

# The craft and coordination of data curation: complicating "workflow" views of data science

ANDREA K. THOMER, School of Information, University of Michigan, USA

DHARMA AKMON, ICPSR, University of Michigan, USA

JEREMY YORK, School of Information, University of Michigan, USA

ALLISON R. B. TYLER, School of Information, University of Michigan, USA

FAYE POLASEK, School of Information, University of Michigan, USA

SARA LAFIA, ICPSR, University of Michigan, USA

LIBBY HEMPHILL, School of Information & ICPSR, University of Michigan, USA

ELIZABETH YAKEL, School of Information, University of Michigan, USA

Data curation is the process of making a dataset fit-for-use and archiveable. It is critical to data-intensive science because it makes complex data pipelines possible, makes studies reproducible, and makes data (re)usable. Yet the complexities of the hands-on, technical and intellectual work of data curation is frequently overlooked or downplayed. Obscuring the work of data curation not only renders the labor and contributions of the data curators invisible; it also makes it harder to tease out the impact curators' work has on the later usability, reliability, and reproducibility of data. To better understand the specific work of data curation – and thereby, explore ways of showing curators' impact – we conducted a close examination of data curation at a large social science data repository, the Inter-university Consortium of Political and Social Research (ICPSR). We asked, What does curatorial work entail at ICPSR, and what work is more or less visible to different stakeholders and in different contexts? And, how is that curatorial work coordinated across the organization? We triangulate accounts of data curation from interviews and records of curation in Jira tickets to develop a rich and detailed account of curatorial work. We find that curators describe a number of craft practices needed to perform their work, which defies the rote sequence of events implied by many lifecycle or workflow models. **Further, we show how best practices and craft practices are deeply intertwined.**

CCS Concepts: • **Human-centered computing** → *Collaborative and social computing systems and tools*; • **Applied computing** → *Document preparation*; • **Information systems** → *Digital libraries and archives*.

Additional Key Words and Phrases: data curation, knowledge infrastructure, craft, coordination, workflows, social science data

## ACM Reference Format:

Andrea K. Thomer, Dharma Akmon, Jeremy York, Allison R. B. Tyler, Faye Polasek, Sara Lafia, Libby Hemphill, and Elizabeth Yakel. 2022. The craft and coordination of data curation: complicating "workflow" views of data science. 1, 1 (February 2022), 27 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

---

Authors' addresses: Andrea K. Thomer, athomer@umich.edu, School of Information, University of Michigan, Ann Arbor, Michigan, USA; Dharma Akmon, ICPSR, University of Michigan, Ann Arbor, Michigan, USA; Jeremy York, School of Information, University of Michigan, Ann Arbor, Michigan, USA; Allison R. B. Tyler, School of Information, University of Michigan, Ann Arbor, Michigan, USA; Faye Polasek, School of Information, University of Michigan, Ann Arbor, Michigan, USA; Sara Lafia, ICPSR, University of Michigan, Ann Arbor, Michigan, USA; Libby Hemphill, School of Information & ICPSR, University of Michigan, Ann Arbor, Michigan, USA; Elizabeth Yakel, School of Information, University of Michigan, Ann Arbor, Michigan, USA.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2022 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

## 1 INTRODUCTION

Data curation – the technical work put into datasets to make them fit-for-use and accessible over the long-term – is critical to data-intensive science [Borgman et al. 2019; Faniel and Zimmerman 2011; Hey T., Tansley S., Tolle K. 2009; Muller et al. 2019; Palmer et al. 2013]. In data science contexts, this work is often referred to as munging, wrangling, or processing, with a particular focus on the working data into a usable format; this work of making data fit-for-use can take up to 80% of a data scientists’ day-to-day work [Wickham 2014]. In institutional contexts – for instance, in large scientific data archives or institutional repositories – data curation is likely to involve the application of standards in addition to data munging, with a particular focus in making data shareable and easy to reuse. In both contexts: the ways in which data are transformed and manipulated prior to analysis have significant impacts on the quality and reliability of a study [Borgman 2015; D’Ignazio and Klein 2020; National Academies of Sciences, Engineering, and Medicine et al. 2019]. Additionally, making data ready to archive or share is increasingly required by both funding agencies and journals.

