

Highlighted Practical Issues and Consensus in Supporting Research on the Information Environment

Institute for Research on the Information Environment

June 17, 2022

Abstract

The Institute for Research on the Information Environment (IRIE) has teamed with 20 partners to conduct 12 exploratory studies. This report aims to identify and highlight practical issues and commonalities. These studies highlighted a lack of diversity in both the content examined and methodologies used in ways that are consistent with the thesis that inefficient engineering practices are slowing knowledge accumulation. In addition, there was an imbalance between which platforms are being studied and which are actually used most worldwide. Generally, the exploratory studies show a wide discrepancy between what a small number of well-resourced organizations and the rest of the field produced. We provide insight into the field's condition, highlight recommendations for an ideal state, note commonalities between exploratory studies, and identify decision points to consider if IRIE moves forward.¹

¹ Results from the exploratory studies are not exhaustive of the entire information environment field. Nonetheless, the exploratory studies help to generate a better understanding of some current commonalities and challenges in the field that IRIE can address.

Funding: This research was generously supported by Microsoft, Craig Newmark Philanthropies, the John S. and James L. Knight Foundation, and the William and Flora Hewlett Foundation. This work was initiated as part of a joint project between the Partnership for Countering Influence Operations at the Carnegie Endowment for International Peace and the Empirical Studies of Conflict Project at Princeton University. **Author contributions:** Jacob N. Shapiro and Alicia Wanless provided conceptualization, editing, and supervision. Jen Rosiere Reynolds wrote the manuscript and managed research activity planning and execution. Abrianna Rhodes produced Figure 1. **Data and materials availability:** All data required to evaluate the conclusions are present in the paper or papers referenced. Additional data related to this paper may be requested from the authors.

Executive Summary

The Institute for Research on the Information Environment (IRIE), alongside 20 partners, has studied 141 organizations and instruments as part of 12 exploratory studies on the information environment. These studies analyzed characteristics of the field, institutional models, infrastructure that could speed analysis, and some of the unique challenges of this field. (For a list of all organizations and instruments examined, please see Appendix A.1.) This report identifies commonalities to guide IRIE's decisions in the next phase.

These studies highlighted a lack of diversity in both the content examined and the methodology used. Generally, the studies show a wide discrepancy between what a small number of well-resourced organizations and the rest of the field produced. Most research analyzed examined Twitter, Facebook, or both, while a much smaller percentage considered YouTube, Instagram, Reddit, Sina Weibo, and Telegram; other platforms were broadly overlooked. The majority of these publications did not use machine learning (77%) or network analysis (81%). The most common analytic techniques were simple econometric and statistical analysis (53%) and graphical or visual analysis (51%). The following most common methods of analysis among reviewed publications were descriptive (42%), qualitative (29%), machine learning (23%), and network analysis (19%).

Our research highlighted many impediments to working with social media data. While APIs were the most common way of obtaining data, these still require some coding skills, and civil society organizations struggled to get access. Many researchers in our sample also used manual annotations or coding for data collection. Data sharing agreements with platforms proved difficult for researchers to obtain. Researchers in our sample struggled with data science support and especially engineering support. Accordingly, researchers cited the need to create training and tools to expand methodology. In comparison, social media platforms appear to provide social listening and social media monitoring companies with extensive data access. If platforms could grant similar access to researchers, it would allow research to progress faster and greatly inform this important field of study.

Our research partners agreed broadly on a need for expanded data access and a thoughtful approach to do so, including a model where certain qualifications allow specific access. Researchers also highlighted operational and transparency reporting and accessibility to encrypted messaging apps, as well as a need for baselines, samples, and standardization. More broadly, researchers recommended addressing the imbalance between platforms used worldwide and those studied.

We next aimed to understand various structural models and processes others have used to tackle similarly complex and serious problems. As a result, we propose a potential access model with the common features of a binding agreement, output review, administrative support during the process, ethics review, a background check, and training.

Finally, we considered optimal funding models for further research. The three common funding models discussed were government, institutional buy-in, and blend. After analyzing the benefits and weaknesses of these models, recommendations for IRIE to explore include:

1. Utilizing government funding;
2. Developing proprietary data and tools to encourage funding;
3. Further research on investment income; and
4. Increased understanding of different financial growth pathways and models.

