# Platform-controlled social media APIs threaten Open Science

*Brittany I. Davidson[1], Darja Wischerath[1]\*, Daniel Racek[2]\*, Douglas A. Parry[3]\*, Emily Godwin[1], Joanne Hinds[1], Dirk van der Linden[4], Jonathan F. Roscoe[5], Laura Ayravainen[1], & Alicia G. Cork[6]*

[1] School of Management, University of Bath, Bath, UK
[2] Department of Statistics, Ludwig Maximilians Universität München, Munich, Germany
[3] Department of Information Science, Stellenbosch University, Stellenbosch, South Africa
[4] Department of Computer and Information Sciences, Northumbria University, Newcastle, UK
[5] BT Applied Research, Ipswich, UK
[1] Department of Psychology, University of Bath, Bath, UK

Corresponding author: Brittany I. Davidson, bid23@bath.ac.uk
\* These authors contributed equally

**STANDFIRST**

Social media data enable insights into human behavior. Researchers can access these data via platform-provided Application Programming Interfaces (APIs), but these come with

39  restrictive usage-terms that mean studies cannot be reproduced or replicated. Platform-
40  owned APIs hinder access, transparency, and scientific knowledge.

41

42  Social media (SM) data hold tremendous value for studying behavioral patterns over time and across
43  contexts at individual, group, and population levels[1,2]. For example, these data can be used to examine
44  where conflict is likely to occur, where to allocate aid in the event of natural disasters, how online
45  polarization or misinformation is impacting voting patterns. SM data are thus relevant to a broad range of
46  disciplines in the social and behavioural sciences.

47

48  Because SM data are constantly changing as users interact and platforms alter the structure of feeds and
49  interactions, it becomes ever more important for researchers to engage with Open Science (OS) practices
50  to ensure that work is <u>reproducible</u> (reusing the same data and methods provides the same results), and
51  <u>replicable</u> (using the same methods on different data produces comparable results). Reproducibility and
52  replicability are essential to ensure knowledge produced about human behavior is robust, valid, and
53  credible.

54

55  Open Science principles discourage academic misconduct (e.g., misreporting of data, problematic methods,
56  improper documentation of results). While this is important for all scientific endeavors, is of importance
57  for SM research due to the power imbalance between academic researchers and SM platforms. Moreover,
58  given the relevance of SM data for a host of important domains (e.g., politics, mental health,
59  misinformation), it is imperative that findings using SM data are trustworthy and can be relied upon to
60  inform policy.

61

62  SM data can be gathered in a variety of ways (Fig. 1), with some facilitated by SM platforms and others
63  falling outside of the methods officially mandated by the SM platform Terms of Use (henceforth: Terms).
64  Because of the inconsistency in Terms between platforms and the changes that platforms make to their
65  Terms, researchers face substantial ambiguity in how they can collect, store, and disseminate SM data.
66  Further, many of these Terms restrict the extent to which SM data can be shared with other researchers,
67  undermining research transparency and hindering the verification of prior results. Compounding this, recent
68  platform changes have removed many widely used data-collection routes previously essential for SM
69  research, inhibiting the replication of prior findings with new data. Here, drawing on our experiences
70  working with SM data, we shed light on impediments to research transparency associated with platform-
71  controlled access to SM data.

72

73                         [insert Figure 1 here]

74

## Data, APIs, and Terms in Flux

76  SM data and the Terms that govern researchers' access to this data are not static, with both users and
77  platforms able to effect changes that alter the data available to be (re)collected from the platform. Users
78  have the capability to remove or edit data, whilst platforms control many of the routes used by researchers
79  to access data. Researchers are therefore at the mercy of any changes that platforms make to their data-
80  access APIs and Terms governing this access.

