

# Using Twitter Data in Research

## Guidance for Researchers and Ethics Reviewers

**Dr Nicolas Gold**

Department of Computer Science, UCL

### Contents

1	Introduction .....	2
2	Scope of Guidance .....	2
3	Status of Guidance .....	3
4	Ethics Issues .....	3
4.1	Principles .....	3
4.2	Stakeholders and Participants .....	4
4.3	Legal Compliance for Ethics .....	4
4.4	Primary vs Secondary Data Analysis .....	5
4.5	Consent and Expectations .....	5
4.5.1	Twitter .....	5
4.5.2	Tweeters .....	6
4.5.3	Subjects .....	7
4.6	Privacy and Anonymity .....	7
4.7	Restrictions on Research .....	7
4.8	Sharing Data .....	8
4.9	Publishing Data .....	8
4.10	Backlash and UCL Reputation .....	8
5	Other Legal Matters .....	9
5.1	Contractual - Synchronisation .....	9
5.2	Contractual – Indemnification and Audit Requirements .....	10
5.3	Data Protection .....	10
6	Engagement with UCL Policies .....	10
6.1	Ethics .....	10
6.2	Data Protection .....	11
7	Summary .....	11
8	Reference Questions for Reviewers .....	11
9	References .....	12

### Version History

Version 0.1: (30/6/2020) Draft circulated to CS Ethics Committee, Legal Services, and UCL DPO for comment.

Version 1.0: (16/7/2020) Final draft incorporating comments received, submitted to UCL REC for comment and approval.

## 1 Introduction

Twitter was founded in 2006 to provide micro blogging services and is an established part of the contemporary social media landscape. It is a corporate entity and its interactions with users (those who post or consume its data) are governed by contracts and agreements.

Twitter data (among other things) consists of Tweets: short messages posted by its users to the site, often using official or third-party apps on mobile or fixed devices. Such messages may be accompanied by, or contain, metadata (data about the message such as the user's geolocation or the time it was posted), and are organised in terms of accounts. Users interact with Twitter via these registered accounts and can determine the public visibility of their Tweets at their discretion. Unregistered users may only view Tweets designated as public. Accounts may be 'followed' to receive all Tweets posted by that account. Accounts may be personal or corporate. Twitter data can be accessed by reading it on Twitter's own website, through the official Twitter app, or via the Developer API that Twitter makes available.

Twitter is widely seen as a useful source of data for many types of social, Information and Communication Technology (ICT), and other research as the data it provides is apparently publicly available, accessible easily and at low cost (at least on a small scale). However, engaging with Twitter data draws a researcher and their organisation into a contractual relationship with the company that can have a significant effect on the costs (financial or temporal) of the research activity, can have implications for the technical implementation of the research method, and that raises ethical and legal issues that must be addressed in the research design.

This document aims to draw out these implications for investigators wanting to use Twitter in their research at UCL, for ethics reviewers scrutinising applications proposing to use Twitter, and for the REC to inform its debate of the issues when proposals come up for discussion. The guidance discusses the ethical, legal, and practical aspects of engaging with Twitter as a secondary data source, the underlying ethical principles that these engage, and how and where the issues interact with UCL Ethics policy.

Potential research users of Twitter data need to consider a range of factors including:

- *Ethics issues* (principles, stakeholders and participants, legal compliance, primary vs secondary data analysis, consent and expectations, public/private online spaces, privacy and anonymity, restrictions on analysis, sharing data, publication, backlash and UCL reputation).
- *Legal constraints* (contractual, data protection) on the collection and use of Twitter data.
- *UCL policy requirements and interactions* (ethics review level, data protection registration).

These factors arise from various ethics frameworks (e.g. Menlo Report (Dittrich & Kenneally, 2012) and Association of Internet Researchers (Association of Internet Researchers, 2019)) and from Twitter's terms and conditions of service. Research design will need to account for these. The legal and ethical matters are in some cases entwined as Twitter's contractual terms themselves incorporate aspects of good ethical practice and require users to comply with this.

The document concludes with a set of questions that ethics reviewers may find helpful when considering a potential study.

## 2 Scope of Guidance

This document was produced in response to requests from the Computer Science Ethics Committee and the UCL REC for guidance on the issues raised by Twitter data analysis, given the large number of Twitter-related studies reported as submitted to the REC.

There is a large and growing literature in the area of social media ethics. Some is cited here to support the guidance and to direct the reader to further information but this is not a formal or comprehensive survey of such work.

The guidance presented here has no formal standing unless and until adopted by the UCL REC (this will be noted in the Version History above if appropriate).