Table of Contents

Abstract	1
Executive Summary	2
Table of Contents	3
Introduction	4
State of the Field's Research	5
Current Hardships	7
Access	7
Use	8
Impediments to the Social Science	8
Ideal State	10
Questions of Structure to Support	11
Structure	11
Access	12
Proposed Access Model	12
Financial options	13
Conclusion	14
Appendix	15
A.1 List of Organizations/Instruments Analyzed	15

Introduction

The current information environment poses a profound challenge to democratic decision-making. Disinformation disproportionately targets marginalized groups, and content that yields strong negative emotions tends to gather the most engagement, providing an economic incentive to create incendiary content. Unfortunately, little evidence exists on which to design improved policies. Many agree that this is a problem, but there is little consensus around its depth, breadth, impact, or the stakeholders involved.

We want to understand if there are solutions and what those might entail. To do this, we aimed to better understand the state of the field, compare different institutional models, identify infrastructure that could speed discovery, and examine some of the unique challenges of studying the information environment. This report identifies and highlights practical issues and commonalities across 12 exploratory studies from 20 partners and analyzes 141 organizations and instruments to guide IRIE's decisions in the next planning phase.

We found a lack of diversity in the content studied and methodologies used. Most methods were also less sophisticated than what is available. Accordingly, organizations struggled to hire technically savvy people to procure and analyze this data; in turn, researchers were less efficient at understanding social media platforms' actions. Researchers needed operational reporting from platforms to understand the environment better, but there was no one model for managing access and vetting researchers to address these issues. A potential model for IRIE is a median model. Finally, multiple financial models existed, each with benefits and limitations. With all decisions, IRIE needs to integrate budget planning and resource considerations into scoping.

This report directly incorporates the following studies' work:

1. "Research Process 1 (RP1): Current Academic Research on the Information Environment," by Nilima Pisharody and Jen Rosiere Reynolds
2. "Research Process 2 (RP2): Social Media Data in Conflict Research," by Jane Esberg and Nejlja Asimovic
3. "Research Process 3 (RP3): A Survey of Public-Oriented Organizations Analyzing Social-Media Disinformation," by Darren L. Linvill and Patrick L. Warren
4. "Research Process 4 (RP4): Scoping the Institute for Research on the Information Environment," by Nils B. Weidmann, Margaret E. Roberts, Zachary Steinert-Threlkeld, and Sebastian Hellmeier
5. "Research Process 5 (RP5): Civil Society Data Access Needs for Social Media Research," by Samantha Bradshaw and Bridget Barrett
6. "Research Process 6 (RP6): Data Requirements for Understanding Monetization and the Information Environment," by Danny Rogers
7. "Research Process 7 (RP7): Accelerating Research with Multi-National, Multi-Platform Image Archives," by Cody Buntain

8. "Research Administration 1 (RA1): Peer Review and Access Models for Large-Scale Scientific Instruments," by Kristen DeCaires Gall and Diego A. Martin
9. "Research Administration 2 (RA2): Financial Models of Large-Scale Scientific Instruments and Organizations," by Kamya Yadav and Jen Rosiere Reynolds
10. "Research Administration 3 (RA3): Existing Initiatives' In-House Technical Capabilities," by Victoria Smith and Jen Rosiere Reynolds
11. "Privacy and Ethics 1 (PE1): Researcher Access to Restricted Government Data," by Jen Rosiere Reynolds, Aditi Bawa and Kamya Yadav
12. "Privacy and Ethics 2 (PE2): Social Listening Companies and Access to Sensitive Data," by Kamya Yadav and Alicia Wanless

