

Public consultation
on access to data
from online platforms
for research

Public consultation on access to data from online platforms for research Institute for Strategic Dialogue (ISD) August 2022

A.1. Have you ever conducted **research using data** from one or more online platforms? If so, how did you collect it (e.g. using APIs, crowdsourcing, etc.)?

ISD uses three overall types of research methods to collect data from online platforms:

- **Systematic Searching:** using technology to extract large amounts of data and metadata directly from online platforms, e.g. web scraping, APIs;
- **Ethnography:** deep and sustained involvement with a community - researchers take a more human approach; joining, participating in, and observing online spaces;
- **Crowdsourcing and Surveying:** crowdsourcing methods involve users of online platforms voluntarily reporting particular forms of content to researchers, and surveying allows insights into user behaviour, attitudes etc.

Platform	Systematic	Ethnographic	Crowdsourcing
Facebook	X	X	X
Instagram	X	X	X
WhatsApp		X	X
Twitter	X	X	X
YouTube		X	X
TikTok		X	X
Reddit	X	X	
Telegram	X	X	
4chan	X	X	
8kun		X	
Discord		X	
Twitch		X	
Dlive		X	
Gab		X	
Parler		X	
Bitchute		X	
Odysee		X	
VK		X	
Minds		X	
Gettr		X	
Soundcloud		X	

A.2. Did you encounter any **difficulties** in collecting this data? If so, what kind? Please provide examples.

Note: This section is drawn from ISD's recent report, *Researching the Evolving Online Ecosystem: Barriers, Methods & Future Challenges*, available [here](#).

Barrier Type 1: Technological:

Platforms may deliberately use technologies which restrict access to data, or they may also have other technological features which inadvertently create barriers for researchers. The technological features of specific forms of content may also restrict researchers' ability to conduct systematic, large-scale data analysis.

Examples of these technologies and the additional challenges they present include:

- **Encryption:** This is a process by which content is rendered incomprehensible to everyone except specified receivers. Systematic data collection for researchers is impossible without access being granted by the sender or receiver.
- **New formats:** Certain forms of content or data are not (yet) as amenable to systematic search and storage. For example, primarily audio-visual platforms such as YouTube and Spotify present additional challenges because video and audio content cannot easily be searched or analysed in the same manner. AR/VR technologies are also increasingly being developed, and these could be used to spread harmful content or harass other users.¹ It may be possible that new forms of content, perhaps AR/VR-based, will prove much more engaging and effective at radicalising audiences, and/or helping harmful content achieve greater spread or impact. The live and ephemeral nature of AR/VR activity also presents challenges for more systematic data collection.
- **AI-generated content:** As demonstrated by "deep fakes", content generated by artificial intelligence is becoming increasingly believable. This could lead to content proliferating faster than it can be addressed. Additionally, more sophisticated AI could go beyond duplication, allowing content to mutate while retaining its original meaning. The speed at which new content can be developed also makes systematic data collection harder.
- **Decentralisation:** This allows platforms to operate without central governance and

¹ For examples of documented harassment and abuse, see Basu, Tanya, 'The Metaverse has a groping problem already', *MIT Technology Review*, 16 December 2021, <https://www.technologyreview.com/2021/12/16/1042516/the-metaverse-has-a-groping-problem/>; Bokinni, Yinka, 'A barrage of assault, racism and rape jokes: my nightmare trip into the metaverse', *The Guardian*, 25 April 2022, <https://www.theguardian.com/tv-and-radio/2022/apr/25/a-barrage-of-assault-racism-and-jokes-my-nightmare-trip-into-the-metaverse>; Robertson, Derek, 'Crimefighting in the Metaverse', *Politico*, 13 April 2022, <https://www.politico.com/newsletters/digital-future-daily/2022/04/13/who-will-protect-you-in-the-metaverse-00025070>. For examples of initial company research and responses, see Blackwell, Lindsay et al, 'Harassment in Social Virtual Reality: Challenges for Platform Governance', *Proceedings of the ACM on Human-Computer Interaction*, 3(100), November 2019, <https://dl.acm.org/doi/10.1145/3359202>; Gleason, Mike, 'Microsoft, Meta tackle harassment in virtual worlds', *TechTarget*, 17 February 2022, <https://www.techtarget.com/searchunifiedcommunications/news/252513581/Microsoft-Meta-tackle-harassment-in-virtual-worlds>.

can limit the ability of administrators to remove content or ban users (especially those users that have been identified as engaging in patterns of harmful behaviour). Decentralisation may also reduce opportunities for more systemic data access for researchers.

