34] Office of the Australian Information Commissioner (OAIC), 'Rights and Responsibilities' <https://www.oaic.gov.au/privacy/the-privacy-act/rights-and-responsibilities> accessed 24 May 2023.

Australian universities. The majority of the Australian universities involved in ADCT across the 28 sources identified in our search were public, except for six sources where researchers were based at ANU. Therefore, the federal framework did not necessarily apply for researchers at those institutions. Instead, state-based laws and university-specific privacy policies regulate privacy requirements for public universities (see Table 3 for relevant state- and territory-based frameworks, and Supplementary Material 2 for detail on who each framework applies to). Given the research linkages between institutions and government agencies, it is important to note that a public university undertaking joint research with the private sector or federal public sector, will be bound by the APPs. This was the case in four of the sources, where researchers partnered with the AFP. Furthermore, while relevant state and territory privacy laws must be considered where applicable to the entity, the federal *Privacy Act* takes precedence over state law to the extent of any inconsistency. Additionally, researchers of any Australian university must also have regard to institution-specific privacy policies.

## Copyright Act

The *Copyright Act 1968* (Cth) ('the Copyright Act') provides guidance on copyright law and is enacted at the federal level, with no specific copyright laws at the state and territory level. In accordance with the legislation, copyright subsists in an original published or unpublished literary, dramatic, musical or artistic work where the author is an Australian citizen, and where the work is published, the first publication took place in Australia.[35] Online data which fall within the categories afforded copyright protection include social media content: social media posts being construed as 'literary work' and photos as 'artistic work'. The caveat here being that content must have sufficient originality: the creation of the work must require some independent intellectual effort.[36]

Compilations of information may attract copyright protection where they can be construed as having sufficient originality, noting a degree of creativity is required in regard to the selection, form, arrangement of a database.[37] Of relevance also is the issue of authorship. Judicial guidance on this point suggests that human involvement (expressed as intellectual effort) is required for the material to be protected under copyright law.[38] It is, therefore, unlikely that unorganized data collected via automated technology and collated in a mechanical manner would infringe copyright law.

## Criminal Codes

Federal, State and Territory-based criminal codes also provide guidance on pertinent criminal laws. In particular, there are two areas of direct relevance to web scraping activities. First, the unauthorized access to computer systems is criminalized in both Federal and State legislation. At a Federal level, an offence is committed where a person accesses or modifies 'restricted data' (data that is held in a computer, and to which access is restricted via an access control system), doing so intentionally, and knowing that such access or modification is unauthorized.[39] All states and territories have comparable 'computer offences' provisions embedded within their criminal codes, which will apply to persons located within their respective borders (see Table 4). While none of the research studies considered here reported breaching a computer system in such a way to collect data, this is nevertheless a serious risk that necessitates flagging as a consideration by future researchers.

---

[35]  Copyright Act 1968 (Cth), s 32.
[36]  *IceTV Pty Ltd v Nine Network Australia Pty Limited* [2009] HCA 14 33.
[37]  ibid.
[38]  *Telstra Corporation Ltd v Phone Directories Co Pty Ltd* [2010] 194 FCR 142.
[39]  Criminal Code Act 1995 (Cth), s478.1

**Table 4:** Applicable Legal Framework Governing Criminal Law

| Level | Framework |
| --- | --- |
| Federal | *Criminal Code Act 1995 (Cth)* |
| State | |
| SA | *Criminal Law Consolidation Act 1935* |
| WA | *Criminal Code Act Compilation Act 1913* |
| VIC | *Crimes Act 1958* |
| QLD | *Criminal Code Act 1899* |
| NSW | *Crimes Act 1900* |
| NT | *Criminal Code Act 1983* |
| ACT | *Criminal Code 2002* |
| TAS | *Criminal Code Act 1924* |

Beyond unauthorized access, it can also be a criminal offence to possess certain kinds of digital content that could be collected by ADCT. The Federal, State and Territory frameworks listed in Table 4 also define the categories of materials that may be prohibited to collect and/ or possess, including, for example, content depicting child sexual abuse (ie videos, images or text describing such activities).[40] These frameworks also specify applicable exemptions to such breaches,[41] which vary from one jurisdiction to the next, but typically exclude activities that are in good faith and (i) are undertaken for the purpose of advancing educational, medical or scientific knowledge;[42] (ii) are undertaken by a law enforcement, intelligence or security officer in the course of their duties; or (iii) in the administration of justice or another formal classification or child protection function.

## INTERPRETING ETHICAL AND LEGAL CHALLENGES

The use of ADCT across all sources identified within the systematic search raised a number of overlapping ethical and legal challenges for those deploying the technology to navigate. The following discussion will analyse these challenges through the framework of the *National Statement,* alongside those relevant legal frameworks, to elucidate the legitimacy of deploying ADCT within a cybersecurity context and to provide guidance on how such challenges may be overcome or mitigated. The discussion will be structured in line with the primary provisions outlined by the *National Statement,* namely the consideration of (i) risks of research and (ii) participant consent.

### Identifying the ethical and legal risks

From reviewing the nature of the research across the 28 sources identified by the systematic search, a number of potential risks of harm were raised for both the subjects or platforms at the centre of the research, as well as the researchers themselves (see Table 5). In 20 sources,

---

[40]  See, for example, Criminal Code Act 1995 (Cth), s 474.22A—A person commits the offence of 'possessing or controlling child abuse material obtained or accessed using a carriage service' if: (i) the person has possession or control of material; (ii) the material is in the form of data held in a computer or contained in a data storage device; (iii) the person used a carriage service to obtain or access the material and (iv) the material is child abuse material. This offence carries a maximum penalty of 15 years imprisonment.

[41]  See, for example, ibid, s 474.24.

[42]  There are a series of caveats associated with this provision. For example, the Criminal Code Act 1995 (Cth), s 474.22A, stipulates a requirement for those seeking an exemption to obtain written approval from the AFP Minister. Other requirements are also stipulated in State and Territory legislation.

**Table 5:** Number of Sources by Reported or Expected Risks to Subjects

| Challenge | No. of sources |
| --- | --- |
| Potential subject privacy violations (protected data) | 20 |
| Disrupt legitimate website traffic | 22 |
| Researcher breach of policies (Terms of Use, robots.txt) | 25 |
| Researcher breach of copyright | 28 |
| Researchers possessing illegal content | 14 |
| Researchers distributing illicit or high value content | 20 |
| Researcher psychological harms | 1 |
| Researcher identification and associated harms | 5 |

potential privacy violations were reported or expected to occur for the subjects, and in 22 sources, ADCT had the potential to block or disrupt legitimate traffic on the websites being examined. For researchers, the vast majority of sources indicated they may be subject to litigation through the breach of the Terms of Use or provisions of the robots.txt file of a website, or subject to copyright infringement. Across 14 sources, they risked the possession of illegal data (eg CSAM) and in 20 they possessed potentially high value content (eg closed-source personal data) which could be at risk of being intercepted and distributed. Elsewhere, one source flagged the risk of psychological harms for researchers associated with the data collection activity, and in five sources, researchers potentially put themselves at risk of being identified by the platforms or users being scraped. These potential risks of harm to both research subjects and the researchers themselves are elaborated upon below.