The focus here is on data available to the public via Twitter: primarily Tweets and public account information. Studies of non-public data (e.g. private profile information, direct messages) and other Twitter services like Periscope (video live streaming) are not within the direct scope although some of the principles here may apply to these too. Whilst many issues raised may also apply in general terms to social media research on other platforms (e.g. Facebook), this guidance applies to Twitter only because of the dependency on Twitter's specific policies and practice.

In the course of developing this guidance, the author has consulted with (and is grateful to) colleagues in UCL Legal Services on matters of contract and data protection law pertaining to Twitter, and is likewise grateful for helpful discussions with, and feedback from, colleagues on the CS ethics committee. The author is not legally trained or qualified and nothing in this document should be construed as, or used as, legal advice. Where UCL legal advice is cited, it is on the basis of the author's best understanding of the advice received; any errors or misunderstandings are the author's. If clarification of the legal matters described herein is required, UCL Legal Services should be consulted.

Bespoke contractual agreements between UCL (or parts thereof) and Twitter directly may conflict with the general guidance offered here. If so, then the bespoke contractual relationship should be respected as the general terms of access may not apply in those circumstances. The ethical issues will still apply but the particular contractual terms may affect the arguments that can be made to resolve them.

### 3 Status of Guidance

This document is prepared with reference to UCL Ethics Policy and Twitter terms and conditions documents (cited later) in force on 18<sup>th</sup> June 2020. If these change, the revised policies and terms should be consulted and will take precedence in any conflict with the guidance herein, unless and until this guidance is updated to reflect the modified positions.

## 4 Ethics Issues

### 4.1 Principles

There are many ethics frameworks in which researchers can situate their work. It is not the purpose of this document to compare and contrast the relative merits of these in relation to Twitter, and researchers are free to select and work within the most appropriate framework for their work, and/or as required by their professional practice. However, it is worth noting the issues raised by frameworks that deal with ICT and social media data in particular, as these may only be implicit in other ethics codes but nonetheless need to be addressed.

Researchers in Computer Science at UCL are often directed to the Menlo Report (Dittrich & Kenneally, 2012) as a foundational ethics framework and indeed this is useful, albeit quite general. Of relevance to this guidance, two key aspects are important: the need to identify stakeholders in the research (who will likely go beyond the direct participants involved), and the need to incorporate legal compliance and the public interest alongside the Belmont principles of Respect for Persons, Beneficence, and Justice. The report also identifies that additional risks can arise because of the scope, speed, and scale of IT systems permitting access to data. The guidance presented in this document is largely situated in the Menlo Report framework.

Recent work has shown that few UK ethics frameworks deal with social media ethics issues directly (Taylor & Pagliari, 2018) identifying those that do as ESRC, BPS, AoIR, and NIHR in the context of research council frameworks.

The AoIR codes (there are three complementary versions) highlight many aspects of social media research and it is strongly recommended that researchers in this area are or become familiar with the issues and discussions identified there.

The key primary ethical principle that is engaged in relation to Twitter data is **autonomy**: the right of participants to be treated as independent moral agents and to determine their own best interests (Dittrich & Kenneally, 2012). This is not to say that other principles of ethics (beneficence, justice etc) are irrelevant and should not be considered, but simply to say that in practical terms, it is the consideration of autonomy (and its associated concepts of **informed**

**consent and participant expectations**) that pose the most specific complexities in the context of Twitter data analysis.

## 4.2 Stakeholders and Participants

The Menlo Report indicates the importance of stakeholder identification to ensure that all those involved in the research are accounted for ethically (Dittrich & Kenneally, 2012), and casts a wide net for this in the context of ICT research to ensure all those impacted by the research are considered. In the case of Twitter data research there are a number of stakeholders to consider:

- Twitter itself,
- those who post on Twitter (hereafter: *Tweeters*),
- those who are posted about on Twitter (hereafter: *Subjects*, may or may not be the same as *Tweeters*),
- those who consume Twitter data through reading its website and/or apps (hereafter: *Consumers*, may or may not be *Tweeters* and/or *Subjects*),
- those who consume Twitter data through its Application Programming Interface (API) (hereafter: *Harvesters*, may or may not be any of the above),
- the organisation that employs any of *Tweeters*, *Consumers*, or *Harvesters* (if relevant),
- the researchers,
- UCL.

## 4.3 Legal Compliance for Ethics

Specific aspects and research implications of compliance with Twitter's terms are dealt with the following sections. The ethical need for legal compliance is established by the Menlo Report (Dittrich & Kenneally, 2012) and since Twitter's terms capture good ethical practice in many ways, compliance with these may also simplify the ethical justification required for a piece of research. Compliance is also required to ensure that UCL is not in breach of contract with Twitter (and thus compliance falls within UCL's Sensitive Research Policy framework clauses on UCL reputational risk owing to potential breach of contract).