State of the Field's Research

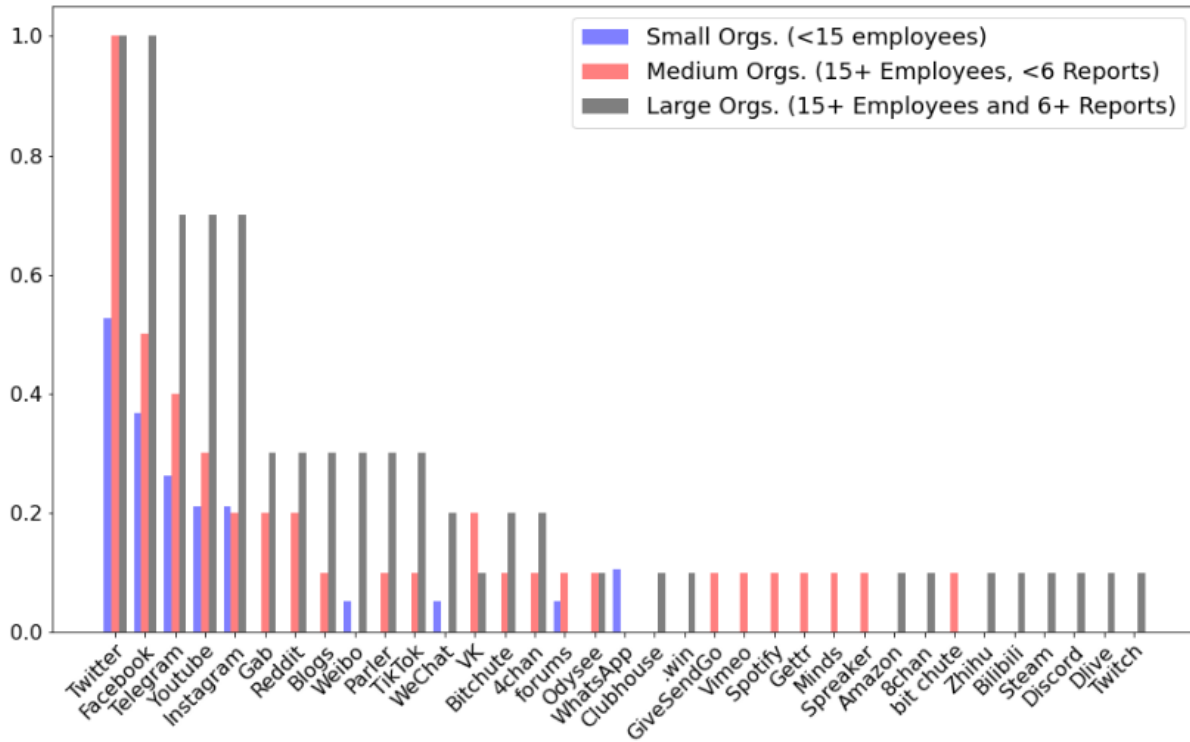
Across the exploratory studies, there was a lack of diversity in the content studied and methodologies used. Most methods were less sophisticated than what is available. Most research examined covers Twitter, Facebook, or both, with few studies focused on other platforms. Linvill and Warren's Figure 2 highlights this.

1. Academic publications from top research journals tended to focus on Twitter (59%); a fair number of these studies (22% of the 59%) do so in conjunction with other platforms. 26% of publications looked at Facebook.²
2. A combined 47% of all research organizations looked at Telegram, YouTube and Instagram.³
3. The subsequent majority of publications analyzed Reddit (7%), YouTube (5%), Instagram (5%), and Sina Weibo (3%).²

² Nilima Pisharody and Jen Rosiere Reynolds, "Current Academic Research on the Information Environment," Carnegie Endowment for International Peace, Princeton University, June 16, 2022: 5.

³ Darren L. Linvill & Patrick L. Warren, "A Survey of Public-Oriented Organizations Analyzing Social-Media Disinformation," Clemson University Media Forensics Lab, 2022:5.

Figure 2. Share of Organizations with a Report Covering Each Platform



Source: Linvill and Warren, "A Survey of Public-Oriented Organizations Analyzing Social-Media Disinformation."

Among publications in top journals, we sampled throughout the quality distribution, at least as judged by normalized citation counts, to learn more about methodology.

1. The vast majority (83%) of papers analyzed social media at the post level, and most used textual analysis to do so (68%).
2. 23% of publications reviewed used machine learning, and 19% used network analysis.
3. 27% used feature extraction, most of that being Natural Language Processing (80%).
4. 12% tracked information across more than one platform.
5. The most common analytic techniques were simple econometric and statistical analysis (53%) and graphical or visual analysis (51%). The next most common method of analyses were descriptive (42%), qualitative (29%), machine learning (23%), and network analysis (19%).
6. 49% of the papers in top journals used a platform API to obtain data, and 34% used manual annotations or coding. The subsequent most popular methods (in descending order) were crowdsourcing or surveys (14%), independent collections (14%), scraping (9%), and 1% used social listening tools or bots the author created.