81

82    Fundamentally, SM data consist of the digital traces that users provide via their engagement and interactions
83    on the platform. Over time, this data shifts and erodes due to (a) the structural changes that platforms make
84    to the interaction features available to users and (b) the ephemeral nature of the action-records as users alter
85    or delete their content and profiles or change their privacy settings. Illustrating this, Pfeffer et al.[3] found
86    that after one year, less than 70% of original tweets were still available, decreasing to ~54% after three
87    years. This can impact some content more than others: political campaigns have extremely high proportions
88    of tweet and user decay (missingness over time)[4], which has implications for reproducing results especially
89    when data sharing is restricted (see below).
90
91    It is likely that data missingness will increase as platforms enact policies that call for the removal of inactive
92    accounts. This will further reduce the extent to which prior findings can be reproduced using the original
93    data collection procedures. For example, X (formerly Twitter), announced (via an Elon Musk tweet) that
94    they will be '*purging accounts that have had no activity at all for several years,*'. Although it was
95    commented that tweets would be archived, no further information was provided. Google similarly
96    announced that it will start deleting Google (and associated YouTube) accounts that have been inactive for
97    two years from December 2023. As these changes are likely to result in the removal of large swathes of SM
98    data, they will have a substantial impact on the reproducibility of findings drawing on older datasets, thus
99    impacting digital archiving and preservation[5].
100
101   Not only does SM data change and disappear but, troublingly, the data-access APIs themselves are
102   constantly changed and updated. These changes are frequently undocumented and poorly communicated[6].
103   Changes to APIs can include the addition of new fields to gather data previously unavailable or, in some
104   cases, the removal of existing fields or changes in functionality. Updates can also include changes in the
105   way metrics are calculated. This means that even if researchers shared their code to query an API, someone
106   re-running it to re-gather the data may find that the code does not reproduce the same results as those
107   generated by the original researcher. For example, the Reddit API once provided the raw number of upvotes
108   and downvotes per post/comment, however this functionality was later removed, with only aggregate scores
109   remaining. While these scores are derived from up- and down-votes, the individual values are now
110   unavailable through the API. This impacts any attempts to reproduce or replicate prior research, as the new
111   'score' metric is not transparently comparable with the 'upvote-downvote count' which, if available, could
112   be useful to understand the popularity/virality of content, user behavioral patterns, or (mis/dis)information
113   spread. Unfortunately, platforms typically do not document these changes. This highlights how crucial it is
114   for APIs to have up-to-date and transparent changelogs with their documentation.
115
116   Alongside changes to the data and the APIs, researchers must also be aware, firstly, of the restrictions that
117   the Terms imply and, secondly, of changes to the Terms by which data can be collected, stored, and
118   processed[7]. Many SM providers state that it is the researcher's responsibility to keep up to date with the
119   Terms (e.g., TikTok: '*However, it remains your sole responsibility to review these Research API Terms*
120   *from time to time to view any such changes*'). Terms are not fixed for a given platform and often differ
121   between platforms. For example, SM platforms tend to adopt different data ownership models (user-owned
122   vs. platform-owned), which can impact the viability of data collection routes. For instance, Reddit deems
123   all user-generated content as user-owned, meaning that data donation would not breach their Terms (in
124   theory), whereas TikTok forbids any data collection outside of their API, which is presently only available

125 in the US and Europe, eliminating any attempts at reproduction by researchers (or reviewers) outside of
126 these regions.
127
128
129

## Raw Data Sharing Restrictions

131 Sharing the original data underlying the findings reported in a study is critical for facilitating reproduction.
132 Unfortunately, alongside the transparency-risks associated with evolving datasets, APIs, and Terms, many
133 platforms restrict the extent to which researchers can share the raw data collected via their APIs[7]. These
134 restrictions undermine the extent to which findings using SM data can be reproduced as researchers cannot
135 rerun analyses on the original data. For example, since 2016, X has restricted the sharing of raw platform
136 data collected via its API, with subsequent versions of the Terms indicating that researchers were only
137 allowed to share Tweet and user IDs (unique identifiers allocated to each tweet/user). In mid-2023, these
138 Terms were revised to allow the sharing of 50,000 raw tweets a day between two researchers with an upper
139 limit of 1.5M tweets. X backtracked again in August 2023, and stated: '*Academic researchers are permitted*
140 *to distribute an unlimited number of Tweet IDs and/or User IDs if they are doing so on behalf of an*
141 *academic institution and for the sole purpose of non-commercial research*'. While the lifting of this
142 restriction, in theory, enables the sharing of SM data for research purposes (e.g., collaboration, verification,
143 reproducibility), in practice it remains restrictive as researchers can only share Tweet/User IDs and not the
144 raw data. To collaborate/verify/reproduce results, IDs need rehydration (recollecting raw data from IDs via
145 the API). This brings challenges with dynamic data and relies on third-parties paying for API access[1].
146 Currently, X's API is both expensive and restrictive regarding data collection, sharing, and thus impeding
147 replication attempts, especially with large datasets.
148
149 These Terms are therefore in direct conflict with reproducibility and replicability because (a) researchers
150 cannot openly share raw datasets and (b) a complete rehydration is not possible due to data deletion from
151 users, potential field changes and updates to the API. This has a disproportionate impact on the extent to
152 which older datasets can be reproduced and highlights the direct impact of constant changes in API Terms
153 on research[5]. See Table 1 for more examples.
154
155 Other Terms essentially restrict data sharing by virtue of the compliance processes that researchers must
156 follow. For example, TikTok has restrictions in place on the use of their data specifically in relation to users
157 who remove or change their content (e.g., account, posts, engagement) in extremely short timeframes. The
158 Terms state: '*You agree to regularly refresh TikTok Research API Data at least every fifteen (15) days, and*
159 *delete data that is not available from the TikTok Research API at the time of each refresh.*' This is
160 problematic, as it can cause research results to become unstable wherein results would likely fluctuate with
161 each data refresh. Further, researchers would need to perform substantial amounts of additional work
162 (recollecting data every 15 days), which would be especially challenging when working with large datasets.
163
164 In addition to restrictions on the sharing of raw data, platform Terms can also impact the reuse of work
165 using their platform data. For instance, TikTok states: '*After you publish any Research outputs, you agree*
166 *that TikTok will have free and unlimited access to and use of your publication and Research outputs.*',
167 noting that researchers must send all research outputs to TikTok (August 2023). This is potentially