Yet, despite its importance, data curation – and data curators – are often overlooked in accounts of data science. Job ads for data scientists frequently call for data “unicorns,” “ninjas,” and “rock stars” to wrangle messy datasets through mythic abilities (not through skill or craft), or “janitors” as if data processing were a rote sanitizing process that requires little specialized expertise or training [D’Ignazio and Klein 2020]. Data science clients similarly think of data work as “magic” [Kross and Guo 2021] – which, while seemingly complimentary, obscures the skill and effort needed for this work. In addition to obscuring the skill needed for curatorial work, this invisibilization obscures the varied judgements, decisions, and data processing steps that go into data processing and have big impacts on the final trustworthiness, reproducibility and auditability of a study.

The work of data curation can also be obscured, somewhat ironically, through attempts to render it visible as part of a regularized workflow. Workflows and curatorial best practices aim to break curation into a discrete set of steps, or show it as one ‘phase’ of work in a project (e.g. Higgins [2008]; Muller et al. [2019]). The goal of workflow representations is to make curation more reproducible and routine – but it comes at the cost of obscuring the skill needed to do these tasks well and furthering the idea that curatorial work is a rote task that just any human can be plugged into.

Obscuring data curation also renders the labor and contributions of data curators invisible [Plantin 2019]. Like other forms of service work, well-executed curation is hidden [Suchman 1995]. Additionally, obfuscating curatorial work makes it challenging to understand the impact of specific curatorial actions, and therefore to efficiently prioritize, plan, or fund data curation (anonymized for review). Without understanding the impact of data curation, the developers of curatorial tools cannot assess or prioritize which features and functionalities will best increase curatorial efficacy or later data reuse.

To better make the work of data curation visible, we conducted a close examination of data curation at a large social science data archive, the Inter-university Consortium for Social and Political Research (ICPSR). ICPSR recently adapted external standards and best professional practices to create robust internal guidelines for curation, and the scale, centrality, and collaborative aspect of curatorial work at ICPSR make it an excellent site for a case study of data curation. ICPSR is the largest social science data archive in the world, and it contains datasets from over 16,000 studies. ICPSR’s professional, in-house curation activities distinguish it from other data repositories such as the UCI Machine Learning Repository<sup>1</sup> or Dataverse<sup>2</sup> where data providers are expected to curate data themselves.

<sup>1</sup><https://archive.ics.uci.edu/ml/index.php>

<sup>2</sup><https://dataverse.harvard.edu/>

105 This research is part of a larger project focused on understanding the impact of curatorial work, and aimed at  
106 developing metrics that better measure and account for the benefits of that work. We aim to make curatorial work more  
107 visible, and thereby easier to account for in budgets and in academic promotion cases. Our methods include interviews  
108 with ICPSR stakeholders, as well as computational analysis of curation logs. Here, we address the following research  
109 questions:  
110

- 111 (1) What does curatorial work entail at ICPSR, and what work is more or less visible to different stakeholders and in  
112 different contexts?  
113
- 114 (2) How is that curatorial work coordinated across the organization?  
115

116 We drafted these questions with the goal of understanding both the visible and invisible work that goes into data  
117 curation; prior studies have shown that much of this work escapes view [Plantin \[2019\]](#), but fewer have sought to specify  
118 what, exactly, is invisible.  
119

