



## Feedback on the Draft Delegated Regulation on Data Access Provided for in the Digital Services Act

Diyi Liu, Manuel Tonneau, Juliette Zaccour

Oxford Internet Institute, University of Oxford

December 2024

We are doctoral researchers from the Oxford Internet Institute, University of Oxford, specialising in the technical, ethical, and socio-legal dimensions of platform governance. Our research focuses on empirically investigating systemic risks posed by digital platforms, with particular emphasis on auditing human-in-the-loop content moderation systems. Our research expertise includes both the effectiveness of platforms' approaches to illegal and harmful content, and the broader implications of data access mechanisms for public interest research. We also conduct research on privacy-preserving data sharing solutions and their impact on research integrity and reliability.

As individual researchers actively engaged in platform and data governance studies, we offer this feedback on the proposed mechanisms and practical implications of the [Draft Delegated Regulation on Data Access under the Digital Services Act](#) (hereinafter the Draft Act).

The Draft Act emerges at a critical inflection point for platform research. Recent developments have created significant obstacles to independent investigation of platform governance and societal impact. These challenges include the widespread shutdown of platform APIs and research tools<sup>1</sup>, severely limiting researchers' ability to study systemic risks arising from digital platforms<sup>2</sup>. More concerning still, independent researchers face increasing legal pressures when conducting studies

---

<sup>1</sup> Perriam, Jessamy, Andreas Birnbak, and Andy Freeman. "Digital methods in a post-API environment." *International Journal of Social Research Methodology* 23, no. 3 (2020): 277-290.

<sup>2</sup> Davidson, Brittany I., Darja Wischerath, Daniel Racek, Douglas A. Parry, Emily Godwin, Joanne Hinds, Dirk van der Linden, Jonathan F. Roscoe, Laura Ayravainen, and Alicia G. Cork. "Platform-controlled social media APIs threaten open science." *Nature Human Behaviour* 7, no. 12 (2023): 2054-2057; Loveluck, L. "How new Twitter rules could hinder war crimes research and rescue efforts." *The Washington Post*. June 20, 2023, accessed Dec 01, 2024, <https://www.washingtonpost.com/technology/2023/06/20/twitter-policy-elon-musk-api/>; "Letter: Twitter's New API Plans Will Devastate Public Interest Research." *Coalition for Independent Technology Research*. Apr 3, 2023, accessed Dec 01, 2024, <https://independenttechresearch.org/letter-twitters-new-api-plans-will-devastate-public-interest-research/>.

related to harmful content and(or) algorithms on certain platforms<sup>3</sup>. These developments underscore the urgent need for robust regulatory frameworks ensuring protected platform data access for research. In this context, Article 40 of the DSA represents not only a crucial mechanism for ensuring platform accountability within the European Union, but also establishes a potential global standard for researcher access to platform data.

In this feedback letter, we address several aspects of the Draft Act requiring clarification before finalisation. Our analysis focuses on three key areas: the appropriateness of data access modalities, the scope and context of accessible data, and the underlying enforcement and coordination mechanisms. Drawing on our empirical research experience, we offer recommendations for strengthening these provisions to ensure meaningful transparency and effective implementation.

### **Data Access Modalities and Conditions**

The determination of data access modalities outlined in Article 9 of the Draft Act requires substantial clarification. While Recital 6 requires data providers to provide data inventory overviews and “*where possible, indicate suggested modalities for accessing them,*” Recital 16 requires applicant researchers to propose preferred access modalities in data access application. This creates potential confusion during the decision-making process for determining appropriate access modalities and the rationale behind such decisions.

Moreover, the current framework may inadvertently encourage researchers to accept whatever access modalities are available in the inventory, which could include Terms of Service from data providers or third-party providers that impose additional restrictions. While Article 15(3) prohibits data providers from imposing “*archiving, storage, refresh and deletion requirements that hinder the research referred to in the reasoned request in any way,*” stronger safeguards are needed for research independence. Data providers often include requirements in their Terms of Service for researchers to submit research outputs before publication, ostensibly to identify potential personal data disclosure. While data protection is crucial, the Draft Act should explicitly protect researchers’ academic freedom to publish findings without prior approval from the data providers. Specifically, data protection measures are to be established through the reasoned request; as Recital 16 suggests, it is the Digital Services Coordinator (DSC)’s role to verify that the access fulfils both data protection and research integrity requirements. Data providers have an opportunity to raise data protection concerns through the dispute settlement procedure (Article 13) and should therefore not intervene at the later stages of the research, so as to preserve researchers’ independence.