- **Blockchain:** This is a technology via which events (e.g. who posted what content and when) are recorded in an unalterable ledger. This allows the current, true state of a system to be determined by consulting the current state of the ledger without the need for human intermediaries. Blockchain can therefore be used to accomplish decentralisation (e.g. platforms such as Riot). It is also often used to support payment in cryptocurrencies and, increasingly, platforms are using this to allow users to directly monetise content rather than relying on advertising (e.g. Odysee and LBRY). From a research perspective, systematically collecting data from blockchain-based platforms without public APIs remains relatively unexplored territory. Particularly strict use of blockchain might make deletion of content by a centralised authority impossible or nearly impossible (e.g. a situation where an offending user would have to consent to the deletion of their content).²

Barrier Type 2: Ethical and Legal

Accessing data from online spaces, and particularly the collection and processing of that data, can raise ethical issues, such as invasions of privacy or the use of data or content without users' consent. This may also lead to contraventions of ethical research practices, platform terms and conditions, or even the law.

This challenge can be particularly extreme for academic researchers who must often pass strict ethical approval procedures, as well as comply with relevant legal requirements. Law enforcement agencies (and intelligence services in many countries) are also subject to additional legal restrictions on their access to and use of personal data. This is desirable for a multitude of reasons, most notably the human right to privacy and ensuring due process. While the right to privacy is not absolute, exceptions need to be justified under the rule of law. Consequently, privacy restrictions can limit the ability to find harmful content. Some researchers have argued that the growth of privacy legislation across the world (most notably the General Data Protection Regulation (GDPR) in the EU and GDPR-influenced laws in other countries) may give platforms additional incentive not to share data.³

Messaging apps like WhatsApp are a pressing, current example. A huge amount of content is exchanged on WhatsApp, including forms of disinformation, incitements to violence and other harmful material. If a researcher is a member of a WhatsApp group, collecting data is incredibly easy; WhatsApp has a simple functionality to export an entire chat history as a text file. But how did the researcher join said group? Did they gain explicit permission from all the members to use the group's content for research (potentially leading participants to self-censor)? Or are the group members unaware of the researcher in their chat, and therefore might they be non-consenting research participants? Did the researcher potentially gain access to the group via deception?

² Jurdak, Raja, Dorri, Ali and Kanhere, Salil S., 'Protecting the 'right to be forgotten' in the age of blockchain', *The Conversation*, 30 October 2018, <https://theconversation.com/protecting-the-right-to-be-forgotten-in-the-age-of-blockchain-104847>.

³

These problems may be even starker for messaging apps which, as a key part of their market offer, explicitly promise greater privacy and security than more mainstream options like WhatsApp. Platforms that promise greater focus on the privacy of their users have also attracted harmful communities. For example, MeWe was founded in 2012 by privacy advocate Mark Weinstein and has since become popular among conspiracy theorists and far-right extremists.⁴ Kik, an anonymous instant messaging service, has reportedly been used to facilitate child sexual exploitation.⁵ As outlined in the above section on technological barriers, these platforms often use encryption. Additionally, such groups are unlikely to welcome a potentially hostile researcher.

As many of these platforms were created in response to increasing regulations and moderation practices in traditional social media, these “alternative platforms” (or alt-tech⁶) are often presented as bastions of “free speech” and therefore can attract communities and ideologies that have been banned in other spaces for breaching community standards and/or hate, disinformation and harassment policies. This means platform moderation (and by extension terms and conditions and general platform activity) may be explicitly opposed to actions such as content takedowns and banning accounts, or even downgrading harmful content in algorithmic recommendations, newsfeeds or search results.

Barrier Type 3: Fragmentation

Much online content, including harmful content, is *theoretically* accessible online without barriers caused by technological structures or ethical and legal issues; however, one still does need to know where to look. Often relevant content is among vast amounts of material that cannot be searched quickly and systematically, for example, via a platform-wide search function or API. We refer to platforms where theoretically accessible content cannot be searched quickly or systematically as “fragmented”.

As the content is publicly visible, fragmented platforms may be seen as a subcategory of open platforms.⁷ Not all open platforms will be fragmented, however, as some of them do offer the ability for researchers to systematically search content. Fragmented platforms are also distinct from closed platforms. While closed platforms also cannot be searched systematically, they cannot be accessed without additional information or permissions either (e.g. passwords or other types of personal identification).