While some APIs restricted access to university researchers, the civil society organizations used APIs when they could; they celebrated CrowdTangle as one of the few avenues they gained data access. In addition, the civil society organizations (CSOs) also highlighted Whatsapp as a primarily manual data

collection and message forwarding platform, rather than a comprehensive system for data collection like the other platforms mentioned.

Lack of Computational and Data Science Expertise

Analyzing the information environment requires significant preexisting knowledge and assets. Unfortunately, organizations struggled to hire technically savvy talent. Out of the existing centers institutes doing repeated social scientific or descriptive analysis on the information environment, 41% did not have dedicated in-house data science staff, and 67% lacked engineering support. In line with these percentages, groups reported finding it slightly easier to obtain data science support than engineering support. Generally, civil society organizations in the field also lacked the expertise or resources to perform computational analysis and data science work. Personnel costs were prohibitive, and outsourcing is difficult due to limited knowledge of and access to the market. Fact-checking organizations did have some support through the International Fact Checking Network, which provided tools, data access, and training.

For comparison, social listening and social media monitoring companies used APIs, third-party cookie crawlers, and AI-powered systems to collect data and help brands improve their sales. The data types collected were content, demographic, identification, and location. Ultimately, it appears that social media platforms provided them extensive access to their data, suggesting that platforms could share similar data with researchers.

Current Hardships

Researchers struggle with accessing and using data and the lack of platform operational reporting.

Access

APIs require some coding skills and are not always available to those without a university affiliation. This lack of accessibility significantly limits effective scholarship. Civil society organizations also have particular needs as their teams frequently lack computational skills and the expertise or resources needed to negotiate data-sharing agreements with platforms, particularly CSOs with different language skills, as platforms' contracts are often written in English. Even if organizations succeed and establish data-sharing arrangements with platforms, there are still many limitations on how their researchers can use the data. For example, contracts generally prohibit platform data from being made public or being used to create other tools for researchers.

Non-profit, private sector, and academic organizations, as well as those specifically studying conflict areas, stressed the difficulty of accessing content that has been moderated or taken down, also known as the "black hole problem." The examined organizations also highlighted the issue of access to and interpretation of user security and data privacy. CSOs and those studying conflict areas mentioned a need for access to encrypted messaging apps. Esberg and Asimovic also noted a lack of access to location metadata, content reach, and recommendations as limitations. Notably, we found that

metadata—whether of individual messages or media content—was seldom analyzed among non-profit, private-sector, and academic organizations, possibly due to lack of access.