---

<sup>3</sup> Nick Robins-Early. “Judge dismisses ‘vapid’ Elon Musk lawsuit against group that catalogued racist content on X.” *The Guardian*. Mar 25, 2024, accessed Dec 01, 2024, <https://www.theguardian.com/technology/2024/mar/25/elon-musk-hate-speech-lawsuit>; Kayser-Bril, N. “AlgorithmWatch forced to shut down Instagram monitoring project after threats from Facebook.” *AlgorithmWatch*. Aug 13, 2021, accessed Dec 01, 2024, <https://algorithmwatch.org/en/instagram-research-shut-down-by-facebook/>.

**We recommend that the Draft Act establishes baseline requirements for appropriate data access modalities, developed through expert consultation during the finalisation of the Draft Act, to guide DSCs in evaluating proposed access arrangements. Additionally, the Draft Act should explicitly address any types of data sharing agreements that impose post-access restrictions on research independence, particularly regarding the publication of research findings. Data protection concerns should be addressed through the established procedures in Article 13 rather than through other provider-imposed restrictions.**

### **Scope of Data Access and Meaningful Transparency**

The Draft Act provides a comprehensive framework in Recital 12 for data access through reasoned requests, particularly *“data related to content moderation and governance, such as data on algorithmic or other content moderation systems and processes, archives or repositories documenting moderated content, including accounts as well as data related to prices, quantities and characteristics of goods or services provided by the data provider.”* While the Draft Act rightfully emphasizes necessity and proportionality in data access applications, its current implementation framework presents significant challenges.

The requirement for researchers to prove that their research purposes *“cannot be achieved by other existing means”* creates an undue burden of proof, particularly regarding data that are deemed available *“through other sources”* (Recital 12). In practice, this could allow platforms to deflect legitimate research requests by pointing to existing transparency disclosures, even when such data lacks the granularity or context needed for meaningful research on systemic risks.

Our ongoing analysis of the DSA transparency reports from Very Large Online Platforms (VLOPs), particularly regarding content moderation workforce and effectiveness, demonstrates the limitations in both scope and context of data available through current transparency practices.

Regarding the moderation workforce, the lack of appropriate normalisation metrics severely hinders meaningful analysis. While the number of moderators is reported by EU language, there is no accompanying information on the volume of content generated in each language, making it impossible to evaluate whether the allocated resources are sufficient for each EU language. Platforms do provide data on monthly average user numbers by EU country, but this data is an imperfect proxy for normalisation due to two key issues. First, languages often span multiple countries, and individual countries may have speakers of various languages, making it hard to map country-level user numbers to language-level moderator numbers. Second, moderators primarily review content rather than users, meaning that user numbers provide an imperfect insight into the actual workload for content moderation. Without robust language-specific content volume metrics, assessing the adequacy of language-level moderator numbers therefore remains highly problematic, despite platforms having access to the necessary data for more robust assessments of resource allocation.

Beyond the question of missing normalisation data, we argue that more justifications are needed from platforms to motivate potential cross-lingual disparities in their moderator workforce. Our ongoing work reveals significant disparities in resources invested by platforms in moderator workforce across EU languages when normalising by academic estimations of the amount of content that moderators need to review. While there may be legitimate reasons for such disparities – such as different performance in AI detection models used for moderation across languages, different prevalence rates of harmful content, or financial incentives and operational constraints across markets – platforms generally provide no justification for these resource allocations. Understanding these rationales is crucial for contextualising whether moderation resources are appropriately distributed to address systemic risks across the Union.

We encounter similar issues when studying moderation enforcement: while platforms report absolute numbers of moderated harmful content for categories such as hate speech and child sexual abuse materials (CSAM), these figures lack baseline data on total volume of harmful content that would help understand the share of all harmful content that is actually moderated. Although platforms conduct internal research on the prevalence and effectiveness of content moderation, such as Facebook’s studies on hate speech prevalence based on representative sampling<sup>4</sup>, there are no mechanisms in the current Draft Act to ensure independent researchers can access such analytical data for verification and further study. In the absence of data sharing from platforms on the subject, getting access to such information can prove prohibitively costly for researchers as illustrated by our recent work<sup>5</sup>: to examine how the prevalence and composition of hate speech varies across eight languages and four English-speaking countries, we annotated a representative sample of tweets posted within one day on Twitter (now X)<sup>6</sup>. The annotation costs were approximately 30,000EUR, which shows the substantial barriers researchers face when platforms withhold contextual data necessary for analysis of systematic risks, particularly the significant resources required to replicate analysis that platforms conduct internally.

**To address these challenges and ensure meaningful transparency while respecting the principles of necessity and proportionality outlined in the Draft Act, we recommend that the Draft Act: (a) clarify the burden of proof for data access requests by establishing clear criteria for what constitutes adequate alternative data sources, (b) require DSCs of establishment to maintain an updated registry of available data sources, and (c) explicitly address the scope of accessible data to include necessary contextual information for raw**

---

<sup>4</sup> Arcadiy Kantor. “Measuring Our Progress Combating Hate Speech.” *Meta Newsroom*, Nov 19, 2020, accessed Dec 05, 2024, <https://about.fb.com/news/2020/11/measuring-progress-combating-hate-speech/>.