It is important that the legal relationships between the various parties involved in research using Twitter data are understood. The **entire basis of legitimate access to Twitter data rests on legal contracts** (note also that UCL's existing guidance on social media research requires adherence to the terms and conditions of a given service (UCL Research Ethics Committee, n.d.)).

Users of Twitter (in all the categories identified above) agree to Twitter's Terms of Service (Twitter Inc., 2020a) when signing up to get an account, or simply by browsing its content via its website or its apps. In doing so they make a legal agreement with Twitter either for themselves (if acting in a private capacity), or if their activity is being done in the course of employment then they do so on behalf of their employer. So **UCL staff accessing/analysing data, or browsing or posting to Twitter for work purposes are doing so under a contract between UCL and Twitter**, even if UCL was formally unaware such a contract had been formed on its behalf [UCL legal advice]<sup>1</sup>.

The **terms of the contract govern the data that is accessible** to the researcher (thereby potentially differentiating data that is legitimately available from that which is technically available). It is also important to note that the **terms restrict what can be done** with the data. The contract (both general terms (Twitter Inc., 2020a) and those specific to API use (Twitter, Inc, 2020b)) reflect Twitter's promises to its users about control and respect for their voice and these values flow contractually to *Consumers* and *Harvesters*.

Where a researcher wishes to undertake work with, or access data, in a way that falls outside the contract terms, it is important to note that UCL may then be at risk of breach of contract and the work would fall under the UCL Sensitive Research Policy (clause 8) and will almost certainly require REC review and approval before starting. Lawful use of the data does not eliminate all ethics issues however.

---

<sup>1</sup> The need for oversight arrangements at UCL for such agreements has been raised but the establishment of such processes at UCL falls outside the scope of this document.

#### 4.4 Primary vs Secondary Data Analysis

From an ethics standpoint, Twitter data usage in research **falls between primary and secondary data analysis**. It is primary data in the sense that there is no formal curatorial activity standing between the participant providing the data and the researcher who wishes to use it, yet secondary data because the data is gathered not from the participants but from an aggregating intermediary (Twitter). The Menlo Report characterises this as an increased logical or physical distance between the researcher and the humans needing protection in research (Dittrich & Kenneally, 2012).

**IMPORTANT: In general, it is important to recognise that although in many cases researchers treat Twitter data as a secondary data source, the Twitter ‘dataset’ is unlike many datasets used for secondary data analysis in that it is dynamic. The contents change regularly, not just by the addition of new Tweets, but also by deletion and other user-driven changes to the status of available information (Tweets, accounts etc). As discussed in more detail in sections 4.5 and 5.1, Twitter expects those who use its information to respect the changes that users make to the availability of their content. This means that the ethical arguments below around consent and privacy must be applied at every use and regularly during retention, not just at the outset of a research study. This has implications for research design since data retrieval and ongoing management must account for this, and particular care taken around consent for publication (see section 4.9). The dynamism may also make reproducibility more difficult since there can be no guarantee of the same data existing at each analytical attempt.**

It is perhaps helpful to see ‘public’ Twitter data as **private data on public display** on the basis of **ongoing consent under contract**, rather than public data that arises from publication. Seeing the data in this way makes clearer the ethical consent arguments (see section 4.5) because it **separates** the notions of **data ownership** from **data disposition**.

An analogy may be found by considering privately-owned land that is made available from time to time for public use through an agency that acts on behalf of the land-owner. The public does not own the land and authority over access to it remains with the land-owner. However, when so-designated by the land-owner via the agency, the public may use the land for the purposes permitted by the owner and agency. That access may be withdrawn at any time. The land therefore has the temporary disposition of ‘publicly available’ at the discretion of the land-owner, but the ownership and control remains private and when access is withdrawn by the agency on behalf of the land-owner, the land is no longer publicly available and may not be used.

#### 4.5 Consent and Expectations

Within the context established in the foregoing discussion, a researcher must **consider the consent requirements** that apply to relevant stakeholders to ensure that research is undertaken ethically. These are influenced by the mixed primary/secondary nature of the data to be accessed, the general societal context of the research, and the contractual permissions and constraints that govern the researcher (and UCL’s) relationship with Twitter in respect of the particular study concerned.

The stakeholders (aside from UCL) concerned here are primarily **Twitter, Tweeters, and Subjects**.