Modern search tools (most notably Google but also platform-specific technologies like CrowdTangle⁸ or the Twitter API) have only recently increased the ease with which

⁴ Dickson, EJ, ‘Inside MeWe, Where Anti-Vaxxers and Conspiracy Theorists Thrive’, *Rolling Stone*, May 2019, <https://www.rollingstone.com/culture/culture-features/mewe-anti-vaxxers-conspiracy-theorists-822746/>.

⁵ Crawford, Angus, ‘Kik chat app ‘involved in 1,100 child abuse cases’’, *BBC*, 21 September 2018, <https://www.bbc.co.uk/news/uk-45568276>.

⁶ Alt-tech describes social media platforms used by groups and individuals who believe major social media platforms have become inhospitable to them because of their political views. This includes platforms built to advance specific political purposes; libertarian platforms that tolerate a wide range of political positions, including hateful and extremist ones; and platforms which were built for entirely different, non-political purposes like gaming.

⁷ While closed platforms cannot be searched systematically either, they also cannot be accessed without additional information (e.g. passwords or other types of personal identification). See Footnote V for a full definition of open and closed platforms.

⁸ CrowdTangle is a tool for searching public content on Facebook and Instagram. It is owned by Meta and over time, the company has limited the available data. Nonetheless, CrowdTangle still allows a quick keyword query to return an enormous range of material.

researchers could quickly and systematically locate content. This ease, however, can be (and often has been) overstated. A huge amount of the web, potentially over 90%, does not appear in Google Search (this is the so-called “Deep Web”).⁹ Furthermore, important forms of social media and online communication (private and/or encrypted messages, emails and closed groups) have always been off-limits to external researchers. Nonetheless, rapid and systematic searching has become vastly more possible as a technique for the discovery of harmful content and behaviour. But two converging trends may be reducing the power of these methods.

The first trend is that many online platforms, both new and established, are reducing the data that can be accessed through APIs or other tools. This means many key areas of platforms are beyond the scope of the API, forcing researchers to adopt older, more labour-intensive and less systematic research methods, such as manually finding and reading material.

While increasing regulatory and public pressures have their benefits in terms of enhancing privacy and data rights, we may see that platform search tools and APIs become more restrictive by default. Many of the newer platforms identified in our scoping do not have platform-wide search functions, even as part of their APIs. While it is still often possible to use relatively old technologies to access relevant data, this may involve more ad-hoc and labour-intensive methods that need to be designed and maintained for specific purposes, including to produce data in a systematic format. In some cases, using such technologies to access data may also break platforms’ terms of service, thereby presenting additional ethical and legal challenges.

A second potential trend is the broader fragmentation of online hate spaces. The increasing willingness of many large platforms to claim they are “acting against harmful content and behaviours” may be driving these communities to seek (or build) a wide variety of alternative spaces. Technical features may also contribute to this trend. Sites like nandbox allow users to easily create new messenger apps with little technical expertise. These types of service could facilitate the rapid fragmentation of potential spaces for hosting extremist content and communities. There is also a range of large, fragmented platforms like Discord, Spotify or DLive on which harmful content could (and already does) go undetected amid a huge mass of other textual or audio-visual content.

Even if harmful content and behaviours are discovered and addressed on one online platform, they can continue to proliferate across a variety of other platforms as users migrate across the online ecosystem. This is a long-standing issue in addressing harmful online activity, and some measures have been developed to address it, for example, “hashing” to aid the removal of illegal child abuse and terrorist content.¹⁰

Nevertheless, even with tools like this, complete removal of such content from the internet remains extremely challenging. For example, if the precise form of the content

⁹ Technically, the Deep Web consists of online material which is not “indexed” by search engines and so will not appear in a search on Google, Bing, DuckDuckGo, etc. This includes a huge range of material that many people use daily, for example, any material which requires a password to access or is behind a paywall. The Deep Web is not to be confused with the “Dark Web”, which can only be accessed through specific browsers and is often used for illegal activity.

¹⁰ See ‘FAQs / Explainers’, *Global Internet Forum to Counter Terrorism*, <https://gifct.org/explainers/>; ‘Image Hash List’, *Internet Watch Foundation*, <https://www.iwf.org.uk/our-technology/our-services/image-hash-list>.

varies or evolves (rather than being directly replicated), then tracking and removing similar or related content can be even harder. Here, the challenges to identifying relevant content posed by fragmentation may be further exacerbated if edited or similar content is spread at scale across a range of different platforms that cannot be searched quickly and systematically.