120 By triangulating accounts of data curation from interviews and records of curation in Jira tickets, we develop a  
121 rich and detailed account of curatorial work at ICPSR. In doing so, we bridge research in CSCW and the library and  
122 archival sciences on data curation. We find that while there are several standard curatorial activities performed at  
123 ICPSR, and well defined standards for different "levels" of curation, considerable craft and coordination are needed to do  
124 this work well; in other words, craft is needed to "work" a workflow. This non-technical work is necessary to facilitate  
125 technical work, and has been less well-defined in prior discussions of data curation. Surfacing the role of craft and  
126 coordination has important implications for curatorial projects' and teams' planning. This defies the rote sequencing  
127 of events implied by many lifecycle or workflow models. We provide a detailed account of curatorial workflows at  
128 ICPSR and explain how workflow-based accounts of data curation can obscure both the individual skilled "artistry" and  
129 coordination necessary in this work.  
130

131 We additionally reflect on the visibility of data curation, both within ICPSR and to data users. As [Plantin \[2019\]](#)  
132 has previously described, much of data curators' work is intentionally kept invisible to the final data consumers – yet,  
133 curators experience their jobs as being hypervisible to their supervisors, via the extensive documentation they create.  
134 We discuss how different kinds of invisible work are at play in data curation at ICPSR and explain how CSCW and data  
135 science can benefit from better understanding the judgements and skill that goes into effective data curation.  
136  
137  
138

## 139 2 PRIOR WORK

### 140 2.1 What is data curation?

141 For the purposes of our research at ICPSR, we have defined data curation as the work needed to make a dataset fit-for-use  
142 over the long term (anonymized for review) A detailed description of data curation activities at ICPSR is provided in  
143 Sections 3.1 and 3.1.2 . Depending on the scholarly community, researchers have described this kind of work in varying  
144 ways and employed several research data management-related terms as synonyms.  
145  
146

147 Much of the CSCW literature discusses curation in terms of fitness for use aimed at a single user or end goal [[Feger  
148 et al. 2020](#); [Kandel et al. 2011](#); [Taylor et al. 2015](#)]. For example, Feger et al. focus on the activities "essential for generating  
149 reproducible artefacts," while Kandel et al. are concerned with "wrangling" data to enable meaningful analysis for the  
150 research at hand. Others, such as [[Muller et al. 2019](#)] define curation as a type of human "intervention" that includes  
151 data-cleaning, converting metadata, and data alignment. In proposing "Datasheets for datasets", Gebu and colleagues  
152 [[2021](#)] argued that data used in machine learning should carry documentation that describes its collection processes  
153 and transformations, two steps in the process of curating data for reuse. For data science workers, data curation is  
154  
155  
156

157 a collaborative activity centered on information exchange and data and code transparency [Zhang et al. 2020]. In  
158 their ethnographic study of the Long Term Ecological Research Network, Karasti et al. [2006] describe “information  
159 managers” data stewardship strategies and strengths, including their ability to turn localized, heterogeneous data into  
160 a networked resource. In this way, data curation, and the development of information infrastructures, is a long-term  
161 sociotechnical endeavor.

162  
163 Libraries, archives, and data repositories conceptualize the aims of curation more broadly, emphasizing supporting  
164 data’s long-term preservation and usability for a multitude of potential future purposes. Researchers in this area have  
165 defined digital curation as “the active involvement of information professionals in the management, including the  
166 preservation, of digital data for future use” [Yakel 2007], while *data curation* refers to active management throughout  
167 the data lifecycle [Palmer et al. 2013].<sup>3</sup> Models of data curation, such as the Digital Curation Centre’s Curation Lifecycle  
168 Model [Higgins 2008] and the Big Data Lifecycle Model [Pouchard 2016], identify different curation needs at different  
169 phases of the lifecycle of data and provide a workflow for curation work performed initially by data producers and  
170 then by the data curators and processors at archives and repositories. Other models take the form of terminological  
171 frameworks, such as the “Data Practices Vocabulary” by Chao et al. [2015], which outlines a taxonomy of terms  
172 describing curatorial work. There are overlaps between CSCW and LIS/archival conceptualizations of data curation; the  
173 differences tend to be in the detail in which curatorial work is described and the focus on a long-term future for the  
174 data in LIS/archival science. For instance, Higgins [2008] description of the entire curation lifecycle mirrors the five  
175 stages of data science work practices by Muller et al. [2019]. Higgins’ model includes 8 sequential phases of work with  
176 data: Conceptualize, Create/Receive, Appraise/Select, Ingest, Preservation Action, Store, Access Use and Reuse, and  
177 Transform. Muller et al.’s data science lifecycle model includes 5 sequential phases of work, with curation in the center:  
178 Discovery, Capture, Curation, Design, and Creation (2019).