A strict position on informed consent would require a researcher to seek individual consent from Twitter and from every person whose Twitter data (or who) they wished to study. This might be considered impractical given the potential number of such people involved in the research and therefore alternatives may need to be developed. Each stakeholder is addressed separately in the following subsections.

##### 4.5.1 Twitter

The permission of Twitter to conduct research using its platform is demonstrated in its provision of a Developer API and the terms that go with it: in essence, **Twitter offers its platform for research (but only under certain conditions)**. **Note that scraping Twitter is not permitted as a method to access its data**<sup>2</sup>. It is important to note

---

<sup>2</sup> “You may not do any of the following while accessing or using the Services: ... (iii) access or search or attempt to access or search the Services by any means (automated or otherwise) other than through our currently available, published interfaces that are provided by Twitter (and only pursuant to the applicable terms and conditions), unless you have been specifically allowed to do so in a separate agreement with Twitter (ctd.)

that Harvesters (i.e. researchers) **must inform Twitter of their intentions at the outset** (and if these change) in order that it can approve the proposed work. As such, Twitter imposes an informed consent process that meets its own acceptability criteria. Nothing more than **active compliance with the terms** would therefore be required from an ethics standpoint (unless a researcher intended not to comply).

#### 4.5.2 Tweeters

There are certain circumstances where active informed consent from individuals is required (see publication issues in section 4.9 below) but typically this would be seen as impractical for the kinds of projects researchers often wish to undertake. Impracticality alone is not a strong ethical defence however and more detailed consideration needs to be given as to whether the active informed consent requirement can and should be waived.

**Twitter makes clear to its Tweeters what will happen to their data** (including that it may be used for research) at the point of signing up to use its services. There is plenty of help and support readily available and the controls for public visibility or otherwise are clear. Twitter is **also clear on the ability of Tweeters to change the visibility of their data and that it (Twitter) will then stop** (or will restart as appropriate) **making it available** elsewhere. The Developer API agreement that Harvesters agree to, and the various supporting documents for its use, are also very clear that **users have control over the public disposition** of their data, and that **this should be reflected in its use by others**<sup>3</sup>.

The reasonable expectation of Tweeters may therefore be assumed to be that **when their data has public disposition it may be viewed and used for research**, in accordance with the agreements that Twitter makes with them and makes with those who use data exposed on the platform. The overall societal understanding of Twitter and its use lends weight to this argument. Counter-balancing this is consideration of the mixed views of private/public spaces evidence in researching online communities (Sugiura et al., 2017).

Overall it seems justifiable (absent other risks) to argue that the normal ethical requirement of individual informed consent could be waived since research using a Tweeter's publicly visible data is within their reasonable expectations: **Ethical consent is arguable as implicit under these particular conditions**. It is crucial though to realise that **this argument rests on two things**: that the **data being used has current public disposition**, and that the **use of that data respects the agreements** under which it was made available to the Harvester (researcher). **If either of these conditions is breached** (e.g. a user deletes or protects a Tweet thus removing its current public disposition, or a Harvester/researcher does not respect the full implications of the promises made by Twitter to the Tweeter) then **the ethical argument no longer holds** and an active individual consent would be required<sup>4</sup>.

**Implicit consent** for data use is thus **only available as an ethical defence when considered against the current state of the dynamic Twitter dataset** at the time of analysis (see section 4.4 for a discussion of the dynamic nature of Twitter data, section 5.1 for the way in which Twitter expresses this contractually, and section 4.9 for the implications in publication). **Retention and use of data no longer in the publicly accessible live dataset on Twitter breaches participant autonomy** since it does not respect their current wishes to withdraw data from availability (in

---

...(NOTE: crawling the Services is permissible if done in accordance with the provisions of the robots.txt file, however, scraping the Services without the prior consent of Twitter is expressly prohibited);" (<https://twitter.com/en/tos>)

<sup>3</sup> "One of Twitter's core values is to defend and respect the user's voice. This includes respecting their expectations and intent when they delete or modify the content they choose to share on Twitter... We believe that business consumers that receive Twitter data have a responsibility to honor the expectations and intent of end users." (<https://developer.twitter.com/en/docs/tweets/compliance/overview>)

"...As business consumers of Twitter data, we have a collective responsibility to honor the privacy and actions of end users in order to maintain this environment of trust and respect... The state of a User or Status can change at any time due to one of the actions above, and this impacts how consumers of Twitter data are expected to treat the availability and privacy of all associated content. When these actions happen, a corresponding compliance message is sent that indicates that the state of a Status or User has changed." (<https://developer.twitter.com/en/docs/tweets/compliance/guides/honoring-user-intent>)

<sup>4</sup> Note that this paragraph is concerned solely with the ethical consent argument, irrespective of whether breaching the data use terms and not respecting the public visibility of data would also lead to a potential breach of UCL's contract with Twitter. Such matters are dealt with elsewhere in this document: even if the legal risks were felt to be acceptable, the ethical consent aspect must still be addressed.

the way in which this was offered to them at the time of data surrender, and the researcher has no consent agreement with the individual concerned that would override this).