<sup>5</sup> Tonneau, Manuel, Diyi Liu, Niyati Malhotra, Scott A. Hale, Samuel P. Fraiberger, Victor Orozco-Olvera, and Paul Röttger. “HateDay: Insights from a Global Hate Speech Dataset Representative of a Day on Twitter.” *arXiv preprint arXiv:2411.15462* (2024).

<sup>6</sup> Pfeffer, Juergen, Daniel Matter, Kokil Jaidka, Onur Varol, Afra Mashhadi, Jana Lasser, Dennis Assenmacher et al. “Just another day on Twitter: a complete 24 hours of Twitter data.” *Proceedings of the international AAAI conference on web and social media* no.17 (2023): 1073-1081.

numbers. Continuing with the example above, necessary contextual data could include moderator count per language, both the total volume of harmful content (prevalence) and the volume that is actually detected and moderated (enforcement), which could strengthen our assessment of moderation effectiveness.

### **Balancing Data Protection with Research Integrity**

While Recital 16 states that the DSC “*should assess whether the access modalities proposed by the applicant researchers in the data access application are appropriate to fulfil the requirements of data security, data confidentiality and protection of personal data and, at the same time, enable the attainment of the research objectives of the research project*”, in making decisions on the appropriate access modalities, Article 9(2) limits the DSC’s considerations to “*the sensitivity of the data requested, the rights and interests of the data provider, including the protection of confidential information, in particular trade secrets, and the security of its service*”. Without explicitly stating the need to balance these considerations with research needs and the public interest, there is a risk of undermining the social value of the proposed research.

**We recommend amending Article 9(2) to clarify the DSCs should balance the risks and benefits of data access in deciding on appropriate access modalities, by considering not only data sensitivity and provider interests but also the value of the proposed research and its alignment with the public interest.**

Previous research has documented instances where platform-provided data or research tools proved inconsistent or misleading<sup>7</sup>. To address these concerns, the Draft Act should require data completeness and quality assurances from the data providers. Without these additional safeguards and specific requirements, the risk of “transparency theatre” – where platforms provide data that satisfy technical requirements without enabling meaningful oversight – remains significant<sup>8</sup>. The Draft Act presents an opportunity to establish not only legal requirements for data access, but also comprehensive standards for meaningful access that enable effective research and accountability.

**We therefore recommend adding provisions that require data providers to attest to the completeness, accuracy, and representativeness of the shared data. This would ensure that access modalities go beyond formal compliance and effectively enable meaningful research.**

---

<sup>7</sup> Pearson, George DH, Nathan A. Silver, Jessica Y. Robinson, Mona Azadi, Barbara A. Schillo, and Jennifer M. Kreslake. “Beyond the margin of error: a systematic and replicable audit of the TikTok research API.” *Information, Communication & Society* (2024): 1-19; Timberg, C. “Facebook made big mistake in data it provided to researchers, undermining academic work.” *The Washington Post*. Sep 10, 2021, accessed Dec 01, 2024, <https://www.washingtonpost.com/technology/2021/09/10/facebook-error-data-social-scientists/>.

<sup>8</sup> Suzor, Nicolas P., Sarah Myers West, Andrew Quodling, and Jillian York. “What do we mean when we talk about transparency? Toward meaningful transparency in commercial content moderation.” *International Journal of Communication* 13 (2019): 18.

Indeed, beyond the modalities of access, the granularity and quality of the data itself are determinant in establishing meaningful access. In particular, the choice of anonymisation techniques is critical in ensuring research integrity while safeguarding privacy<sup>9</sup>. Unilateral decisions by the data providers related to data anonymisation, such as using privacy-enhancing technologies like differential privacy and synthetic data generation, hold a risk of (voluntarily or involuntarily) significantly altering the data and undermining the reliability and validity of research findings<sup>10</sup>.

**Therefore, we recommend (a) amending Article 9 to clarify the extent to which each party is responsible for determining appropriate anonymization techniques (where relevant), and (b) amending Article 15 to require that data providers adequately document the use of data processing and anonymization techniques, and provide validity guarantees to researchers. By extension, we recommend extending Article 15(3) to specify that data providers should not impose data processing and aggregation that was not set out by the reasoned request and that hinder the research in any way.**

## Coordination Mechanisms

The Draft Act's emphasis on harmonisation and consistency across DSCs of establishment raises concerns regarding implementation and oversight. While assigning substantial responsibility to DSCs of establishment and other national regulators, the Draft Act lacks robust mechanisms for ensuring consistent application of standards across jurisdictions. A primary concern lies in the handling of delays and determination of access modalities. Recital 7 requires DSCs of establishment to notify principal researchers of delays in processing reasoned requests, particularly when "*data access applications imply international data transfers*" or where risks to "*security of the Union*" are detected. Furthermore, Article 9(2) grants DSCs considerable discretion in determining the appropriateness of access modalities, requiring consideration of "*the sensitivity of the data requested, the rights and interests of the data provider, including the protection of confidential information, in particular trade secrets, and the security of its service.*" Without specific guidance, this could lead to inconsistent interpretations and applications across jurisdictions, service providers, and individual requests.