A.3. If so, have you ever **abandoned all or part of a research project due to inability to access data** from online platforms? If so, was this the consequence of access being refused? Please provide examples.

Typically, ISD would not abandon an entire project due to data access challenges, but would instead be forced to only use manual/ethnographic methods for the relevant platform. This does restrict the types of research questions we are able to answer. See table above for examples of platforms where only ethnographic methods are possible. This may be because the platform does not provide technical means to access data at scale, or because although technical means are available, their use would contravene the platforms' Terms of Service (and therefore ISD could be exposed to legal risks associated with breaking contract law).

A.4. If not, which **factors** do you think enabled you to successfully collect this data? Did you have the **cooperation of the platform** studied to access this data? If so, how did this materialise? Please provide examples.

ISD's ability to employ multiple different types of research methods (see above) – i.e. our expertise in ethnographic and manual digital research methods allows us to access publicly available data from a wide range of online platforms. We do not typically cooperate directly with platforms when using this type of research approach.

Where we access platform data at scale via 'systemic' technical means, we do not have more privileged access compared to other researchers.

ISD maintains comprehensive Data Protection, Research Ethics, Ethnographic Research, and Researcher Safeguarding & Wellbeing policies to ensure research is conducted in-line with data protection requirements (e.g. GDPR), and in an ethical, legal and safe manner. ISD would be willing to share these policies with ARCOM on request.

➤ *Specific questions for online platforms:*

The policies for making data available for research differ significantly from platform to platform. The following questions aim to gain a better understanding of **their respective policies** and the determinants of these policies: nature of the service, technical or legal specifics, or assessment of specific risks that data sharing could pose.

A.5. Do you have a **policy on sharing your data** with third parties for research purposes?

i) If so:

- **how long** has it existed?

- does it concern one or more specific **categories of recipients** (researchers, NGOs, businesses, etc.)?
 - are there any **criteria for selecting** these recipients? If so, which?
 - what **type(s) of data** does this policy cover?
 - does it include a **control or monitoring component** regarding the use of the data provided?
- ii) If not, for what **reasons** have you not initiated such a policy? These may include legal, regulatory, technical, financial and other risks. Specify your assessment of these risks resulting in the decision to not open up your data.

B. Governance

➤ *Definition of actors:*

Access to data useful to society raises the question of **opening it up to all research stakeholders**. While the academic world appears to be the main beneficiary of more open access, the contribution by **think tanks, journalists and civil society** to the knowledge of issues related to online platforms deserves reflection¹¹. The question of **actors' neutrality**, given the funding they may receive from certain platforms, also arises.

B.1. Should we **define and possibly limit further up the line the types of actors** that can receive access to data: researchers, journalists, NGOs, think tanks, civil society, etc.?

- i) If so, based on what **criteria** (possibly combined with the nature of the research or the objectives pursued)?
- ii) Should they have the **same access possibilities** or should these differ according to the type of actor?

B.2. Should there also be a **minimum level of access for the general public** (or a broader category of recipients than academic researchers), such as the mandatory provision of a certain amount of anonymised data in an open data format?

The starting point for data access **ought to be public**, however, certain semi-private data should be limited to **vetted researchers** to ensure the information is not used by content creators or nefarious actors to "game the system". Semi-private data is often situated in a grey area such as limited-access or closed groups with multiple participants (Discord channels, Facebook 'closed' groups, Slack Channels, large WhatsApp groups, invite-only Telegram channels and groups, and Google documents). In such cases, a channel should be first assumed to be private, and its 'public' nature must be proven. To determine whether user data or platform features should be considered more public or more private, criteria such as size, purpose, accessibility, and nature of relationships should be considered. A **tiered access structure** by which a regulator or institutions accredited by a regulator or other body have increased access to data may be advisable in light of data

¹¹ One of the modes of these contributions is participatory science and research. These are "*forms of scientific knowledge production in which civil society actors participate, alongside researchers, in an active and deliberate way*" at all stages of the research continuum, such as data collection, analysis and interpretation of results (Source: [Participatory research · Inserm, Science for Health](#)).

protection and privacy requirements.

Transparency should empower a broad base of **vetted researchers** whose independent scrutiny is vital to holding platforms accountable. The research community has built records of expertise, research methodologies, knowledge and resources to monitor, identify and analyse new trends and harms of online platforms, making it indispensable for regulators. **Vetted researchers** may include not-for-profit bodies, journalists, civil society organisations or associations that are independent from commercial interests and represent the public interest.