## 184 2.2 What renders data curation invisible?

185  
186 What may not be clear from all these definitions of curation work is that when curation is done well, it is invisible  
187 to the data’s users. Successful data curation enables reusers to readily access and use datasets and does not highlight  
188 the data transformations or metadata generation that make the data ready for use. Curation work can be done by the  
189 data producers themselves as part of their research data management process. This curation work includes cleaning,  
190 organizing, and storing their data for their own localized research needs [Wallis et al. 2008]. Often, however, curation  
191 tasks—the data manipulation, cleaning, documentation, preservation, and other work discussed below—are carried out  
192 by data curators or processors within the repository or archive that has selected the data for inclusion in its holdings  
193 [Johnston et al. 2018].

194  
195 The role of the data archive is two-fold: to enable the researcher to share their data with the scientific community,  
196 and to ensure the data’s long-term accessibility and preservation [Green and Gutmann 2007]. The data producer deposits  
197 their data with the archive, and then at some future moment, the data appear in a standardized form, ready for use,  
198 with inconsistencies smoothed over and issues addressed, presented to the data user without the explicit traces of the  
199 curation work that are only visible to those within the data archive [Kervin et al. 2014; Plantin 2019]. By producing  
200 data according to professional standards, for instance, curators purposefully render themselves and their work invisible  
201 [Plantin 2019].

202  
203  
204  
205  
206 <sup>3</sup>“Digital curation” is often used as a broader, more general term for management of any collection of digital objects, whereas “data curation” refers to the  
207 long-term care and management of data specifically. We focus on data curation in this paper, but draw on research on digital curation and preservation  
208 where relevant.

209 The invisibility of the work makes it hard to classify. Prior scholars have explained myriad ways that work can  
210 be invisible. For instance, Nardi and Engeström [1999] identify four types of invisibility at work: 1) work done in  
211 invisible places, such as the highly skilled behind-the-scenes work of reference librarians; 2) work defined as routine or  
212 manual that actually requires considerable problem solving and knowledge, such as the work of telephone operators;  
213 3) work done by invisible people such as domestics; and 4) informal work processes that are not part of anybody's  
214 job description but which are crucial for the collective functioning of the workplace, such as regular but open-ended  
215 meetings without a specific agenda, informal conversations, gossip, humor, storytelling.  
216

217 Similarly, Star and Strauss [1999] propose three forms of invisible work: where "the act of working or the product  
218 of work is visible to both employer and employee, but the employee is invisible"; where the "workers themselves are  
219 quite visible, yet the work they perform is invisible or relegated to a background of expectation"; and, when "both work  
220 and people may come to be defined as invisible" according to particular indicators. Curators possess different types of  
221 invisibility. For example, D'Ignazio and Klein [2020] recognize the highly skilled behind-the-scenes work prevalent in  
222 what they term "data cleaning" at the same time recognizing that others discount the intellectual work required. Kross  
223 and Guo [2021] report on the black-boxing of curation that leads clients to deem the results "magic."  
224