#### 4.5.3 Subjects

The situation for Subjects is more complex. In many cases, there will **be little to no evidence that a Subject has consented** to their information being made available to Twitter and thereby Consumers and Harvesters. The nature of **disclosure may be implicit** (e.g. 'my brother') but nonetheless poses risks to identification and privacy when combined with other extant data and real-world knowledge of the individuals concerned. A strict application of consent requirements would require each person involved to be identified and asked for consent. However, this would itself be a potentially unjustifiable intrusion of privacy, a breach of Twitter's contractual limits on off-Twitter identification and matching, and an over-burdening of potential participants by the researcher.

Pragmatic arguments may thus need to be adopted: a balanced ethical position might be to apply the same rules to Subject data as to Tweeter data in terms of overall data access and management, with additional risk mitigation applied to any data involving Subjects (e.g. no reporting of Subject information at the individual level under any circumstances).

Implicit in the above is the need for **any retained data to be regularly synchronised** to the state of the online Twitter data set **on the grounds of the consent** arguments. Twitter provides mechanisms to do this (and contractual requirements to do so), discussed further in Section 5.1.

#### 4.6 Privacy and Anonymity

Participant protection with respect to privacy and anonymity is similar to any other research study where personally identifying data is handled. However, it is worth noting that social media data is **almost impossible to anonymise** because internet search engines can resolve the content back to its source (and thus the identity of the participant). It may also contain GDPR special-category data and/or high-sensitivity data from an ethics standpoint. Since the researcher sampling the API stream of Tweets has little control over the content of the received Tweets at the point of reception, the **data should be treated as if high-sensitivity/special-category at the point of ingestion**, and appropriate protocols and data protection management put in place to address this.

#### 4.7 Restrictions on Research

Twitter places **restrictions on the nature of derivation/inference** that may be undertaken using data retrieved from it. In particular the terms state (emphasis author's):

*"You should be careful about using Twitter data to derive or infer potentially sensitive characteristics about Twitter users. **Never derive or infer, or store derived or inferred, information about a Twitter user's: Health (including pregnancy), Negative financial status or condition, Political affiliation or beliefs, Racial or ethnic origin, or beliefs, Sex life or sexual orientation, Trade union membership, Alleged or actual commission of a crime"***

(<https://developer.twitter.com/en/developer-terms/more-on-restricted-use-cases>)

Researchers intending to use Twitter data **for this kind of analysis** should **seek specific legal and ethical advice** on the interpretation of this clause **at the time they are designing their research to ensure that their proposed analysis will not breach it**. Concerns may be raised if, for example, a project intends to ascribe a political affiliation to a particular user (irrespective of whether that user had stated their affiliation directly) for the purposes of making a decision (like whether to follow their account in the research). Prima facie, that would be deriving information about them.

Note that the page goes on to say:

*“Aggregate analysis of Twitter content that does not store any personal data (for example, user IDs, usernames, and other identifiers) is permitted, provided that the analysis also complies with applicable laws and all parts of the Developer Agreement and Policy.”*

<https://developer.twitter.com/en/developer-terms/more-on-restricted-use-cases>

Thus keyword searching in a Tweet stream for mentions of the above information would likely be acceptable because the aggregate analysis (e.g. counting) would not store personal data that relates an individual to the restricted-analysis status.

**Other restrictions** on data use include **off-Twitter matching** (using data from Twitter (and/or elsewhere) to identify or otherwise associate a Twitter user with their identity elsewhere – this requires specific opt-in consent), and **prohibitions on surveillance, profiling, facial recognition, monitoring sensitive events and groups, vetting, credit or insurance risk analysis, individual profiling or psychographic segmentation** and so forth. There are other **restrictions on the use of geographic information** accompanying or within Tweets, **automated messaging**, and **measuring the service** itself. All of these restrictions are fully documented in Twitter’s various agreements to which an API user (Harvester) agrees.