Public interest research activities should reflect the purposes of “scientific research” as indicated in the EU’s GDPR. Although this term is not explicitly defined, the GDPR recognises that it should be “interpreted in a broad manner” and that it includes “studies conducted in the public interest in the area of public health”. Thereby, results of research can provide the basis for the “formulation and implementation of knowledge-based policy” with regard to long-term correlation of a number of social conditions.

➤ *Modes of granting access to data:*

The modes of granting access and the possible **criteria** based on which research projects would be selected should also be taken into account. Indeed, although the legitimacy of the use of data for research purposes is not at issue here, the implementation of this principle raises many issues. For example, the **respective roles of European or national institutions** that might be involved in selecting research projects needs to be considered.

B.3. In your opinion, is a data access model based on formulating access requests to a **trusted third party** relevant?

- i) If so:
 - should this trusted third party be a **European or national** public actor? In this case, what would be its **interactions with other authorities**, for example those responsible for personal data protection?
 - What could be the **modes of organising a targeted and supervised** data access protocol?
 - Should the modes of involvement of the trusted third party be defined according to the **level of risk** associated with the **data**?

The model of a trusted third party for formulating access requests to data should be considered. The trusted third party ought to be an independent intermediary entity, and operate on a combination of both national and European levels. This combined approach would allow the intermediary to share responsibilities and resources across levels. It would allow specific national concerns to be identified while facilitating and fostering research on pan-European issues and trends (cross-border research). Regarding the coordination with other relevant bodies (e.g. data protection regulators, national research bodies and others), an Advisory Board or Group approach could be considered. Data access request applications could then be referred to this Advisory Board if coordination with other relevant bodies is needed. In applications for data access deemed “higher risk”, these could be automatically referred to the Advisory Board. Lastly, the Advisory Board could also play a role in considering appeals or platform objections to data access requests.

- ii) If not:
 - for **what reasons?** These can be diverse: legal, academic, logistical, etc.

- Do you think a model of **direct interaction** between the platform and researchers is preferable? If so, why?

B.4. In the eventuality of a mode of regulation that would involve the intervention of a **trusted third party** in opening up data for research projects:

- i) who would be responsible for **monitoring the implementation** of the application protocol?

National media regulators should be responsible for the monitoring of the implementation of an application protocol, with oversight from the national parliament/legislative body, and/or the European Parliament and Commission.

- ii) what **safeguards** could be put in place to ensure access to data that satisfies the stated need?

Any such data transfers between the independent regulator (administrator) and the executive (government ministries, law enforcement or public prosecutor) must be well justified and, where possible, subject to parliamentary scrutiny to avoid misuse of platform data by the executive branch. The guiding principles of the GDPR of purpose limitation and data minimisation should function as useful guidelines in this regard. As such, data gathered by the regulator for the purpose of identifying whether platforms are compliant with relevant digital regulation should only be used for other purposes by the executive branch if this is well justified.

- iii) how should the **transparency of the decisions** by the access protocol organisers be guaranteed?

From the outset, the criteria against which the requests are evaluated, as well as the submission and decision-making process, must be made publicly available. Applicants should be able to request feedback if their application is denied, stating reasons (e.g. lack of relevance/public interest, lack of specificity/data minimisation, insufficient data protection safeguards etc.). To allow researchers to share best practices and save time when preparing applications, and to ensure legitimate requests for data are more likely to be successful, a public database should provide summaries of previously successful applications. There should also be an opportunity to revise and resubmit applications (within a certain limit), and/or appeal decisions made.

- iv) what position and roles should each of the **stakeholders** have, especially the platforms?

Regulators should be able to request applications on particular research topics/themes/areas of concern. Platforms should have the opportunity to contribute to the process, including raising concerns around the relevance/necessity of requested data, but this should take place in a formal/structured way, and be covered by transparency requirements (i.e. any objections platforms make, and their reasoning and justifications are public).

- v) do you identify any **risks inherent** in this model? Which ones?

Risk aversion and enforcement of this model are key risks.

1. Risk aversion: this model could prevent some research projects from smaller, new or less established or experienced researchers from gaining access to data, which could disproportionately impact researchers from, or researching topics relevant to marginalised communities (e.g. ethnic and religious minorities, LGBTQ+, or economically disadvantaged communities, etc.). Risk aversion could also limit the identification of emerging issues not yet fully understood by the "third party actor" and

other stakeholders (such as platforms).

1. Enforcement – for example, how does the third party intermediary ensure data is not retained unduly, is stored securely, and is not shared beyond approved researchers.