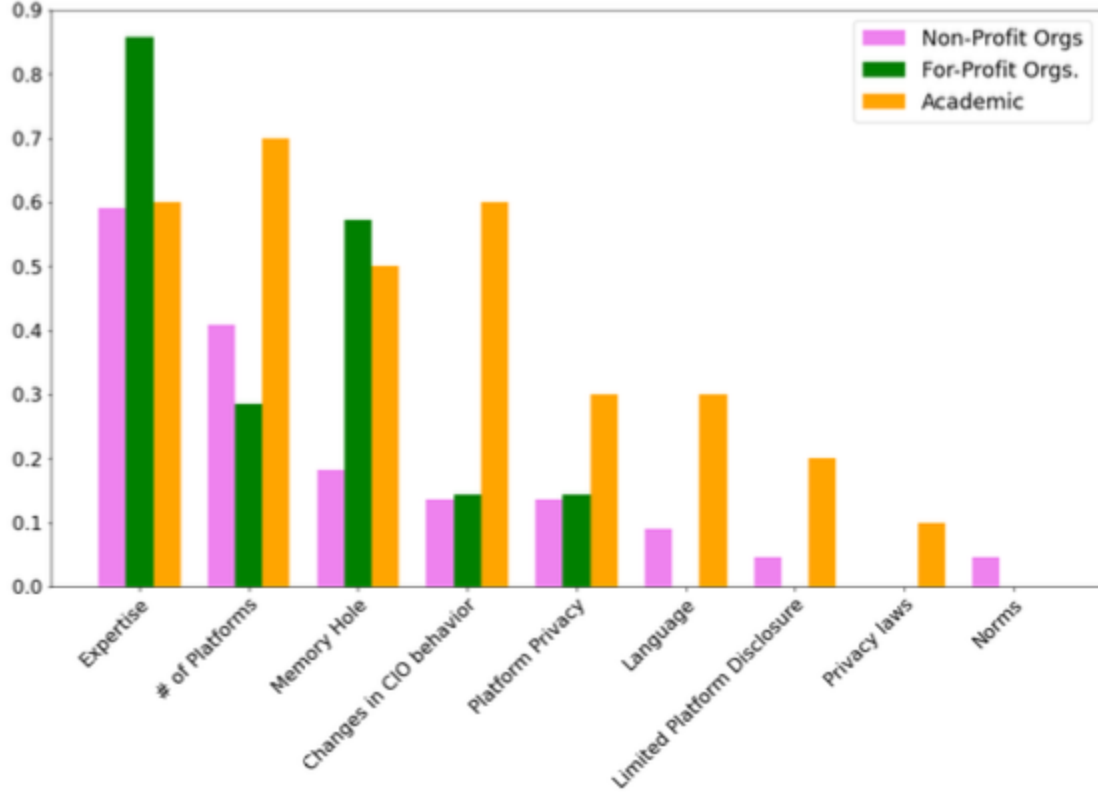
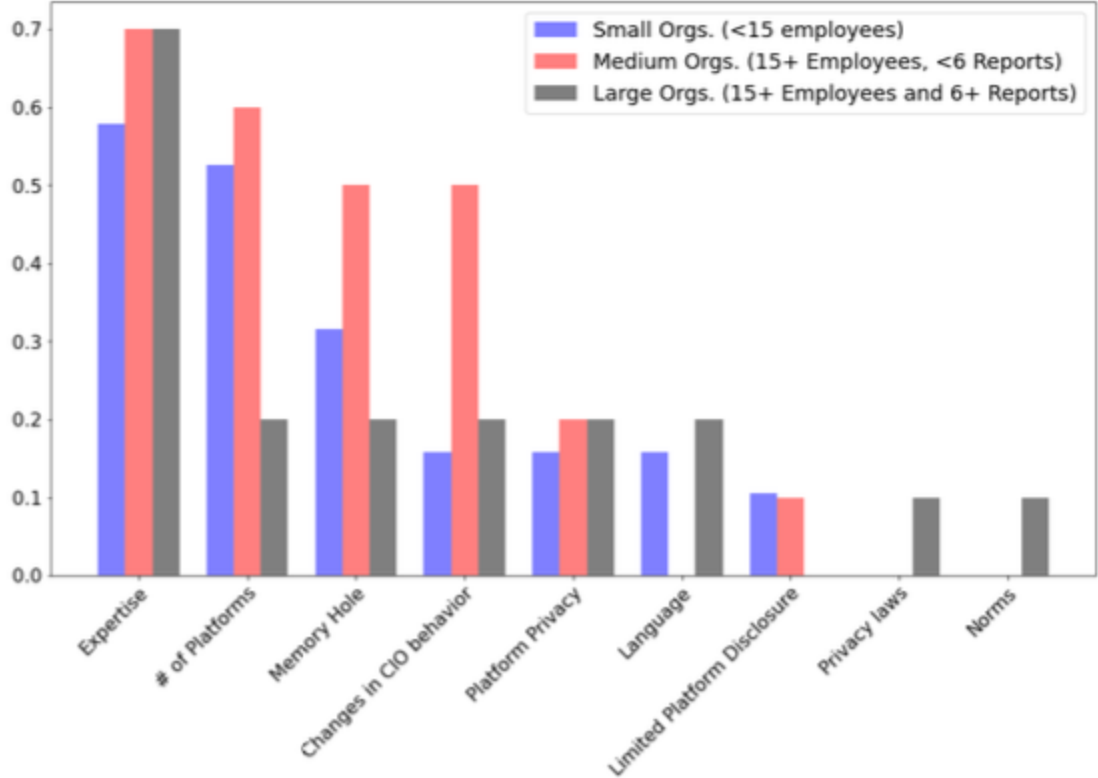
Use

Another significant limitation relates to the use of data once obtained. Linvill and Warren assessed that many non-profit, private sector, and academic organizations “do not seem to know what they do not know, perhaps driven by the relative nascence of the field.” To maximize effectiveness and impact, quantitative methodologies must be married with qualitative knowledge; without that broader knowledge to explain the data, it is very difficult to translate them into meaningful messaging or policy change. We found that overall, statistical analyses were scarce in public-oriented organizations’ publications; most were simply limited to descriptive statistics.