Many ethical risks of sensitive data analysis using Twitter data are thus potentially reduced by the terms of legitimate data access. If a researcher intends to undertake this kind of analysis in breach of the terms, this should be regarded as high-risk research (see section 6.1). It is **important that researchers and reviewers are aware of the breadth of these constraints** to ensure that UCL is not inadvertently placed in breach of contract, and also to give appropriate consideration as to whether the benefits of research outweigh that aspect of ethical risk.

#### 4.8 Sharing Data

Data sharing falls within the restricted use cases. Twitter **permits the sharing of Tweet IDs and User IDs** in a dataset for others to use (for academic research the number is currently unlimited). The recipient must then ‘re-hydrate’ those IDs into Tweets using Twitter’s API (thus ensuring that anything that has been deleted on the platform is not shared because the ID will not resolve to anything). A limited number of ‘hydrated’ Tweets can be shared, but only privately, and **the dataset creator who is sharing their dataset must ensure that the recipient has agreed to the Twitter terms** before doing so (ensuring that the recipient must then also respect the synchronisation and other requirements).

#### 4.9 Publishing Data

Twitter places particular **restrictions on the form in which Tweets may be published**, requiring certain items of data to be retained in the published form. The forced retention of this material may pose a challenge to privacy. In addition, its synchronisation requirements apply to published material, meaning that should a user delete or protect a Tweet that has been quoted in a paper, that paper would need to be modified to remove it.

This is a well-recognised issue in the ethics literature and there is general guidance to **seek individual informed consent from the user whose Tweet is intended to be published** so that even if they withdraw the Tweet from Twitter, the researcher has clear ethical consent to publish it (a more nuanced approach to publishing Tweets can be found in (Williams et al., 2017)). **Where relevant, this should be included in ethics documentation** submitted for review since it will require ethical consent forms and careful participant briefing, given that their data may have been used without their knowledge to undertake the research to this point. Relevant aspects of the GDPR must also be considered in these situations.

#### 4.10 Backlash and UCL Reputation

Since Twitter reflects a segment of opinion on current local and global events, it is likely that researchers working in this area may be publishing results that could be considered controversial. Particular attention should thus be paid to the aspects of ethics concerning the protection of researchers to ensure that where appropriate: any publications arising from the analysis are appropriately considered for their controversial potential, the risk of attracting



unwanted attention to the researchers is managed, that researchers are appropriately briefed and aware of relevant personal risks, and that UCL is aware of this potential, and the potential for press interest. One route to this might be to seek REC or Research Integrity advice and/or approval for any potentially controversial publication in order that the researchers can “lean” on UCL’s formal support and prior knowledge in the event of problems.

The potential for breach of contract is also present in this area and thus UCL’s reputation should be considered if a particular proposal may require this.

## 5 Other Legal Matters

### 5.1 Contractual - Synchronisation

Data synchronisation and management is an important aspect of Twitter data use. The ethical argument for synchronisation of locally-retained data is made in section 4.5, but there is also a straightforward contractual requirement to do so.

In general terms Twitter’s position on compliance is reflected thus:

*“Any developer or company consuming Twitter data via an API holds an obligation to use all reasonable efforts to honor changes to user content. This obligation extends to user events such as deletions, modifications, and changes to sharing options (e.g., content becoming protected or withheld). Please reference the specific language in the [Developer Policy](#) and/or your Twitter Data Agreement to understand how this obligation affects your use of Twitter data.”*

(<https://developer.twitter.com/en/docs/tweets/compliance/overview>)

In particular, the Developer API agreement states (emphasis Twitter’s):

*“Content compliance*

***If you store Twitter Content offline, you must keep it up to date with the current state of that content on Twitter. Specifically, you must delete or modify any content you have if it is deleted or modified on Twitter. This must be done as soon as reasonably possible, or within 24 hours after receiving a request to do so by Twitter or the applicable Twitter account owner, or as otherwise required by your agreement with Twitter or applicable law. This must be done unless otherwise prohibited by law, and only then with the express written permission of Twitter.***

*Modified content can take various forms. This includes (but is not limited to): Content that has been made private or gained protected status, Content that has been suspended from the platform, Content that has had geotags removed from it, Content that has been withheld or removed from Twitter”*

(<https://developer.twitter.com/en/developer-terms/policy>)

In addition to making this requirement, Twitter provides two methods by which this may be achieved. One is an on-demand API for checking the current state of a particular user or tweet (free but rate-limited in terms of requests per time period), the second requires a subscription to the Compliance Firehose, in which Twitter provides real-time updates on content status. The former is recommended by Twitter for display, 1:1 engagement, distribution to a 3<sup>rd</sup> party by file download, and extended temporal storage. The latter is recommended for those consuming and storing large quantities of Twitter data for extended periods of time.