*Privacy*

The distinction between public (open-source) and private (closed-source) information can become blurred when posting online, and expectations of privacy may vary between users and across contexts. For example, in eight of the sources, data were publicly accessible for scraping (eg information posted on a public social media account or open forum where clear rules state that the posts are made public).[43] In 20 of the sources, investigators were required to create accounts to access and scrape information that was hidden from public view (eg closed forums, blogs or marketplaces).[44] While there may be no restrictions on the creation of an account to access and collect that closed information (eg only requiring a valid email address to sign-up to a website), this still has potential privacy implications and poses an ethical dilemma.

Although definitions around what constitutes public, or open-source, information online are, at best, unclear, scholars suggest that the social context and norms of the community being studied should inform the ethical appropriateness of the data collection—regardless of whether the data are immediately publicly available online or hidden until an account with the website is created.[45] For example, when scraping posts from a community of users with 'public' accounts who are sharing their opinions and using hashtags to contribute to a topic, these users may likely consider their content as public. This may also apply to communities of dark web cryptomarket vendors who wish to advertize their illicit products to a large audience in order to make a sale in

---

[43] See, for example, Frana-Katica Batistic and others, 'Analysis of Google Trends to Monitor New Psychoactive Substance. Is there an Added Value?' (2021) 326 FSI <https://doi.org/10.1016/j.forsciint.2021.110918> accessed 24 May 2023.

[44] See Neda Afzaliseresht and others, 'From Logs to Stories: Human-centred Data Mining for Cyber Threat Intelligence' (2020) 8 IEEE Access 19089.

[45] Décary-Hétu and Aldridge (n 6); Martin and Christin (n 4); Sergio Pastrana and others, 'Crimebb: Enabling Cybercrime Research on Underground Forums at Scale' (2018) Proceedings of the 2018 World Wide Web Conference 1845.

an act of 'crypto-anarchy'.[46] However, some communities may be less willing to consider their data 'public'. For example, while account creation may provide researchers with access to data in closed web forums, some users within said forums may view this as a breach of their privacy, particularly where there is discussion of delicate content or illegal activity (eg in a victim support forum;[47] or a hacker forum[48]).

Guidance on this ambiguous ethical risk may be best informed by the legislation. The *Privacy Act* and various state and territory privacy frameworks focus on data that is considered *personal* or *sensitive*. As evidenced through the sources identified within this paper, ADCT for cybersecurity research may deliberately or incidentally collect both personal information (ie 10 sources collected digital communications and media files that may identify an individual through usernames, communication content, or media depiction), and sensitive information (ie those same 10 sources may include sensitive information pertaining to those identifiable individuals (eg race, sexual orientation, political beliefs, biometrics).

These privacy frameworks do not place a blanket prohibition on the collection of personal (or sensitive) information. No breach occurs where the collection of personal information is reasonably necessary for, or directly related to, the organization's functions or activities and collected via lawful means. Personal information must be collected from the individual, unless it is unreasonable or impractical (or another exception applies), noting that sensitive information can only be collected with the individual's consent (unless an exception applies).[49] The topic of consent, and difficulties associated with gaining consent through the use of ADCT, is further outlined in the *Consent* discussion section below. However, regardless of consent or the waiver of consent, a central legal issue underpinning the automated collection of data relates to the *treatment* of personal or sensitive information collected online, whether it be from closed- or open-source locations. The below discussion focuses on the legal requirements researchers and practitioners must have regard for after collecting such information, in line with the Australian Privacy Principles (APPs) outlined in the *Privacy Act*.

Where APPs apply, it is imperative to comply with the legal requirements regarding the collection of personal information. The Office of the Australian Information Commissioner (OAIC) guidelines clarify that 'collecting information' includes gathering, acquiring or obtaining personal information by any means and from any source, including individuals, other entities, surveillance cameras where an individual is identifiable or reasonably identifiable, information associated with web browsing (ie that collected by cookies) and by biometric technology, such as voice or facial recognition.[50] Judicial guidance on the meaning of 'personal information' was provided in *Privacy Commissioner v Telstra Corporation Limited*.[51] The court opined that:

> [I]n every case it is necessary to consider whether each item of personal information requested, individually or in combination with other items, is about an individual. This will require an evaluative conclusion, depending upon the facts of any individual case, just as a determination of whether the identity can *reasonably* be ascertained will require an evaluative conclusion.[52]

However, only personal information collected for inclusion in 'a record' (a document or electronic device)[53] or 'generally available publication' (such as in a magazine, book, article,

---

[46] See, for example, Décary-Hétu and Aldridge (n 6).
[47] Tully O'Neill, ''Today I Speak': Exploring how Victim-survivors Use Reddit' (2018) 7 Int. J. Crime Justice Soc. Democr. 44.
[48] Pastrana and others. (n 45).
[49] Privacy Act 1988 (Cth), APP 3; Office of the Australian Information Commissioner (OAIC), 'Guide to Data Analytics and the Australian Privacy Principles' <https://www.oaic.gov.au/privacy/privacy-guidance-for-organisations-and-government-agencies/more-guidance/guide-to-data-analytics-and-the-australian-privacy-principles> accessed 24 May 2023.
[50] OAIC, *Australian Privacy Principles Guidelines*, (2019) B.27.
[51] [2017] FCAFC 4 (19 January 2017).
[52] [2017] FCAFC 4 (19 January 2017) [63].
[53] Privacy Act 1988 (Cth), s 6(1).

newspaper or other publication that is, or will be, generally available to members of the public—regardless of form or payment of an access fee)[54] is encapsulated under the federal privacy laws. For those collecting personal information via ADCT for research purposes, there is often an emphasis on dissemination of results—with likely outputs involving the publication of results in academic journal or government report form (as was the case in all 28 sources). Additionally, in the recent case of *Facebook Inc v Australian Information Commissioner*,[55] the court provided clear guidance on the fact that 'a user's device upon which personal information has been stored by means of a cookie is a record which contains personal information.'[56] It is possible to construe that the data gathered by ADCT forms a record by virtue of the fact that results are compiled on an electronic device.

In practice, this means that where researchers or practitioners are undertaking joint research involving APP entities, there is no contravention of the federal privacy framework where ADCT is used to collect information which is not 'personal' or 'sensitive' as defined in the *Privacy Act*, or where such information has been de-identified. Further, where personal or sensitive information has been collected (scraped), it will not breach federal privacy laws where the collection of that information is reasonably necessary for, or directly related to, the organization's functions or activities (ie for research purposes), and is collected via lawful means (ie the data are not restricted computer data accessed without authority)[57] from individuals (where practicable). Where these requirements have not been complied with and the information is published (such as in a magazine, book, article, newspaper or other publication that is, or will be, generally available to members of the public—regardless of form or payment of an access fee) or is a record (data stored on an electronic device), a breach is likely.