Impediments to the Social Science

Esberg and Asimovic cited non-human activity and manipulations, such as bots and trolls, and bias in content moderation in social media companies as a hindrance to social science. While it is not possible to eliminate these elements from the information environment, a lack of operational reporting on these activities impedes researchers from understanding human users’ authentic experiences and opinions.

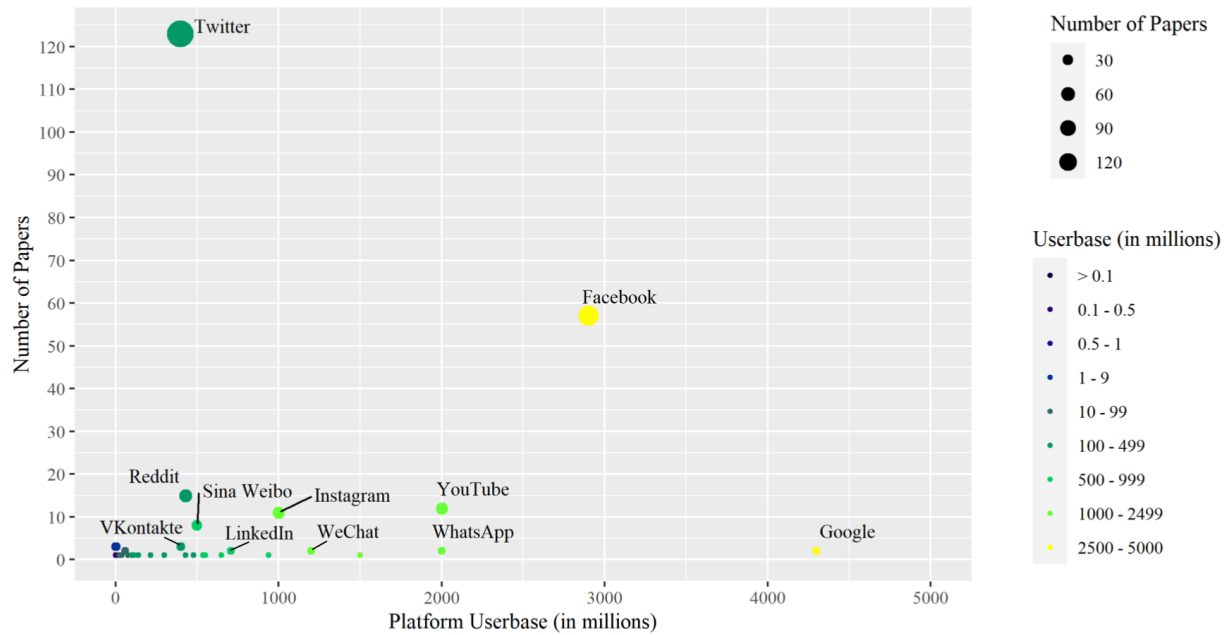
Figure 8. Share of Organizations with a Report Subject to Each Limitation



Source: Linvill and Warren, "A Survey of Public-Oriented Organizations Analyzing Social-Media Disinformation."

A final and transcending limitation is the imbalance between platforms used worldwide and those studied. Both the top journal publications and the public-oriented research publications heavily focused on Western populations, at the expense of other groups and platforms. See Figure 1 for relative number of papers by platform on size of user base from Pisharody and Rosiere Reynolds' "Current Academic Research on the Information Environment" sample.

Figure 1. Number of Papers Written on Platforms vs. Platform Userbase Size



Ideal State

Researchers within the counter-influence community have frequently cited data access as a problem. However, "data access likely means different things to different researchers."⁴ As a result, significant

⁴ Victoria Smith, "Existing Initiatives' In-House Technical Capabilities," Carnegie Endowment for International Peace, June 7, 2022.

work has been published on transparency to guide conversations around these topics.^{5 6 7 8 9} This group of studies identified the following data access requests:

1. First, a deep discussion about the circumstances under which making data public is appropriate;
2. Data on exposure;
3. Data on impressions;
4. Day-to-day use of specific hashtags and viral posts;
5. Engagement;
6. Expanded access to APIs;
7. Expanded data features extracted through APIs;
8. Expanded political ads libraries; and
9. Increased access to more platforms.