168 problematic, as Terms like these could conflict with publisher agreements, alongside employer Terms
169 relating to reuse of employer name and IP.

170

171 **Final Thoughts**
172 Platform-controlled APIs can threaten the reproducibility and replicability of SM research. The perspectives
173 provided in this comment are grounded in our experiences using these APIs in our research and, while we
174 have studied the Terms that bound this conduct, we do not and cannot provide legal advice. This is a
175 complex area with dynamic policy and regulatory events, and specific legal counsel may be required to
176 guide each study. Notably, at the time of writing, various regulatory bodies are considering whether and
177 how to compel large online platforms to provide data access for research purposes[8]. While this has the
178 potential to address some of the challenges for reproducibility and replicability, other challenges inherent
179 to the data will remain, and the implementation of any policy reform will face substantial regulatory,
180 institutional, platform, and infrastructural challenges. We also acknowledge that while we are encouraging
181 transparency and data sharing, there are numerous ethical and privacy challenges that come with sharing
182 SM data[1,5,9]. Making responsible decisions is a complex process that remains an open question within this
183 landscape.

184

185 Whilst we have used a handful of SM platforms as illustrative examples, we have argued that, broadly,
186 these data collection routes and the Terms that govern their use pose substantial restrictions that not only
187 threaten the transparency of our research but, more fundamentally, risk restricting the advancement of our
188 knowledge on human behaviour[10,11]. Specifically, we have highlighted challenges arising due to (a) the
189 evolving nature of platform data, APIs, and Terms, and (b) the restrictions that platforms place on how data
190 can be accessed, stored, processed, and shared. Alongside these elements, over the preceding years, a
191 growing number of platforms have either removed their data-access APIs, restricted the nature and amount
192 of data available through the API, or placed their APIs behind exorbitant paywalls. Despite the challenges,
193 the constant changes to these APIs will continue to place restrictions on attempts to reproduce and replicate
194 prior SM research and, in doing so, hinder scientific progress.

195

196 **References**
197 1. Walker, S., Mercea, D. & Bastos, M. *Information, Communication & Society* **22**, (2019).
198 2. Lazer, D. et al. *Nature* **595**, (2021).
199 3. Pfeffer, J. et al. *International AAAI Conference on Web and Social Media* **14**, (2023)
200 4. Bastos, M. *American Behavioral Scientist* **65**, (2021).
201 5. Acker, A. & Kreisberg, A. *Arch Sci* **20**, (2020).
202 6. John, N. A. & Nissenbaum, A. *The Information Society* **35**, (2019).
203 7. Freelon, D. *Political Communication* **35**, (2018).
204 8. Leerssen, P., Heldt, A. & Kettemann, M. C. *Digital Communication Research* **12**, (2023)
205 9. Sen, I., Flöck, F., Weller, K., Weiß, B. & Wagner, C. *Public Opinion Quarterly* **85**, (2021).
206 10. Bruns, A. *Information, Communication & Society* **22**, (2019).
207 11. De Vreese, C. & Tromble, R. *Political Communication* **40**, (2023)
208 12. Ohme, J. et al. *Communication Methods and Measures*, (2023)
209 13. Krotov, V. & Silva, L. *AMCIS* **17**, (2018).
210 14. Fiesler, C., Beard, N. & Keegan, B. C. *ICWSM* **14**, (2020).