APP 6 provides that personal information may only be used or disclosed for the purpose for which it was collected (the 'primary purpose'), or for a secondary purpose if an exception applies.[58] 'Use' is not defined in the statute, however, the OAIC Guidelines suggest that an APP entity uses personal information when it handles and manages that information within the entity's effective control.[59] Examples provided include the entity accessing and reading the personal information; the entity searching records for the personal information; the entity making a decision based on the personal information; the entity passing the personal information from one part of the entity to another; the occurrence of unauthorized access by an employee of the entity.[60] In regard to use of data, it has been suggested that 'use' includes data analytics—running analysis on the data or seeking to match it with other data sets.[61] Analysis of data collected via ADCT occurred across all 28 sources, and is an imperative part of the research process.

Consent to data use (eg where express or implied consent is provided) or disclosure for a secondary purpose vitiates breach.[62] Further, where sensitive information is directly related to the primary purpose or the information is not sensitive information related to the primary purpose, there is no breach. Exceptions of particular relevance to a cybersecurity research context (eg where exposure to criminal activity may occur through data scraping) relate to use or disclosure of the information where required or authorized under Australian law (in accordance with statutory requirements) or a court/tribunal order,[63] or where the APP entity reasonably believes use or disclosure of the information is reasonably necessary for law enforcement activities conducted by or on behalf of an enforcement body.[64]

---

[54]  ibid.
[55]  [2022] FCAFC 9 [75].
[56]  ibid [161].
[57]  See, for example, Criminal Code Act 1995, Div 478.
[58]  Privacy Act 1988 (Cth), APP 6.1.
[59]  OAIC (n 50) B.146.
[60]  ibid.
[61]  Bennett Moses and others (n 9) 35.
[62]  Privacy Act 1988 (Cth) APP 6.1 (a).
[63]  ibid, (b).
[64]  ibid, (e).

In light of the ease with which information can be transferred and disclosed to overseas based servers/internet platforms, it is highly relevant to consider requirements for disclosure of personal information to overseas recipients. Similarly 'use' and 'disclosure' are not defined in the legislation and assumes its ordinary meaning, namely, to release or share. Further, inherent in the meaning of disclosure for the purposes of the privacy framework, is the concept of loss of control: disclosure of personal information to an overseas recipient means that the APP entity no longer has control over how that information is subsequently handled by the recipient.[65] Given the nexus between disclosure and loss of control of the information, there are particular requirements an APP entity must comply with to protect the privacy of an individual's personal information when disclosing that information to overseas recipients.[66] Notably, an APP entity must take reasonable steps to ensure the recipient does not breach the APPs in relation to that shared information.[67] De-identifying the information is one way an entity can demonstrate the taking of reasonable measures. However, taking reasonable steps does not preclude the entity from being held accountable for practices by the overseas entity which would breach the APPs.[68] There is, thus, an element of legal risk pertaining to disclosing personal information with cross-border recipients.

An exception to this privacy legislation, which is highly relevant for researchers or practitioners working jointly with APP entities on projects where information is shared with overseas recipients for research purposes, provides that reasonable steps ensuring the overseas recipient's compliance with the APPs do not need to be taken where a reasonable belief that the information is protected in a substantially similar way in which the APPs protect the information can reasonably be formed.[69] This greatly limits the legal risk where disclosure of personal information occurs between Australian researchers (working jointly with an APP entity) and countries such as Canada or those belonging to the European Union where personal information is protected in a similar manner as per the Australian privacy framework.[70] International collaboration is commonplace in research. For example, in only 14 of our 28 empirical sources, the researchers were based solely in Australia. The remaining 14 sources involved collaboration between Australian researchers and those in Canada and the EU (Switzerland, Sweden), as well as the UK, US, China, Singapore and Bangladesh. Consent to the disclosure by the individual to whom the data belongs/pertains to is another exception.[71] Further, and most relevant to practitioners and law enforcement are the following exceptions: where disclosure is required or authorized under statute or a court/tribunal order;[72] where disclosure is required or authorized under an international agreement relating to the sharing of information to which Australia is a party;[73] where the APP entity reasonably believes that disclosure is reasonably necessary for an enforcement related activity conducted by, or on behalf of, an enforcement body and the recipient is a body that performs functions, or exercises powers, that are similar to those performed or exercised by an enforcement body.[74] These exceptions give APP entities, such as the AFP, wide scope to disclose personal information relating to alleged crimes (eg cybercrime). As previously discussed, four of our empirical sources involved collaboration with the AFP.

---

[65] OAIC (n 49), 2.6.
[66] Privacy Act 1988 (Cth), s 16C and APP 8.
[67] ibid, APP 8.1.
[68] ibid, s 16C and APP 8.1.
[69] ibid, APP 8.2(a).
[70] Noting this is a generalisation derived from examining the legal frameworks in a cursory manner rather than detailed comparative legal analysis.
[71] Privacy Act 1988 (Cth) APP 8.2 (b).
[72] ibid, APP 8.2 (c).
[73] ibid, APP 8.2 (e).
[74] Privacy Act 1988 (Cth), APP 8.2 (f).

*Technological and financial harms*

The deployment of ADCT have the potential to cause technological harm to websites (disrupting or blocking legitimate traffic) and, as a consequence, financial harm to website administrators.[75] To combat this, websites (and their administrators) often insert provisions within their Terms of Use or via a robots exclusion standard (robots.txt)—that stipulate conditions around downloading, server request rates, and bypassing CAPTCHA to avoid overloading a website.[76] Only one of the sources identified in our search reported on directly navigating such requests through programming the ADCT to abide by any identifiable robots.txt file appearing on the website being examined.[77] However, we anticipate that a significant proportion of the included sources would have encountered and been forced to deal with this issue. This issue is discussed further in the below *Legal risks to the researcher* section.

*Legal risks to the researcher*

**Breaching website policies (Terms of Use and robots.txt)**

Some websites expressly prohibit web scraping, data mining and other ADCT capabilities through the use of a robots exclusion standard (robots.txt) protocol or directions within their Terms of Use policy. Where researchers ignore or misinterpret these prohibitions, they may be subject to litigation through a potential breach of contract claim where (i) the Terms of Use are sufficiently clear in prohibiting such activity, and (ii) the agreement is an enforceable agreement. Given the rich data sources comprised on social media sites, Terms of Use for sites such as Instagram, Facebook, Twitter, LinkedIn and YouTube have been crafted in a manner which makes the terms sufficiently clear, ensuring that a reasonable user is made aware of a prohibition against scraping of content unless express *consent is* sought and *provided.*[78] Further regulation of conduct which may be associated with ADCT (such as creating fake profiles or collecting data in unauthorized ways) may also be encompassed within broader terms prohibiting certain conduct more generally.[79] However, the legitimacy of litigation risk stemming from a breach of Terms of Use are hotly contested in the literature.[80] That is, many researchers believe these guidelines are not legally enforceable[81]—particularly when they appear on websites where illegal activity is being conducted (ie illicit marketplaces, CSAM forums, etc.).