Researchers also want more operational reporting, including:

1. Algorithmic recommendations;
2. Amounts of money going to particular accounts (e.g., YouTube channels)
3. Baseline measure for both the number of active users in a country;
4. Content and accounts that platforms have removed;
5. Content moderation process;
6. Content moderation staff;
7. Coordinated inauthentic behavior;
8. Country-level or language-level breakdowns in reporting;
9. Internal platform research;
10. A standard structure of transparency reporting;
11. Traffic numbers to and the number of bid requests from each site; and
12. Who was paid to promote and the amount they were paid.

Two studies recommended models of providing access at different levels of granularity and aggregation, particularly concerning location information and moderated or taken-down content.

⁵ Alex Abdo et. al. , “A Safe Harbor for Platform Research,” Knight First Amendment Institute at Columbia University, January 19, 2022, <https://knightcolumbia.org/content/a-safe-harbor-for-platform-research>.

⁶ Hultquist, John, “Anticipating Cyber Threats as the Ukraine Crisis Escalates,” Mandiant, January 20, 2022, <https://www.mandiant.com/resources/ukraine-crisis-cyber-threats>.

⁷ GIFTC, “GIFTC Transparency Working Group: One-Year Review of Discussions” Global Internet Forum to Counter Terrorism, July 2021, <https://gifct.org/wp-content/uploads/2021/07/GIFTC-WorkingGroup21-OneYearReview.pdf>.

⁸ Caitlin Vogus and Emma Llansó, “Making Transparency Meaningful: A Framework for Policymakers,” Center for Democracy and Technology, December 14, 2021, <https://cdt.org/insights/report-making-transparency-meaningful-a-framework-for-policymakers/>

⁹ Heidi Tworek and Alicia Wanless, “Time for Transparency from Digital Platforms, But What Does That Really Mean?” Lawfare, January 20, 2022, <https://www.lawfareblog.com/time-transparency-digital-platforms-what-does-really-mean>.

Once a researcher obtains social media data, understanding and processing it often also proves challenging. Therefore, multiple research groups recommended creating training and tools to expand methodology. For example, in the case of images, Buntain suggested the creation of user-friendly tools that "allow downloading or visualizing data in contexts where technical infrastructure may not support the extraction and the analysis of large amounts of raw data."¹⁰

Finally, some broad recommendations for the field included cultivating a more global focus and creating baselines, samples, and standardization. This could look like a baseline collection of images, a sample set of accounts and images from one platform for a specific country, and the development of standard definitions for the field.

Questions of Structure to Support

There is currently no one model for managing access and vetting researchers studying the information environment. IRIE would fill this gap by providing a balanced model. To create this model, we examined the organizational structure, peer review, and access models for 17 large-scale scientific instruments from 13 different research organizations across multiple fields and 31 access procedures to access restricted data from selected government institutions in five countries. We aimed to better understand various structural models and processes of how others have tackled problems of considerable seriousness and an ambitiously large scale.

Structure

47% of the large-scale scientific instruments had a multinational leadership model, and 29% were housed at academic research centers. Except for the instruments with public access, all large-scale scientific instruments incorporated scientific experts into their governing body, access process, or external peer review committees. Of the instruments with publicly available governance models, 75% had a council, board, or committee overseeing administrators' daily operations and decisions. Administrators have responsibility for the organization, operations, and budget. In contrast, 12% of the instruments are governed and administered via a member committee. Membership requirements varied, and 80% of the instruments had some national or citizenship condition. However, all had exceptions or resource allotments reserved for nonmembers. Another 12% noted that members directly elect the board.

Access

If IRIE moves forward, it must determine an appropriate access model, as resources will be constrained to some degree. In addition, the solution IRIE builds may need to be restricted for privacy reasons. To help these decisions, we examined 17 large-scale scientific instruments and 31 access procedures for restricted government data.

¹⁰ Cody Buntain, "Accelerating Research with Multi-National, Multi-Platform Image Archives," Institute for Research on the Information Environment, June 16, 2022.

Governments restricted data access using one or more protection mechanisms. Detailed eligibility requirements were prevalent but not universal. For example, about a third of these processes called for some alignment (such as in values or mission) between the researcher and the data-owning entity. About 30% had requirements for education, experience, or skills needed to work with the data, such as several courses in data science. About 20% mandated the project be generally for the public good. All access processes we studied required approval, and the majority required a binding agreement and output review and provided an advisor before or during the process. Further frequent commonalities included ethics approval, a background check, and a researcher fee.