## Authorship Contributions

Conceptualization: BID

Writing Original Draft: BID, DW, DR, DAP

Writing Original Revision: BID, DR, DAP, JH

Writing Review & Editing: BID, DR, DAP, DW, EG, JR, JH, DL, AGC, LA

Writing Review & Editing Revision: BID, JH, DR, EG, DAP, AGC, DL, LA

**Figure 1. Infographic of common routes for SM data access on LHS, RHS includes a description of each route and the types of data one typically obtains from each data access route.**
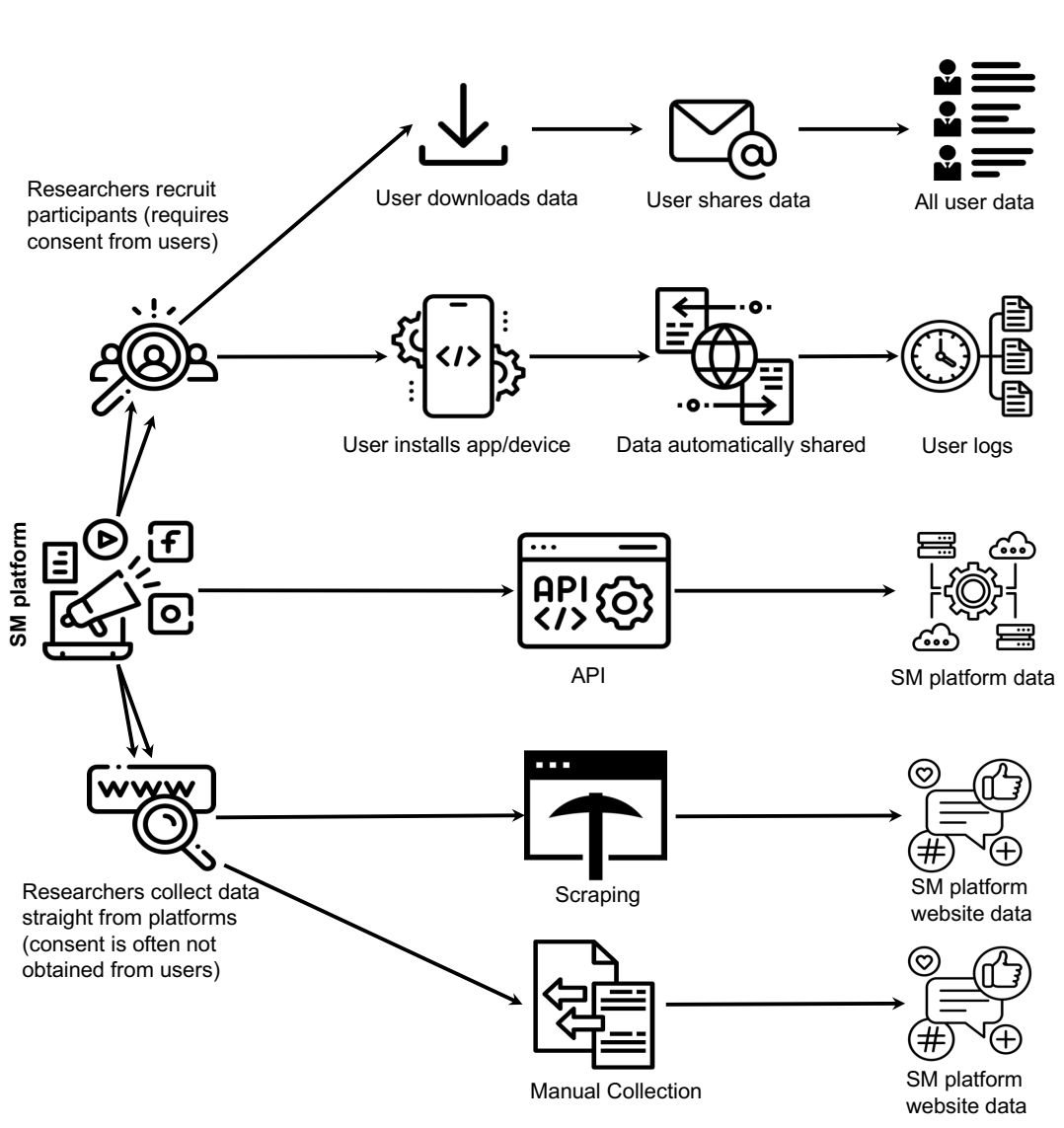
*Table 1. Illustrative Examples of Terms that Impact Open Science*