C. Construction of scientific projects

Recent and future transformations of online platforms raise the question of **researchers' ability to identify their data needs** in order to shed light on a social, economic, political or cultural phenomenon. The risk of **information asymmetries** between researchers and platforms is high, and **support for a scientific project by an external committee or regulator** could be a way to facilitate the development of research protocols.

C.1. When preparing their access request(s), how can we foster **researchers' knowledge of the data** from the platforms that they might contact for their studies?

Provide comprehensive code book / list of available (and potentially unavailable) metrics and types of data for each platform.

Provide searchable database of previously conducted research (applications / outputs) to provide examples of types of data available, and which research questions it has been used to answer, which methods have been employed etc.

C.2. Who would define the **scope of the research projects** and their connection to one or more missions of public interest and preside over the identification of the data to which access would be necessary? Should the data concerned be restricted to particular fields of research? If so, which ones? For example, combating information manipulation, hate and online piracy.

Researchers define the scope of their projects (e.g. key research questions, platforms, time period, language, geographic context etc.), and submit an application that includes a justification for why the research is in the public interest, and an explanation of how the requested (types of) data will be used to answer the research questions.

The regulator and/or third party body would then assess the application and justifications and explanations to assess whether they are reasonable, and in-line with data minimisation principles to ensure only necessary data is shared.

We would not recommend specifically restricting pre-determined types of data/metrics for different topics/fields of research, but instead have a consistent set of rules/principles that define how decisions will be made around what types of data can be requested/shared.

C.3. How would requests for access be **formulated** by interested researchers? For example, through calls for project tenders on predefined and/or ad hoc themes, after identifying relevant study topics?

We would recommend a standardised and consistent application process to reduce the burden on researchers, and a combination of calls for research on themes of particular

interest/concern for the regulator, and an open/rolling application process that allows researchers to submit suggested topics/themes (as they may be better placed to identify emerging issues etc). In either instance however, the research would have to relate to topics/themes that would be within the scope of regulation.

➤ *Assessment of access requests and award criteria:*

The questions in this section are based on the assumption that research projects requiring access to data from online platforms have been defined through formalised requests (e.g. to a trusted third party). The question of assessing their **scientific quality** arises. How **innovative projects are and their level of contribution to the scientific literature** are aspects that could influence the modes of data opening. Examining requests in light of these issues would require **the involvement of independent expert committees** to assess requests, based on a clear protocol and transparent criteria. These could take different forms depending on the discipline, **while remaining within a previously defined theoretical authorisation framework.**

C.4. Do you think it is appropriate for a **committee to assess and monitor** access requests?

See response to Question B.3.i above regarding "Advisory Board / Group"

i) If so, how should this **assessment committee** be composed (e.g. an international scientific committee)? Should one or more **regulators** have a position and role on it and, if so, which?

This assessment committee should be made up of a combination of relevant regulatory bodies (e.g. data protection regulators), research expertise (both academic and civic society researchers), as well as potentially an industry representative.

ii) If not, why not? What mechanisms would you consider more able to meet researchers' access requests?

C.5. To what extent would the more or less **binding** nature of the **obligations for platforms to open up their data** require their presence on the assessment committees? Should platforms also have a **right of return** in relation to researchers' requests or even a **right of refusal**?

See response above – platforms should have the opportunity to object to requests, which would be considered by the third-party (and appeals via the Advisory Board/Group), but not an outright ability to block research unilaterally. Justifications would have to be provided and independently assessed, and should be subject to strict transparency requirements.

C.6. What would be the **criteria for granting access**? For example, is it necessary to have a research project involving interdisciplinary teams, possibly from structures located in at least two EU countries, in order to be selected?

We would recommend a more principles-based criteria:

- Public interest

- Relevant to 'harms' covered by regulation
- Proportionality of data requests (i.e. data minimisation)
- Specificity of data requests (i.e. directly necessary to answer research questions)
- Data protections / cybersecurity policies and track record of researchers/organisation (e.g. any examples of misuse of data etc.)
- Suitability of researchers to conduct project (necessary expertise, experience etc.)

We would not recommend more arbitrary requirements (e.g. interdisciplinary teams, multiple organisations from different Member States etc.) as this may hinder certain types of research (e.g. topics that are very country specific, or narrow/specific projects that do not require an interdisciplinary team).

C.7. Should a **time dimension** be included in the assessment of calls for project tenders, so that only those with a short or long duration are selected?