There are three main types of internet agreements which are relevant when considering whether an internet agreement, such as those of social media sites, is enforceable. 'Click-wrap' agreements require the user to actively agree to the Terms of Use (usually by clicking on an 'I accept/I agree' button acknowledging terms), thereby providing the user with reasonable notice of the terms (meaning the terms are enforceable).[82] 'Sign-in wrap' agreements require users to click a button or sign-in in order to access the site after receiving notice of the terms via a hyperlink. Access to the terms prior to sign in is regarded as providing reasonable notice of the terms (terms are validly incorporated into the agreement) meaning that the terms of the

[75]    Brewer and others. (n 13); Thelwall and Stuart (n 4).
[76]    Capriello and Rossi (n 4); Filippo Menczer, 'Web Crawling' in Biong Liu (ed), *Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data* (Data-Centric Systems and Applications, Springer, 2011).
[77]    Ball and Broadhurst (n 2).
[78]    Instagram Terms of Service, General Conditions, cl 10 <https://help.instagram.com/478745558852511; Facebook Terms of Service cl 3.2.3> accessed 24 May 2023; Facebook Terms of Service <https://www.facebook.com/terms.php; Twitter Terms of Service.> accessed 24 May 2023; Twitter Terms of Service ≤https://twitter.com/en/tos> accessed 24 May 2023; LinkedIn Terms of Service <https://www.linkedin.com/legal/user-agreement#obligations> accessed 24 May 2023; YouTube Terms of Service, cl 4H <https://www.youtube.com/t/terms> accessed 24 May 2023.
[79]    Instagram Terms of Service, Basic Terms, cl 10 <https://help.instagram.com/478745558852511> accessed 24 May 2023.
[80]    Amber Zamora, 'Making Room for Big Data: Web Scraping and an Affirmative Right to Access Publicly Available Information Online' (2019) 12 JBEL 203.
[81]    See, for example, Martin and Christin (n 4).
[82]    *eBay International AG v Creative Festival Entertainment Pty Ltd* [2006] 170 FCR 250.

agreement are enforceable.[83] Alternatively, 'Browse-wrap' contracts do not require the user to actively confirm acceptance of the terms and simply provide links to subpages within the website or hyperlinks to a separate webpage containing the terms of use. Practically, this means that a user may not be provided with notice of the Terms of Use prior to entry into the agreement making the terms unenforceable (since the terms do not form part of the contract).[84]

Very little case law provides judicial guidance on the violation of a website's Terms of Use through use of automated data collection technologies. To date, no Australian cases have specifically dealt with this issue. A case which lends some insight is *hiQ Labs, Inc. v. LinkedIn Corp* ('hiQ Labs case'),[85] a matter which has been ongoing over the past five years. In April 2022, the U.S. Ninth Circuit Court of Appeals ruled that data scraping public websites is not unlawful and that LinkedIn cannot stop its competitor, hiQ Labs, from scraping LinkedIn users' publicly available data. Although the case is not binding in Australia, it is instructive for a couple of key reasons, namely, (i) there is a distinction between data which is publicly accessible (open-source data in its purest form) and data requiring a log-in (closed-source data), and where Terms of Use prohibit scraping (or other means of automated data collection), and (ii) the finding that publicly accessible data can lawfully be scraped does not preclude a breach of contract claim in the relevant jurisdiction.

### Breach of copyright

Across all 28 sources, authors extracted data using ADCT that may be subject to copyright, and therefore placed those researchers and practitioners at risk of litigation. That is, across sources, the data collected using ADCT had the potential to include original sales listings, media files or written posts published online (see Table 2 for type of data extracted using ADCT). Copyright in a work is the exclusive right to reproduce, publish, communicate, perform and adapt a work.[86] Copyright in works is infringed where a person not being the owner, and without licence of the owner, does an act comprised in the copyright.[87] Where work is reproduced in a material form (includes any form whether visible or not of storage of the work),[88] copyright is infringed where a *substantial part* of the work is reproduced.[89] This means that use of data falling within the scope of a work in which copyright exists may infringe copyright where a significant part of the data is replicated or adapted. In particular, it has been suggested that the copying of social media feeds in the context of data collection, copying an individual's social media or web profile to use as evidence, and replicating social media data to create and train analytic tools would prima facie infringe copyright.[90] This particular type of data was collected through ADCT across eight of the empirical sources identified in this paper.

Importantly, for individuals undertaking research involving automated collection of data, the statutory exemption of fair dealing for the purposes of research or study may apply.[91] Within the *Copyright Act 1968* (Cth) ('the Copyright Act'), sections 40(2) and 103C(2) stipulate that use or adaption of copyright items will constitute a fair dealing for the purpose of research or study only if the dealing includes the purpose and character of the dealing; the nature of the work or adaptation; the possibility of obtaining the work or adaptation within a reasonable time at an ordinary commercial price; the effect of the dealing upon the potential market for, or value of, the work or adaptation; and in a case where part only of the work or adaptation is

---

[83] *Dialogue Consulting Pty Ltd v Instagram, Inc* [2020] FCA 1846.
[84] *Specht v Netscape Communications Corporation* 306 F 3d 17 (2nd Cir 2002).
[85] 938 F.3d 985 (9th Cir. 2019).
[86] Copyright Act 1968 (Cth), s 31.
[87] ibid, s 36.
[88] ibid, s 10.
[89] ibid; IPC Global Pty Ltd v Pavetest Pty Ltd (No 3) [2017] FCA 82.
[90] Bennett Moses and others. (n 9) 44.
[91] Copyright Act 1968 (Cth), s 40.

reproduced—the amount and substantiality of the part copied taken in relation to the whole work or adaptation.