**Researchers intending to accrue Twitter data should ensure that the dataset they intend to collect can be synchronised using one of these methods.** This may either incur financial cost to secure the compliance subscription, or may bound the size of the retained dataset to that which can be synchronised regularly through the free API.

Further guidance may be found in the Compliance area of Twitter’s website which also contains good practice guidance for data users and specific guidance on how to respond to compliance events.

## 5.2 Contractual – Indemnification and Audit Requirements

The Developer contract includes an agreement that the developer (or their organisation where appropriate) of Twitter's API and data **agrees to indemnify Twitter** against claims made against it if the developer's usage is not consistent with the agreement.

It also requires agreement that the user (or organisation) will **allow a compliance audit** at any time by Twitter to inspect all material in the user's possession relating to its use, and that Twitter may request and be supplied with a **written report from a signed representative** listing the current deployment of the materials and content (raising again the issue of oversight flagged earlier).

Neither of these areas necessarily have relevance to a specific project, but are included since ethics review (at whatever level is appropriate to the work e.g. departmental, and/or Head of Department, and/or UCL REC) should ensure that UCL is agreeable to these terms in each case.

## 5.3 Data Protection

Although Twitter appears to make no formal data protection agreement with API users' or their organisations in the terms and conditions of access, UCL would be seen as a Data Controller in respect of personally identifiable data gathered from Twitter. Legal advice also suggests that the constraints of Twitter's terms and conditions are the way in which Twitter's own GDPR duties are fulfilled (e.g. a user's withdrawal of data from Twitter is in effect the user exercising a contractual right (contract being the lawful basis on which Twitter processes their data) to require cessation of processing by Twitter and thus by Twitter's contracted parties (users of its API)). Given the indemnification clauses referred to previously, it seems likely that the potential exists for Twitter to try and recover its costs from UCL in the event of an enforcement action caused by a researcher not complying with the Developer agreements.

Researchers should also be aware that **undertaking analysis that might be considered as profiling (“automated processing of personal data to evaluate certain things about an individual” (ICO, 2020))** may fall under separate areas of the GDPR with **additional consent and privacy requirements** (e.g. “...send individuals a link to our privacy statement when we have obtained their personal data indirectly.” (ICO, 2020)).

# 6 Engagement with UCL Policies

## 6.1 Ethics

It is likely that all uses of Twitter data will engage UCL's Ethics Policy (and Sensitive Research Policy). The data involved comes from humans, is a record of human activity, and (depending on the particular data) may reveal information derived from humans. It therefore falls within the scope of UCL Ethics Policy.

It is **almost impossible to anonymise Twitter data** (because removing direct identifiers does not prevent a search engine resolving the content of a Tweet back to its originator) and therefore is unlikely to meet the criteria for Exemption 2 (UCL Research Ethics Committee, 2018) from formal REC review.

Whether the Exemption 1 (UCL Research Ethics Committee, 2018) criterion is met depends on the interpretation of 'publicly available information' in the context of Twitter data, and whether one regards it as 'information' or a 'dataset'. The view taken in this guidance document is that the data is not 'publicly available' in the sense that the Exemption 1 intends (see section 4.4), and is not 'information' because the data is highly structured for machine processing: it is a dataset. As an exemplar to illustrate the contrast: an Ofsted report is publicly available information not data (it is not structured for machine processing and it has been made available publicly (published) in perpetuity) whereas Twitter data is structured for machine processing if accessed in accordance with the terms, and its public nature is temporary and only in that state by ongoing consent (see section 4.5).

**Most activity involving Twitter data will therefore require REC approval under the current UCL Ethics framework.**

Resolving ethics issues is a process of balancing the potential harms and burdens of research on participants with the potential benefits of the research itself (Whitney, 2016). Investigators should therefore not be discouraged from challenging the positions set out in this document if they believe their research activity warrants setting aside some

of the guidance, particularly if they can support their positions using the ethics literature in this area. This may increase the risk level of the proposed research (particularly if, for example, the proposed activity may place UCL in potential breach of contract through non-compliance with the terms of access and use). It is therefore important that proposed deviations are explained and justified fully in order that the REC can make an informed judgement as to the risk/benefit balance when reviewing the proposed work.

## 6.2 Data Protection

Since personally identifying data is likely to be processed in the course of Twitter-based research (at the point of ingestion even if not afterward), such activity will **usually require registration with the Data Protection Office** at UCL in accordance with UCL's Data Protection Policy. This will require considering **appropriate privacy notices** to be provided to those whose data is to be used. Where it is impractical to do this (e.g. because of the number of potential participants or for other reasons), there are provisions in the GDPR for alternative approaches to be used. However, **before any such determination of impracticality can be or is made, a Data Privacy Impact Assessment must be undertaken** and fully documented to ensure that such an approach is lawful [UCL legal advice].