| Selected Terms (August 2023) | Reproducibility | Replicability |
|---|---|---|
| Reddit<br><br>Noting Reddit states user content is owned by users, [one cannot] 'u*se the Data APIs to encourage [...] violation of third party rights (including using User Content to train a machine learning or AI model without the express permission of rightsholders in the applicable User Content)*'[a] | If data cannot be shared, researchers cannot reproduce original results from articles as they cannot recollect the same dataset.<br><br>If the researchers had used ML/AI and did not share their trained models, then reproducing (retraining the model) violates Terms and thus reproducing the work is not possible. | Datasets may be replicated provided the API provides the same fields and the ways in which metrics are calculated remain the same.<br><br>If the analysis used ML/AI then these analyses cannot be replicated as this violates current Terms. |
| X (formerly Twitter)<br><br>'*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, Religious or philosophical affiliation or beliefs, Sex life or sexual orientation, Trade union membership, Alleged or actual commission of a crime*.' –however at an aggregate level, this is acceptable. | These restrictions mean any prior papers looking at any of these areas at an individual level cannot be reproduced. | These restrictions mean any prior papers looking at any of these areas at an individual level cannot be replicated. |
| TikTok<br><br>*TikTok Research API Data shall not be kept for longer than is necessary for Research approved as part of your application. You agree to provide TikTok with written certification of data deletion upon TikTok's request.* | This means that data cannot be shared so the exact analysis cannot be reproduced.<br><br>It is also vague as to what 'longer than necessary' means (e.g., end of analysis, publication, end of grant?). This is likely at odds with many university or funder data retention policies, too. | NA |

| LinkedIn [b]<br><br>[You agree not to…]<br>'*Sell, rent, lease, disclose, distribute, share (with the exception of making the Content available to Users through the Application), transfer, sublicense, communicate, or otherwise make available, any Content, directly or indirectly, to any third party (e.g. you may not sell access to an aggregated collection of Member profiles, the most relevant Members for a position, or any social activity, such as posts, likes, or shares by Members)*' | This means that data cannot be shared so the exact analysis cannot be reproduced. | NA |
|---|---|---|

[a] This Term is vague and places researchers in a difficult position of not knowing what they can and cannot do, especially when terms such as 'ML' and 'AI' are incredibly broad and work at different scales (e.g., training a task-specific decision tree versus training a general large language model (LLM)). Furthermore, the wording relating to 'training' is ambiguous, and raises the question of whether researchers are free to apply other pre-trained ML models to Reddit data. For example, can researchers use a pre-trained model on Reddit data but not use the same data for training or tuning? This causes other issues, for example, the use of models that are not adapted to a specific SM dataset may negatively impact the models' accuracy and the inferences that can be drawn.

[b] Note: LinkedIn also has strict Terms regarding storing data, where no data are allowed to be stored, unless you have explicit consent from Members.

231

232

**Type of Data Collection** | **Types of Data Collected**

Researchers recruit participants (requires consent from users)

User downloads data → User shares data → All user data

**Data Donation**: A user-centric approach, where users share their data with researchers. Users download their SM data from platforms and share with researchers[12] or install a smartphone app that logs data that is shared with researchers (e.g., UsageLogger). This can be legally complex depending on SM platform Terms.

Typically, all user data (content, posts, direct messages, etc.). Number of users will be limited to the ethics forms, how many sign up, etc.

User installs app/device → Data automatically shared → User logs

**Tracking**: A user-centric approach similar to data donation. Tracking involves capturing logs of data typically from web browsers, apps, or devices via plugins or apps and is automatically transferred to a research server. Contrary to data donations there are no user interactions or requirements for users to actively share their data[12]. This can be legally complex depending on SM platform Terms.

Typically, user logs (interactions with SM apps, other apps, typically no content from SM platform used). Number of users will be limited to the ethics forms, how many sign up, etc.

SM platform

API → SM platform data

**Application Programming Interface** (API): APIs are often the only official way to access data and are tied to SM platform Terms of Service.

A variety of data as permitted from the SM platforms (e.g., post content, likes, images, URL links, language). Number of posts/users will be limited by the SM platform itself.

Scraping → SM platform website data

**Scraper/Spider/Crawler**: A tool created to gather data on websites (or as explicitly programmed) for data analysis[13]. These do not use an API and are often legally complex[14].

Publicly available data can be collected (e.g., post content, images, number of likes). Number of posts/users is large (can be limited for technical reasons, e.g., rate limits)

Researchers collect data straight from platforms (consent is often not obtained from users)

Manual Collection → SM platform website data

**Manual Collection**: An approach where researchers manually collect data (e.g., copying and pasting data from SM feeds). This can be legally complex depending on SM platform Terms.

Publicly available data can be collected (e.g., post content, images, number of likes). Number of posts/users is often limited due to time required for collection.