225 In social computing, curation work is sometimes invisible because it occurs during data collection or generation.  
226 Machine learning, computer vision, and social media studies often use "found" data [Hemphill et al. 2021; Jo and Gebru  
227 2020; Paullada et al. 2021] and render curatorial decisions such as "what data should be available," "in which format(s)  
228 should data be provided," or "how should this data be sampled" invisible. For instance, datasets scraped from the web  
229 (such as Flickr photos [Scheuerman et al. 2021; Zhang et al. 2015] or Wikipedia talk pages [Wulczyn et al. 2016, 2017])  
230 suffer from biases in representation [Jo and Gebru 2020]. The kinds of curation activities that occur in archives could  
231 address those biases by adjusting samples, weights, or documentation. Annotation processes in which humans add  
232 labels to data that can then be used in machine learning tasks (e.g., facial recognition, hate speech detection), are  
233 another data generation step that are often minimized in reports about the research that depend on them. For instance,  
234 Scheuerman et al. [2021] explain that reference datasets used in computer vision tasks in papers are not well-described,  
235 and the details of the annotation process and potential biases introduced are missing. They argue that the value of  
236 "efficiency" is responsible for this pattern and that explicitly working toward other values such as "care" could improve  
237 data curation practices in computer vision.  
238

### 244 2.3 Craft in data work

245 Throughout the CSCW literature, there has been discussion of the craft needed in technical work; recent papers have  
246 begun to apply this framework more specifically to work with data. Barley and Orr [1997]'s well known volume on  
247 the topic argues that "technical work sits at the intersection of craft and science, combining attributes of each that are  
248 normally thought to be incompatible." These attributes include the use of complex technologies; a reliance on contextual  
249 knowledge and skill; the development of abstract conceptual representations to guide work; and a grounding in a  
250 community of practice [Barley et al. 2020]. Rosner et al. [2018] argue that appreciation of craft work in computer science,  
251 though, is marred by "gendered narratives" about the value of such labor, which, "both haunt and inform HCI's ideas of  
252 technological belonging, participation, and differentiation." In the context of data work, scholars have focused on how  
253 craft practices are used to process and interpret data. Mentis et al. [2016] examine surgeons collaborating remotely  
254 over image data to craft a shared interpretation. More recently, Muller et al. [2019] view the data science pipeline  
255 process as one of crafting the data, in which workers combine technical skill, expertise working with abstraction and  
256

261 representations and decision-making to accommodate unexpected issues and application of more routine techniques  
262 and automated scripts.

263 Within the information sciences, discussion of craft has largely focused on its role as part of librarianship and  
264 archival practice. Archivist Trevor Owens brings these conversations forward to a digital context, in his discussion of  
265 digital preservation (an aspect of data curation) as a craft, “best understood as part of an ongoing professional dialog on  
266 related but competing notions of preservation that goes back to the very beginnings of our civilizations” Owens [2018].  
267 He further writes,  
268

270 “digital preservation must be a craft and not a science because its praxis is; 1) grounded in an ongoing and  
271 unresolved dialog with the preservation professions and 2) it must be responsive to the inherent messiness  
272 and historically contingent nature of the logics of computing.”  
273

274 When this craftful work is done collaboratively, considerable articulation work and coordination are needed to do it  
275 successfully. Articulation work, “consists of all the tasks needed to coordinate a particular task, including scheduling  
276 subtasks, recovering from errors, and assembling resources” [Gerson and Star 1986], whereas coordination is the  
277 “process expertise” entailed in said scheduling and assembly [Barley et al. 2020]. Articulation and coordination work  
278 have been shown to be critical in data curation for multiple reasons. They are needed in maintaining a knowledge  
279 infrastructure’s stability [Karasti and Baker 2004]; facilitating the selection and enactment of data curation protocols  
280 [Darch et al. 2020]; refactoring data structures and vocabularies [Thomer et al. 2018a]; supporting infrastructure design  
281 [Baker and Millerand 2007]; enabling navigation of information during the process of scientific discovery [Palmer 2006;  
282 Palmer et al. 2007]; and are a core component of the “data labours” of building and sustaining data collections [Nadim  
283 2016]. Erickson and Jarrahi [2016] describe the articulation work needed by knowledge workers, such as data curators,  
284 to configure infrastructural solutions to overcome technical and contextual constraints in tools and workplaces. A  
285 recurrent theme in these papers is the lack of tools to support this coordination and articulation work; curators must  
286 coordinate their work often in spite of these tools, rather than through them.  
287  
288  
289  
290