The Australian Copyright Council explains that to determine whether use of copyright material is considered fair dealing, it is necessary to look at the purpose, which must be one of the purposes set out in the *Copyright Act* and the use must be fair.[92] Whether a particular use is fair will depend on the case and the surrounding circumstances. It has been suggested that a number of general factors ought to be taken into account in the determination of 'fair dealing'. This includes, but is not limited to:

> [T]he amount of the copyright material used in comparison with the length of the copyright material; the extent of the use made of the copyright material by the defendant; the motives of the defendant; whether the copyright material is confidential or has not been disclosed to the public; whether the parties are in commercial competition with each other in respect of the use of the copyright material in question and the way in which it has been used; whether the use of the material is reasonably appropriate, rather than necessary, for the permitted purpose; and the relevance of, and weight to be given to, any industry practices or agreements between commercial rivals as to what constitutes a 'fair' dealing.[93]

The purpose of research or study is not explicitly defined within the legislation. The case of *De Garis v Neville Jeffress Pidler Pty Ltd*[94] established that within the context of the *Copyright Act*, research and study should take their respective ordinary meanings. The judgment utilized the Macquarie Dictionary definition of research defined as 'diligent and systematic enquiry or investigation into a subject in order to discover facts or principles'.[95] Beaumont J also used the Macquarie Dictionary definition of study, defined as:

> [A]pplication of the mind to the acquisition of knowledge, as by reading, investigation, or reflection; the cultivation of a particular branch of learning, science, or art; particular course of effort to acquire knowledge; a thorough examination and analysis of a particular subject…[96]

### Possession and distribution of illegal or high value content

While ADCT can be programmed to crawl specific websites and types of files, researchers have no control over the type of material contained on those websites or available files, potentially placing researchers at risk of a number of harms when automatically scraping data. In particular, the use of ADCT may lead to the identification and incidental possession of content that could be classified as illegal (eg CSAM,[97] suicide related material)[98]—particularly when deployed on the dark web where this material is more likely to appear (as was the case in 14 sources).[99] Conversely, the identification and collection of these types of materials may be within the

---

[92] Australian Law Reform Commission, *Copyright and the Digital Economy,* (Issues Paper 42 (IP 42), 2012) 64 [240]; Australian Copyright Council, *Fair Dealing in the Digital Age: A Discussion Paper* (1998), 20.

[93] Michael Handler and David Rolph, '"A Real Pea Souper": The Panel Case and the Development of the Fair Dealing Defences to Copyright Infringement in Australia' (2003) 27 MULR 381.

[94] [1990] 37 FCR 99.

[95] ibid [25].

[96] ibid [32].

[97] Criminal Code Act 1995 (Cth), s 273.6; 471.20; 474.22A; 474.23.

[98] ibid, s 474.29B.

[99] Ball and Broadhurst (n 2); Dalins, Wilson and Carman (n 1).

specific research project's aims—leading to deliberate possession of illegal content (as was the case for one source where CSAM was collected).[100]

Along with the possession of illegal material comes the risk that 'high value' material (eg closed-source information on closed websites, CSAM) may be inadvertently or involuntarily distributed by the researchers. This can occur in situations where the researchers' ADCT is identified or located by individuals or websites with malicious intent who wish to intercept such data. Regardless of the intent behind the collection of this content, there are various associated risks for those carrying out the collection, including legal.

### *Risks of harm to the researcher*

Where graphic or illegal content is possessed as part of the research, there is a risk that psychological harm may occur when reviewing materials. Exposure to material depicting serious harm, or a corpus of text describing such harms, may result in the researcher experiencing distress, or potentially severe mental health implications.[101] This risk was directly present in one source, but may be an incidental risk for researchers deploying ADCT on platforms where there is high risk of exposure to graphic or illegal content (ie the dark web).

ADCT may include an identification string (ie crawler/scraper's name, researcher contact information or links to information about the research) to allow transparent communication with the websites being crawled—informing them that such activity has occurred and providing an opportunity to contact the researchers.[102] Five of the sources identified in our search explicitly reported providing such details. However, it is acknowledged that the practice of web scraping is not always welcome, and there may be circumstances where the researchers wish to avoid identification. For example, some website users (administrators, vendors, buyers) may take issue with data collection activities, and either change their behaviours (leading to biased data), or potentially seek retribution via threats or offline/cyber abuse.[103]

### Identifying issues related to consent

Potential ethical issues pertaining to gaining consent from the subjects at the centre of the research involving ADCT were either reported, or expected based on review of methodology, in all 28 sources (see Table 6). That is, we submit that all of these sources used ADCT to extract data created by, or pertaining to, human subjects, and depict circumstances where the collection of informed consent was not possible. These two aspects are elaborated upon in turn.

**Table 6:** Number of Sources with Reported or Expected Consent Issues

| Ethical challenge | No. of sources |
| --- | --- |
| Research involves human subjects | 28 |
| Obtaining informed consent is impractical | 28 |

---

[100]  Dalins, Wilson and Carman (n 1).
[101]  Brian Pitman and others, 'Social Media Users' Interpretations of the Sandra Bland Arrest Video' (2019) 9 Race and Justice 479.
[102]  Ball and Broadhurst (n 2); Menczer (n 76).
[103]  Décary-Hétu and Aldridge (n 6); Thomas J Holt and others, 'Advancing Research on Hackers through Social Network Data' in Catherine D Marcum and George E Higgins (eds), *Social Networking as a Criminal Enterprise* (Taylor Francis, 2014); Martin and Christin (n 4).

### Human subjects

According to the *National Statement*, voluntary and informed consent is the foundation upon which ethical research involving *human* subjects can take place. The specific nature of a research project may determine the manner by which researchers are required to obtain consent from subjects. This is often influenced by relevant 'codes, laws, ethics and cultural sensitivities of the community in which the research is to be conducted'.[104] Research involving ADCT does not typically engage directly with human subjects, but does involve the large-scale collection and analysis of data previously produced by humans and made available online.[105] As indicated above, all 28 sources scraped and analysed human-produced data that had been posted online and was accessible by ADCT.

The *National Statement* outlines that a subject's participation in research must be in accordance with their voluntary and informed consent, and that the public accessibility of such information does not indicate that an individual subject has consciously provided permission for their data to be used within research. However, the *National Statement* also acknowledges potential difficulties in gaining the consent of human subjects when collecting this type of online data – which may be considered a 'secondary use of data or information'.[106] Therefore, exemption from the requirement to gain consent from human subjects may be possible in some lines of research where it is ethically and legally justified. Accordingly, it is recommended then that researchers planning to use ADCT need to determine (i) whether their research uses human subject data, and (ii) whether obtaining informed consent from those human subjects is possible/feasible (described further below).

### Determining whether a waiver of consent is appropriate

The sources identified in our systematic search included a wide collection of human subject data including forum discussions, social media posts and product listings, with thousands of different data points from various users collected (see Table 7). Considering the nature of this data, it is reasonable to suggest that researchers across all 28 sources would have faced significant barriers in obtaining the direct consent of subjects. This was due to the lack of direct interaction between researchers and subjects, the asynchronous nature of data extracted, and the sheer volume of individual users whose data were collected. In addition, in the vast majority of these cases, users were not identifiable (ie use of an alias/username), and may not have been currently active on the website. As a result, expectations of obtaining consent for the collection and analysis of such data, before or after the research occurs, are impractical in reality.[107]

These limitations are explicitly acknowledged in the *National Statement*, which highlights that a waiver for consent may be possible, given the aforementioned impracticalities, but this must be considered by the relevant reviewing body.[108] According to the *National Statement*, researchers may be exempt from seeking consent of human subjects, so long as they abide by the following Conditions: (A) the research involves minimal risk of harm to subjects, (B) the benefits of the research justify any risks of not seeking consent, (C) it is impractical to obtain consent (eg the aforementioned issues identified as part of ADCT research), (D) there is no likely reason that subjects would not consent if provided the opportunity, (E) subject privacy is protected, (F) confidential data are protected, (G) results are made available to the subjects where they may have significance, (H) subjects would not be deprived of financial benefits they may be

---

[104] NHMRC (n 11), 16 [2.1].
[105] Batistic and others, (n 43).
[106] NHMRC (n 11), 36.
[107] Décary-Hétu and Aldridge (n 6).
[108] NHMRC (n 11), s 2.3.9.