**We recommend developing detailed criteria and guidelines for applying these exemptions to ensure the implementation of these provisions across DSCs.**

## Dispute Resolution

---

<sup>9</sup> Gadotti, Andrea, Luc Rocher, Florimond Houssiau, Ana-Maria Crețu, and Yves-Alexandre de Montjoye. "Anonymization: The imperfect science of using data while preserving privacy." *Science Advances* 10, no. 29 (2024): eadn7053.

<sup>10</sup> Stadler, Theresa, and Carmela Troncoso. "Why the search for a privacy-preserving data sharing mechanism is failing." *Nature Computational Science* 2, no. 4 (2022): 208-210; Hauer, Mathew E., and Alexis R. Santos-Lozada. "Differential privacy in the 2020 census will distort COVID-19 rates." *Socius* 7 (2021); Offenhuber, Dietmar. "Shapes and frictions of synthetic data." *Big Data & Society* 11, no. 2 (2024).

Article 13 of the Draft Act outlines provisions for dispute settlement and mediation following amendment requests, which present several procedural concerns. The Draft Act leaves critical aspects of the mediation process undefined, which raises questions about its effectiveness and fairness. First, data providers are required to propose mediators when initiating the mediation (Article 13(3)), and “*shall be solely responsible for covering the costs of the mediation*” (Article 13(4)). This raises concerns about whether mediators can be impartial and independent if they are both appointed and funded by data providers while being expected to balance provider and researcher interests in the mediation.

The Draft Act doesn’t specify what constitutes “*undue delay*” in the mediation process. More critically, Article 13(8) states that if no agreement is reached by the time limit set by the DSCs, the mediator shall declare the mediation closed, but remains unclear on the subsequent course of actions. This creates uncertainty about whether data access should proceed as specified in the original reasoned requests or be denied entirely.

**We recommend that the Draft Act clarifies the appointment procedure for mediators, the scope of mediator authority, the binding nature of mediator decisions, and the default outcome when mediation closes without agreement.**

More concerning is the limited role granted to affected researchers in the mediation process. Article 13(5) states that principal researchers will be invited to join mediation only “*where appropriate*”. This conditional involvement suggests researchers may have minimal influence during dispute resolution, despite being directly affected by outcomes.

**We recommend that the Draft Act explicitly define circumstances requiring researcher participation and establish clear protocols for keeping researchers informed throughout the mediation process.**

### **Vetting Process and Potential Impacts on Fair Representation**

Lastly, the Draft Act, while advancing crucial data access provisions, requires careful consideration of its potential downstream impacts on academic research and knowledge production. The vetting process and institutional requirements could disproportionately advantage well-resourced research institutions, particularly those within the EU. This raises important questions about equitable access to platform data and its implications for global knowledge production (e.g., Article 9 (4)(d) sets the condition “*that the computing power at the disposal of the vetted researchers is appropriate and sufficient for the purposes of the research project*”). This could concentrate platform research within a select group of institutions, potentially limiting diverse perspectives and approaches<sup>11</sup>. Without careful consideration of non-EU researcher access, the Act might inadvertently create a two-tier

---

<sup>11</sup> Nagaraj, Abhishek, Esther Shears, and Mathijs de Vaan. “Improving data access democratizes and diversifies science.” *Proceedings of the National Academy of Sciences* 117, no. 38 (2020): 23490-23498.



system of platform research. Previous written evidence submitted by colleagues at the Oxford Internet Institute also discussed data access outside the EU<sup>12</sup>. The issue is particularly concerning given that platform impacts and risks often transcend geographical boundaries, requiring diverse global perspectives for comprehensive understanding of systemic risks.

**We recommend that the Draft Act incorporate specific provisions to (a) ensure technical and infrastructural requirements do not create unnecessary barriers to entry for smaller institutions, (b) create clear guidelines for international research collaboration or consider the establishment of shared research infrastructure for researchers from non-EU institutions.**

---

<sup>12</sup> Ibrahim, Lujain, Rocher Luc., and Valdivia, Ana. "Oxford experts call for better access to platform data." *Oxford Internet Institute*, June 8, 2023, accessed Dec 01, 2024, <https://www.oii.ox.ac.uk/news-events/oxford-experts-call-for-greater-access-to-platform-data/>.