## 7 Summary

This document aimed to provide sufficient context for researchers and reviewers to respectively undertake and challenge appropriate research design, justification, and activity when using Twitter data. The length of the guidance reflects the complex and intertwined landscape in which this kind of work is carried out. Nonetheless it is important that those undertaking the work and those with responsibility for overseeing it are aware of the complexities and interactions between ethical and legal issues, and therefore that UCL is then in a position to justify the balance between risk and benefit when the work is approved and undertaken.

## 8 Reference Questions for Reviewers

Whilst by no means a comprehensive list, ethics reviewers may wish to consider the following questions as a starting point when reviewing studies involving Twitter:

- Has Twitter been informed of the proposed work through the Developer API signup mechanism and has it agreed?
- Is any access to Twitter proposed that does not use the API for retrieval (e.g. prohibited methods like scraping)?
- Has the research team read and understood all of the terms that apply to their work?
- Is data to be retained locally? If so, how will this be synchronised sufficiently frequently to reflect the live Twitter data?
- Is data to be shared? If so, does the research protocol for doing so comply with the terms of access?
- How is the informed consent of Tweepsters to be gained (or what argument is being used to justify waiving this requirement)?
- Does the informed consent of Subjects need to be considered and if so, how has it been?
- Is publication of individual Tweets anticipated? If so, what forms, process, and information will be used to seek informed consent for publication from Tweepsters (and if need be, Subjects mentioned in the Tweets)?
- Is the intended analysis likely to derive or infer information in the restricted use cases?
- Are any other restricted activities (off-Twitter matching, profiling etc) proposed? If so, how is this justified to balance the potential breach of contract?
- Has a DPIA been undertaken and data protection registration put in place?
- Does the work qualify as profiling under the GDPR?
- Has the safety of the researchers been appropriately considered given the topic and methods?
- Where a proposal being reviewed appears to be close to the edge of what is permitted, has legal advice been sought on contract compliance?
- Is there a reputational risk to UCL through potential breach of contract?

## 9 References

- Association of Internet Researchers. (2019). *Internet Research: Ethical Guidelines 3.0*. <https://aoir.org/reports/ethics3.pdf>
- Dittrich, D., & Kenneally, E. (2012). *The Menlo Report: Ethical principles guiding information and communication technology research*. US Department of Homeland Security: Science and Technology. [https://www.dhs.gov/sites/default/files/publications/CSD-MenloPrinciplesCORE-20120803\\_1.pdf](https://www.dhs.gov/sites/default/files/publications/CSD-MenloPrinciplesCORE-20120803_1.pdf)
- Information Commissioner's Office (2020), *Rights related to automated decision making including profiling*, <https://ico.org.uk/for-organisations/guide-to-data-protection/guide-to-the-general-data-protection-regulation-gdpr/individual-rights/rights-related-to-automated-decision-making-including-profiling/>
- Sugiura, L., Wiles, R., & Pope, C. (2017). Ethical challenges in online research: Public/private perceptions. *Research Ethics*, 13(3–4), 184–199. <https://doi.org/10/gf2rb2>
- Taylor, J., & Pagliari, C. (2018). Mining social media data: How are research sponsors and researchers addressing the ethical challenges? *Research Ethics*, 14(2), 1–39. <https://doi.org/10/gfpbgj>
- Twitter Inc. (2020a). *Twitter Terms of Service*. <https://twitter.com/content/twitter-com/legal/en/tos.html>
- Twitter, Inc. (2020b). *Developer Agreement – Twitter Developers*. <https://developer.twitter.com/en/developer-terms/agreement>
- UCL Research Ethics Committee. (n.d.). *Guidelines for recruitment via social media and for use of social media data*. <https://ethics.grad.ucl.ac.uk/forms/guidelines-on-the-use-of-social-media.pdf>
- UCL Research Ethics Committee. (2018). *Exemptions*. <https://ethics.grad.ucl.ac.uk/exemptions.php>
- Whitney, S. N. (2016). *Balanced ethics review: A guide for institutional review board members*. Springer International Publishing.
- Williams, M. L., Burnap, P., & Sloan, L. (2017). Towards an Ethical Framework for Publishing Twitter Data in Social Research: Taking into Account Users' Views, Online Context and Algorithmic Estimation. *Sociology*, 51(6), 1149–1168. <https://doi.org/10/ggtf28>