**Table 7:** Number of Sources by Subject Type

| Subject type | No. of sources |
|---|---|
| Darknet vendors | 13 |
| Darknet customers | 6 |
| Darknet administrators | 18 |
| Darknet users | 6 |
| Surface web vendors | 3 |
| Surface web customers | 1 |
| Surface web administrators | 6 |
| Surface web users | 2 |
| Social media users (incl. forums, blogs) | 6 |
| Social media administrators (incl. forums, blogs) | 6 |
| Pornography actors | 1 |
| Unknown individuals depicted in media | 1 |

entitled to through commercialization of the data, and (I) the waiver of consent is not prohibited by federal, state or international law.

Jurisprudential review highlighted that *the DATA Scheme* further supports a number of conditions giving rise to a waiver of consent. Under *the DATA Scheme*, 'Research and Development' is an authorized purpose for data sharing,[109] noting that the data must not include personal information about an individual unless: by consent (and only a minimum amount of personal information is shared) *or* in circumstances where the project cannot proceed without the information, it is in the public interest, only the minimum amount of personal information is shared *and* a permitted circumstance exists.[110] A permitted circumstance for the purpose of research and development is that it is unreasonable or impracticable to seek consent.[111] Importantly, this requirement has a high threshold test—it must not be merely inconvenient to gain consent of a very large number of individuals.[112] If a claim that an unreasonable or impracticable assertion exists, the data custodian must provide a statement that personal information is being shared without consent of individuals because of its impracticality and an explanation of the reasons for concluding so.[113]

Furthermore, the *Privacy Act* and APPs require that an individual's consent be obtained during the collection of 'sensitive' information *unless* an exception applies.[114] The recent case of *Commissioner initiated investigation into Clearview AI, Inc. (Privacy)* ('*Clearview case*')[115] referenced the critical importance of this requirement in regard to facial recognition technology which was construed as being biometric information and, thus, sensitive information.[116] The exceptions are highly relevant to law enforcement agencies which fall within the ambit of an APP entity. Collection of sensitive information can occur where collection is required or authorized by Australian law or a court order, or there is a permitted general situation in relation to the

---

[109] Data Availability and Transparency Act 2022 (Cth), s 15(1)(c).
[110] ibid, s 16B(3)(a)(b).
[111] ibid, s 16B(4)(a).
[112] ibid, s 16B(4).
[113] ibid, s 16B(7)(a)(b).
[114] Privacy Act 1988 (Cth), APP 3.4(a).
[115] [2021] AICmr 54 (14 October 2021).
[116] ibid, [137].

collection of the information.[117] Further, collection of sensitive data can occur if the enforcement entity reasonably believes that collection is reasonably necessary for, or directly related to, one or more of the entity's functions or activities.[118] Across our empirical sources it is unclear how many researchers may have inadvertently scraped sensitive data. However, within two sources, the researchers explicitly identified that sensitive material in the form of adult pornography, and CSAM, were collected, in partnership with AFP personnel using ADCT.[119] In both of these cases, gaining the subject's consent would have been impossible for the researchers—particularly given the materials were likely initially uploaded without the original subject's consent.

### Other considerations

Australian researchers and practitioners need to be mindful of cross-jurisdictional legal frameworks that may have bearing on ADCT activities. Collection and use of data collection in an overseas jurisdiction may be governed by the originating jurisdiction as well as relevant Australian law and policy requirements (ie data transparency, privacy, Terms of Use and copyright). By way of example, our systematic search identified two sources that used ADCT based in Canada. Collection and use of this data is governed by the *Privacy Act*, RSC 1985[120] (governing the Federal government's access to and use of personal information), the *Personal Information Protection and Electronic Documents Act,* SC 2000[121] (governing the collection and use of personal information by private parties), the *Copyright Act,* RSC 1985[122] (governing the use of copyrighted material), as well as other Provincial privacy and data management related legislative requirements. In addition to legal frameworks, research being undertaken by Australian researchers or practitioners containing personal or sensitive information will likely need to comply with Australian and other originating jurisdiction's ethical requirements. In these Canadian source examples, this would require the researcher to ensure that the data collection and use practices are consistent with Australia's *National Statement on Ethical Conduct in Human Research,*[123] as well as Canada's *Tri-Council Policy Statement: Ethical Conduct for Research Involving Humans (2018).*[124] In particular, and as noted in the above sections, researchers must be sensitive to matters pertaining to consent as well as potential risks/harms to researchers and research subjects.

## CHARTING A PATH FORWARD: MITIGATING ETHICAL AND LEGAL CHALLENGES

We conclude by offering strategies to combat these ethical and legal challenges for Australian-based researchers and practitioners deploying ADCT to carry out research within a cybersecurity context (see Table 8 for the mitigating measures identified as taken by each of the 28 empirical sources). Importantly, we highlight that the use of ADCT may be undertaken both ethically and legally provided that various measures are taken at all stages of data collection, storage and reporting, and that these measures are considered in the context of the research aims and the location where data collection is undertaken. Depending on the nature of the deployment

---

[117] Privacy Act 1988 (Cth), APP 3.3. A permitted general situation is explained in s 16A, where there is a list of general conditions that can apply—including the need to collect, use or disclose to lessen or prevent serious threat to life or health and safety, issues of unlawful activity, and issues of international concern such as war operations, peacekeeping or humanitarian assistance.

[118] ibid, APP 3.4(d)(ii).

[119] See Dalins and others. (n 1); Dalins, Wilson and Carman (n 1).

[120] Privacy Act 1988 (Cth), RSC 1985, c. P-21.

[121] c. 5.

[122] c. C-42.

[123] NHMRC (n 11).

[124] Canadian Institutes of Health Research, Natural Sciences and Engineering Research Council of Canada, and Social Sciences and Humanities Research Council of Canada, *Tri-Council Policy Statement: Ethical Conduct for Research Involving Humans*, December 2018.