Similar to above, we would encourage flexibility and would not recommend more arbitrary restrictions – the time taken to conduct research should be assessed in terms of whether it is appropriate for the project in question.

It would likely be beneficial to have a mix of research – some larger scale/more long-term studies, combined with shorter more agile research that is able to quickly investigate emerging risks/topics.

➤ *Scientific production and showcasing:*

C.8. Should the work resulting from the analysis of these data be **externally certified**? If so, what form might this take?

We would recommend that this is optional – i.e. data access should not only be open to academic researchers. The application process should ensure that approved researchers are able to produce high quality, credible research. Outputs should be publicly available, enabling scrutiny of non-peer reviewed studies. If sensitive data has been used, other researchers could request access to the data to scrutinise an existing research output and/or attempt to replicate findings.

While additional reviews of research can play an important role in quality assurance and ensuring findings are replicable, there can also be disadvantages of this approach, particularly in the context of the rapidly evolving online ecosystem, and the range of harms that can result. Certification could slow the research process, and therefore be less well suited to responding to emerging risks and 'harms' in this fast-moving online ecosystem. Additionally, peer review or certification processes themselves are not infallible, for example through delays in finding appropriate reviewers, missing potential methodological or ethical issues, or can lack independence from the original research team. These issues can especially occur in particularly specialist or emerging areas of research.

C.9. What precautions should be taken concerning the **publication of the studies carried out**, for example concerning the sensitivity of the data that would have been used? How can the implementation of these precautionary measures be reconciled with the fundamental principle of **researcher independence**?

All research should be published publicly (including findings derived from non-sensitive/aggregated data), but any sensitive and/or personal data could be held back

and only shared in full on submission of a specific application process (i.e. less onerous as a new application, as data would already be determined, but would still have same level of safeguards as recommended above).

D. Data protection and technical considerations

➤ *Identification of relevant data and construction of materials:*

D.1. Given that research projects relying on platform data may favour an angle of analysis that would require a specific database format (variables, granularity, etc.):

- i) how can we enable the **creation of specific or unique databases** that would be built to meet specific needs?
- ii) to what extent would certain research projects enable the **construction of innovative indicators or measures** that could contribute to collective knowledge on the issues studied?

D.2. Can and should data access be **jointly constructed** on an equal footing between governance actors, researchers and platforms, based on the model of INSEE's CASD¹²?

D.3. How can the **data access framework** – governance, types of data identified in relation to missions, etc. – **be made long-term** to ensure it remains suited to the regular innovations of and changes to platforms?

A principled based-approach would help ensure consistency in application process over time. This would require:

- Ensuring governance structures provide the right balance of power or influence between different stakeholders, and the right incentives for good faith cooperation. In addition, all decision-making processes should aim for transparency by default.

- Regular mandatory updates to code books / lists of available metrics should be required as platforms evolve. For example, new platform functions, internal changes to how data is collected, how certain metrics are defined, or the establishment of new platforms or platforms which fall into the scope of regulation.

- Establishing and maintaining ongoing multi-stakeholder working groups to identify and address emerging technical, data protection, or cyber-security challenges, in effect establishing a basic early warning system to identify emerging and future challenges.

Modes of access and storage:

D.4. **What modes of access** should be preferred for online platforms' data? What are their different advantages and disadvantages? Should these differ according to the data collected? If so, why?

Those functions of the platform that are public, and have a reasonable user expectation of publicity, should be computationally transparent. Computational transparency requires platforms of scale to provide APIs, or application program interfaces, that facilitates various functionalities for public interest research purposes. A functional API should enable real time scrutiny of a system's inputs and outputs to verify function and impact. As an indicative list of requirements, APIs should facilitate:

- Searchable public identifiers including the keywords of textual content, the URL of a link or piece of media, or the handle used by an account.
- Search functionality by the text of the content itself, by the content author or by date range;
- Metrics that indicate the top-performing content by geographic region and language;
- All images, videos, and other content in a standard machine-readable format;
- Live and historical data for real-time, and longitudinal trend analysis;
- Fit-for-purpose download options;
- Reliable information about the reach of a piece of content (i.e., how many times was this content viewed or presented to users in their timelines).

APIs should remain consistent, so that long-term studies are not negatively impacted by changes or limitations in access.

A significant gap exists in the public's understanding of platforms' processes, especially with regard to amplification systems, content moderation and redress mechanisms. User data is used for micro-targeting of advertising, including location, IP, browsing data, and information collected from devices. There are risks associated with the use of automated marketing, where individuals or groups are tracked, measured and targeted by machines, potentially using machine-generated content, for advertising purposes.