Similar to the government process, out of the large-scale instruments examined, 65% had some process for access or project approval, and all began with a request or application. These processes were much less onerous. The only publicly available requirements were that 73% of the large scale instrument processes included a project proposal as part of the request, and 45% included an explicit peer review process. Only 27% included output review, and 36% a contract, whereas most government processes required both.

Proposed Access Model

One proposed access model for IRIE could include the common features of a binding agreement, output review, administrative support during the process, ethics review, a background check, and training. For example:

1. A researcher submits a research proposal that includes an ethics committee approval. This ethics committee may be unique to IRIE or part of a partner institution, such as a university.
2. The proposal undergoes a committee review (possibly just meaning that more than one IRIE staff member reviews). It is accepted or rejected, or researchers are requested to clarify, revise, and resubmit.
3. Once the project has all necessary approvals, the researchers undergo a background check to verify identity, educational and professional credentials, and personal and professional references.
4. Researchers must then attend an orientation session, during which they sign a contract.
5. All researchers must take annual training in proper data stewardship.
6. When the researcher has completed work, an IRIE advisor will conduct a disclosure or confidentiality review of all project outputs to protect data confidentiality. This advisor could be a faculty affiliate or IRIE staff.

Financial options

Multiple financial models existed, each with unique tradeoffs and promising features for IRIE. The organizations' work and size greatly influenced how they allocate funding. Very few similarities emerged due to the diverse instruments and organizations we analyzed. Three common models were:

1. Government: We found that government funding was common for instruments and organizations across various fields of study. A government or multiple governments either partially or wholly funded all the organizations studied.
2. Institutional buy-in: This is another attractive funding model. Here, academic institutions and individuals can buy into the instrument for early data access and contribute to the instrument's construction, management, and operation. Both the Sloan Digital Sky Survey and part of the Dark Energy Spectroscopic Instrument use the institutional buy-in model.
3. Blend: Other funding sources included foundations, universities, private corporations, and investments. Half of the instruments and organizations we studied were supported by some combination of government(s), charities, universities, and other funding.

How instruments and organizations in our dataset distributed their funds also varied significantly. For example, astrophysics instruments primarily constructed, developed, maintained, and operated themselves. In contrast, social science instruments' budgets funded data collection. Finally, the research institutions reviewed mostly spent on project grants or personnel.

Additionally, around 46 percent of all organizations and instruments distributed grants. These grants were either project-based or for institution building.

If IRIE moves forward and establishes its research mandate and scope of operations, the team can confirm a more nuanced decision on financial models. These decisions include the proposed access model and the staff and logistic support detailed above. Accordingly, we propose the following financial recommendations for IRIE to explore:

1. Secure government funding;
2. Develop proprietary data and tools to encourage funding;
3. Explore an investment portfolio or endowment; and
4. Further research growth pathways and aligned funding models, including institutional buy-in.

Conclusion

While this field is new, it is not unstudied. In 2018, a group of scholars convened to discuss the challenges of creating a framework for sharing sensitive online data. They examined how large-scale social and digital data could be collected, shared, and used by researchers while protecting the rights and privacy of individuals represented by the data. Their suggestions included:

1. Differential privacy;
2. Tracking queries and how data are used;
3. Creating standardized definitions of sensitivity and who can access what types of sensitive data; and

4. Establishing partnerships between companies and researchers, whether through an institution that acts as a bridge or companies hiring researchers during their project.¹¹

These suggestions mirror our exploratory studies and add tangible tactical implementation tools to use if IRIE moves forward.

The 12 exploratory studies across 20 partners highlighted a lack of diversity, an imbalance between platforms used and analyzed, and a wide dissimilarity between a small number of organizations and the rest of the field. They illustrate ample opportunity to advance the study of the information environment and enable evidence-based policymaking.

¹¹ Lazer, David, Joshua A. Tucker, Jonathan Nagler, Liliana Mason, and Adam Berinsky. ISSOD Proceedings. Sloan Foundation, 2018. <https://securelysharingdata.com/overview.html>.

Appendix

A.1 List of Organizations/Instruments Analyzed