**Table 8:** Number of Sources Taking Mitigating Measures

| Mitigating measure | No. of sources |
|---|---|
| Program ADCT to avoid certain data/websites | 17 |
| De-identify data | 28 |
| Report aggregate data | 28 |
| Implement strong data security practices | 10 |
| Utilize API where possible | 2 |
| Ensure benefits of research outweigh potential risks | 28 |
| Seek research approval from relevant ethics committee | 8 |
| Partner with law enforcement | 4 |
| Scrape metadata | 7 |
| Target lower value (closed-source) data | 7 |
| Use automated classification models for categorizing data | 6 |

of ADCT (ie the purpose of research and the platform being scraped), there are certain risks of harm that become more pertinent for consideration. For example, privacy concerns for research subjects may be greater where personal and sensitive data is likely to be collected (eg where digital communications and media file data are being targeted); and researchers may be placed at greater risk of psychological harm or litigation when scraping platforms known to carry illegal materials (eg dark web marketplaces). Additionally, the applicable legal framework depends on the researcher affiliation (ie whether involved with a public or private entity determines the application of federal or state-based laws), the location of the researchers within Australia (ie for adherence to the relevant state-based laws or guidelines), as well as whether the project involves any international collaborators (ie the relevance of international laws).

The considerations and directions raised in the current review are fairly comprehensive. However, it is important to note that given the speed at which technology and online platforms are constantly evolving, it is unclear how these ethical and legal issues may become shaped in the future regarding the use of ADCT. Therefore, ongoing consultation with the relevant ethical boards to remain up-to-date on any further considerations for achieving ethical and legal conduct in research using ADCT is recommended.

### Direction to mitigate against privacy violations

Sometimes it is unclear whether data collection of certain websites may be inappropriate, or whether the technology may automatically crawl unknown websites and potentially violate the privacy of individual users. Therefore, one method to mitigate the risk of unwanted collection is to program the crawler to not 'generate user credentials, complete captchas or otherwise obtain access to restricted sites'.[125] In 17 sources, researchers explicitly reported employing such methods to restrict the ADCT, or restricted their use of ADCT to one specific area of a website.

In general, researchers should avoid collecting data which may be considered personal or sensitive. However, where that type of data must be collected, it will be imperative for those deploying ADCT to demonstrate that the data collection is needed to pursue particular research outcomes and show, by way of a detailed research proposal, how the data will be assessed as relevant and not excessive to the organization's functions—in line with the *Privacy*

---

[125] Dalins, Wilson and Carman (n 1).

*Act*. Furthermore, it is essential that APP entities are clear about the primary purpose of ADCT before it can be used.[126] Where data have been purchased under a contract or license agreement, the terms ought to include specifics on purpose of data use and how it may be handled (eg which tools may be utilized when handling/processing the purchased data).[127]

Elsewhere, researchers can significantly reduce any risk of harm to subjects and protect their privacy by de-identifying and removing all personal information from any data that is collected, stored, analysed and reported—and ensuring similar practices are undertaken by any overseas agencies the data are shared with. Such practices are commonplace throughout the literature, and were carried out across all 28 sources. This included the removal and anonymization of any usernames or pseudonyms, website names or identifying quotes or information. We acknowledge, however, that anonymization may not be universally possible across all stages of ADCT research—for example, where media files are deliberately or incidentally downloaded.[128] Therefore, in these cases, and in addition to those studies where anonymization and de-identification occurs, only aggregated data should be reported on (a practice which was implemented across all 28 sources). Considerations around the storage and transmission of sensitive data is also of paramount importance. Best practice advocates that scraped data should be stored separately from the crawler itself during data collection, with strong security measures implemented to protect against any intrusions from external parties (eg encryption, limited remote access, IP-based restrictions).[129] The extent to which such practices are implemented by researchers using ADCT is unclear, with only 10 of the identified sources explicitly reporting strong storage security practices.

It is important to flag that subject privacy may need to be breached in cases where the collected data may uncover illegal activity. In Australia, researchers are mandatory reporters, and must ensure that any criminal behaviour that is exposed during the research (ie the identification of CSAM), be disclosed to the relevant authorities.[130] Alternatively, researchers may be subject to orders to disclose information to relevant authorities following aggregated reporting of results.[131] Researchers should be sensitive to such requirements, as they may involve the re-identification of subjects to assist in law enforcement investigation. As outlined in the *Privacy Act* APP 6, a breach of privacy is permissible where it is necessary for the pursuit of disclosing identified criminal activity for law enforcement purposes (even if that is a secondary purpose for the use of data). Of particular relevance to the use of ADCT within the cybersecurity contexts focused on in this paper, are the requirements outlined in Chapter 4 of the *National Statement*: that research must be reviewed by an approved Human Research Ethics Committee ('HREC') if it is intended to study or expose illegal activity, or likely to discover it.[132] Therefore, it is recommended that ethical approval is sought from the relevant review board, highlighting where the aim of the research is to examine illegal behaviour from the data scraped using ADCT, or where the location the ADCT is deployed may be likely to contain illegal material (ie the dark web).

### Direction to mitigate against technological or financial risks

We recommend here that researchers should proceed with caution and take the context into account when designing their tools so as not to unduly harm a system. Nevertheless, in making this recommendation, we acknowledge that Terms of Use and abiding by the provisions in the

---

126  Privacy Act 1988 (Cth), APP 6.
127  ibid.
128  See, for example, Dalins and others. (n 1).
129  See, for example, Dalins, Wilson and Carman (n 1).
130  NHMRC (n 11).
131  ibid.
132  ibid, 75 [4.6].

robots.txt may not be legally enforceable, and discuss these aspects further in the *Direction to mitigate against litigation associated with breaches of Terms of Use* section below.

### Direction to mitigate against litigation associated with breaches of Terms of Use

Researchers and practitioners who wish to utilize ADCT should consider obtaining express consent to undertake ADCT from the platform being targeted for collection. Where this is not possible, it should be considered whether a website's Terms of Use prohibit automated collection of data. If that is the case, be aware that the company may enforce the term (commence civil litigation) where (i) the website becomes aware of the breach; (ii) the website owners decide it is in their interest to initiate legal action (they may not instigate proceedings for a plethora of reasons). Furthermore, consideration of the type of internet agreement specific to the website is important, as this can influence whether terms are enforceable. If the agreement is a browse-wrap internet contract, it is unlikely to be enforceable by the website. On the other hand, click-wrap and sign-in wrap internet agreements are likely to be enforceable by the website. Assessing the legal risk of breaching the Terms of Use—based on the recent *hiQ Labs* ruling—legal risk is likely to be minimized where the data being scraped is publicly available and where access is not protected by login requirements.

As was the case for mitigating risks associated with a waiver of consent, a breach of Terms of Use may be justified on the grounds that the benefits of the breach of such guidelines for research purposes (eg to detect/understand illegal activity) outweigh the risks.[133] In particular, websites containing illicit content are highly unlikely to enforce any Terms of Use. However, as outlined in the above *Risk to the Researcher – Possession and distribution of illegal content* section, researchers will need to be alive to potential criminal laws which may apply to the possession of data comprising of illicit content, such as CSAM.