As an indicative list of requirements, a **searchable ad library** should include:

- Paid political ads and issue-based ads, without limiting access on the basis of pre-selected topics or keywords;
- Targeting criteria used by advertisers to design their ad campaign as well as information about the audience that the ad actually reached;
- Number of impressions that an ad received within specific geographic and demographic criteria (e.g., within a political district, in a certain age range), broken down by paid vs. organic reach; number of engagements that an ad received;
- How much an advertiser paid for ad placement;
- Any additional platform functionality, whether the ad was a/b tested and the

different versions of the ad;

- If the ad used a lookalike audience and if so, the features (age, gender, geographic, etc.) used to create that audience;
- If the ad was directed at platform defined user segments or interests; or if the ad was targeted based on a user list the advertiser already possessed;
- Fit-for-purpose availability of up-to-date and historical data.

D.5. How can a **secure access mode** be guaranteed, particularly when the data is **not anonymised** and/or concerns **business secrecy** issues?

D.6. How should these data be **stored** to ensure the **protection of personal data** and, where appropriate, **business confidentiality**?

D.7. What would be the role and scope of intervention of **data protection authorities** (national and the European Data Protection Centre) in assessing the risks associated with access to this data?

See response above on B.3.i. regarding an Advisory Board, alongside existing responsibilities/powers (i.e. ability to enforce data protection requirements on researchers).

D.8. Should research projects receive **support** from the structure granting access, e.g. of a technical, financial or other nature?

Yes, potentially in the form of a technical 'help-desk' type function, funded via the regulator and/or platforms (and platforms should also have a similar function and/or designated contact).

This could also serve as hub for digital research – e.g. through the provision of a database of research conducted, resources on methodologies and tools, data protection, and legal/ethical concerns. It could also offer a 'matchmaking' function to connect researchers with interests in similar areas, capacity building for smaller or newer research organisations to share experience and expertise across different research fields, sectors and geographic contexts.

ISD would not recommend the third-party structure making decisions about data access also provides funding, as this could lead to a lack of clarity or misaligned incentives. The regulator (and potentially platforms) could provide funding for research in key areas of interest for them, but researchers should also be able to apply for data access for projects funded by other entities (e.g. foundations, philanthropy, community groups representing minorities and/or marginalised communities etc.).

E. Feasibility of access and incentives

➤ *Support for researchers:*

E.1. How can **researchers be supported** in building their research projects and complying with the GDPR and the standards set by the mechanism?

See response above re. role that third-party could play as a hub for digital research, sharing best practices and capacity building.

E.2. What mechanisms could be used to mitigate the **funding and technical**

capability differences between academic institutions, which could lead to a small number of research teams capturing projects?

See above – partnerships should be encouraged between organisations, especially partnering more experienced and/or well-funded research organisations with smaller and/or less established organisations working on similar topics, themes or methodologies. In addition, a central hub of resources should be provided. Finally, it should be a stated objective of any funding mechanisms (via governments, regulators, platforms etc.) to expand the research field and encourage new entrants.

➤ *Platform incentives:*

E.3. How can effective and balanced **incentives** be put in place to ensure that platforms are part of the open data dynamic? How can these actors be integrated into the system in a coherent way and how can best practices be promoted?

See above – Regulation should require platforms to provide data access when suitable applications are approved. Platforms should be involved in a largely advisory capacity, and also contribute to developing best practices and capacity building activities (this can be assessed when required under regulation, or other commitments to cooperate with regulators and researchers, e.g. in the EU Disinformation Code of Practice). Platforms should have the opportunity to object to data access requests under pre-determined circumstances / for specific types of reasons (e.g. privacy concerns, misuse/abuse of data concerns etc.).

E.4. Would the involvement of an **external audit committee** be relevant:

i) further up the line, in the assessment of approval decisions based on the CESP model in the field of statistical surveys in France, for example?

Individual applications may not be relevant (except for appeals), however an external audit committee could play a role in a more macro level assessment to ensure the third party body is acting consistently and not discriminating against certain types of applications (e.g. on a given topic, researchers from particular backgrounds/countries etc.).

ii) further down the line, in the review of the platforms' responses to access requests?

An external audit committee could be relevant in assessing whether the regulator is effectively and consistently assessing company compliance with data access requirements.

E.5. What procedural safeguards could be put in place in relation to **business secrecy** issues?