Alternatively, one way of preventing potential breaches of a Terms of Use agreement and circumventing the need for express consent to access data, is to utilize an Application Programming Interface (API) approved by the relevant site which authorizes a user to access certain data. Use of an API is likely to require the user to agree to a terms of data use which affects treatment of the data. For example, Facebook's 'Meta Platform Terms' prohibit user attempts to 'decode, circumvent, re-identify, de-anonymize, unscramble, unencrypt or reverse hash, or reverse-engineer Platform Data that is provided'.[134] It is suggested then that the risk of litigation may be mitigated by using such software when it is available[135]—which was the case in two of the sources identified in our search.[136]

### Direction to mitigate against breaches of copyright

Noting the judicial guidance on this point, it is likely that data collected by automated means for legitimate research purposes, such as those projects approved by university ethics boards; those which have approved internal or external project funding; and those which have a purpose which aims to advance knowledge in a particular way, will fall into the fair dealing for the purposes of research or study exemption. As discussed further in the below *Direction to mitigate risks related to consent* section, the research carried out across all of the empirical sources identified in our search would arguably provide valuable advancement of knowledge from a cybersecurity perspective. While only eight sources explicitly reported approval from the relevant

---

[133] Deen Freelon, 'Computational Research in the Post-API Age' (2018) 35 Polit. Commun. 665; Martin and Christin (n 4).

[134] Facebook Platform Terms, <https://developers.facebook.com/terms/> accessed 24 May 2023.

[135] Oliver C Stringham and others, 'A Guide to Using the Internet to Monitor and Quantify the Wildlife Trade' (2021) 35 Conserv. Biol. 1130.

[136] Guohui Li and others, 'Misinformation-oriented Expert Finding in Social Networks' (2020) 23 WWW 693; Md. Saidur Rahman and others, 'An Efficient Hybrid System for Anomaly Detection in Social Networks' (2021) 4 Cybersecurity <https://doi.org/10.1186/s42400-021-00074-w> accessed 24 May 2023.

university ethics board within the published source, this was likely much higher in reality—with the majority, if not all, of the sources likely to fall within the exemption. Risk of falling foul of the copyright law is, arguably, minimal.

### Direction to mitigate against illegal possession and distribution risks

Researchers must carefully consider the nature of the data being collected, and whether such collection, possession or analysis of said data may constitute a break of Federal or State criminal codes. In cases where this is likely to occur, the mitigation of the risk of criminal charges may be achieved through careful planning by the researchers. This may involve partnering with law enforcement agencies, or transparent notification/request for permission from Commonwealth agencies (ie Minister for Justice, Attorney General's Department, state authorities), in addition to consultation with the researcher's relevant ethical review board.[137] To this end, researchers in four sources reported successfully partnering with law enforcement to carry out the research— which included one source that was explicitly searching for illegal material online.

We acknowledge that this risk is more difficult to mitigate where the nature of research does *not* explicitly involve the collection of illegal material, and the possession of such material is not predictable. It is therefore recommended that researchers adopt a proactive approach, regardless of the research aims, by assessing the likelihood of incidental possession and planning to take the appropriate measures. This may include those aforementioned measures of notifying the relevant authorities of the research involving ADCT, and possible risk of possessing illegal materials as a result—particularly where known illicit websites are being scraped. Additionally, where possible, scraping metadata (ie filename, URL) rather than the media itself could serve as a solution to this issue.[138] Such practices were reported across seven sources. Additionally, we echo calls by other researchers[139] that the risk of distribution may be mitigated through only targeting lower value material (ie material that is publicly accessible) to deter those who may want to intercept valuable material, or separating the ADCT from securely stored data. Such practices were employed across various sources, including seven that targeted low value data, and 10 that implemented robust storage security practices.

### Direction to mitigate against psychological harm to the data collector

This issue was flagged in one of the sources identified in our search[140] which mitigated against the need to review CSAM by working with law enforcement partners who were appropriately trained to review this type of material. Elsewhere, six of the sources identified in our search used classification models that automatically categorized the data without the need for manual review.[141] Accordingly, it is recommended that researchers working with potentially distressing data should seek to implement techniques to obviate the need to manually review data—such as aforementioned partnerships with law enforcement or use of automated classification models.

### Direction to mitigate against inadvertent researcher identification

Revealing the presence and identity of researchers during the data collection process carries some risk. In light of this, researchers are advised to consider these circumstances, and where retribution or bias is anticipated, seek guidance and approval from their institutional review

---

137 Ball and Broadhurst (n 2); Dalins and others. (n 1).
138 Dalins and others. (n 1).
139 Dalins, Wilson and Carman (n 1).
140 Dalins, Wilson and Carman (n 1).
141 Ball and Broadhurst (n 2); Dalins and others. (n 1).

boards to operate discreetly.[142] However, as one of the identified sources flagged, such risks may, in reality, be overstated.[143]

### Direction to mitigate risks related to consent

Condition B for a waiver of consent outlined in the *National Statement* highlights that the benefits of the research must outweigh the risks. The potential benefits of research using ADCT in a cybersecurity context—such as those sources included in the current research—may involve important outcomes for understanding the state of certain crimes, and motivations for criminal behaviour, to enhance public protection and security. For example, all 28 sources examined human data for a range of outcomes relating to the advancement of cybercrime/cybersecurity knowledge. These involved developing new classification models for assisting law enforcement investigations into CSAM online, understanding the impact of policing the dark web to disrupt the trade of illicit goods, and understanding public reactions to criminal behaviour to inform the response of policy makers and practitioners. However, while the sources identified through our search all demonstrated significant benefits, we acknowledge that this may not be true for all future studies intending to employ ADCT.

Given the subjectivity of weighing up the potential benefits of research in comparison to potential risks, it is required that any research involving human subjects or data derived therefrom should involve ethical approval from an authorized institutional review board. Beyond considering the potential risks versus rewards (Condition B), we also provide treatments of Conditions A, C-F, and I (which relate to minimizing subject harm) in our above directions for mitigating risks to research subjects sections. In particular, in compliance with *the DATA Scheme* legislation, researchers must remove personally identifiable information when sharing data unless they have received express consent from the individual, only a minimum amount of personal information is being shared, or in circumstances where the project cannot proceed without the information. For biometric data in particular, this may not be shared unless consent has been provided by the individual to whom the data relates. Furthermore, in accordance with the *Privacy Act,* a waiver of consent for the collection of sensitive data may be satisfied if that collection is required or authorized by court order or Australian law, or if it is considered a reasonable requirement for the function of the particular agency carrying out the data collection. We submit that Conditions G and H (regarding providing subjects with significant results and financial benefits) are beyond the scope of this paper, given that consequential results and financial benefit for research subjects are both irrelevant to the sources identified within the systematic search, and likely not applicable to future research employing ADCT in a cybersecurity context.

## SUPPLEMENTARY MATERIAL

Supplementary material is available at *International Journal of Law and Information Technology* online.

---

[142] Stringham and others. (n 135); Michael Cheng-Tek Tai, 'Deception and Informed Consent in Social, Behavioral, and Educational Research (SBER)' (2012) 24 TCMJ 218.
[143] Ball and Broadhurst (n 2).