### 291 3 METHODS

292 In this paper, we report on a mixed methods study to examine different aspects of the data curation process. We leverage  
293 two bodies of data: 1) semi-structured interviews with stakeholders across ICPSR; and 2) records of curation work in  
294 Jira tickets, a subset of the internal ICPSR documentation that records data curators’ work.  
295  
296  
297

#### 298 3.1 Research Site

299 ICPSR, founded in 1962, is one of the oldest and largest curated social science data archives in the world. It not only  
300 collects, curates, and disseminates data in a broad range of disciplines including political science, sociology, demography,  
301 education, criminology, public health, among others, but it is also a leader in repository infrastructure, data curation  
302 standard setting, and innovation in data curation. ICPSR’s archives include over 16,000 studies containing nearly 6  
303 million variables. ICPSR’s collections are organized into separate archives representing different subject areas and often  
304 sponsored by federal agencies and foundations. We selected ICPSR for three main reasons: 1) curation processes are  
305 well articulated and documented; 2) the volume of data curation is large enough for patterns to emerge, and 3) we were  
306 given access to both documentation and staff to conduct an in-depth study of the curation process.  
307  
308  
309

310 ICPSR’s organizational structure also makes it possible to study data curation in depth. Several years ago, ICPSR  
311 centralized curation into one unit. Curators previously worked for individual archives within ICPSR, reporting to  
312



313 a project manager, who, in turn, reported to an archive director; now curation staff, project management staff, and  
314 archive directors sit within their own distinct organizational units (e.g. the Curation unit, the Project Management  
315 unit). Part of this re-organization also involved a redesign of curatorial standards. As of 2018, datasets are assigned to  
316 one of three standard "level"s of curation which articulate specific curatorial actions that vary according to the amount,  
317 intensiveness, and complexity of effort required as well as the end product delivered. These levels provide a standard for  
318 curation actions and expected outputs, which are assigned based on the format, size, and level of preparation performed  
319 by the data creator prior to deposit [ICPSR 2020]. Higher levels of curation are intended to improve the usability of data  
320 products. All data deposited with ICPSR receive a base level of curation ("Level 1 Curation"), meaning that curators  
321 remediate personally identifiable (disclosive) information and create a metadata record, a Digital Object Identifier (DOI),  
322 statistical files, a webpage, and a codebook explaining the variables in the data collection. "Level 2 Curation" includes  
323 all "Level 1" actions, plus additional data transformations, completeness checks, and preparation of the data for online  
324 analysis. "Level 3 Curation" is intensive and includes custom documentation, attaching survey question text to variables,  
325 and indexing variables for search. Non-tabular data, such as qualitative or spatial data, typically require "Level 3"  
326 curation. For example, the "TransPop, United States, 2016-2018" study shown in Figure 1 is curated at Level 3, meaning  
327 that additional curatorial tasks have been assigned and more time has been budgeted for intensive curation including  
328 extensive disclosure review and remediation and creating searchable question text.  
329  
330  
331  
332  
333

334 *3.1.1 Interviewees and semi-structured interviews.* The internal stakeholders in ICPSR curation extend beyond the  
335 curation unit itself. In order to better understand the impact of curation within the data repository, we conducted  
336 in-depth, semi-structured interviews with 37 ICPSR stakeholders comprising six staff groups: archive directors, project  
337 managers, curation supervisors, curators, user support, and bibliographers. Each archive is led by a director who  
338 spearheads collection development efforts, secures funding, interacts with archive sponsors, and attends disciplinary  
339 conferences and meetings to expand the reach of the archive. When ICPSR ingests data, a project manager shepherds  
340 the data through curation and dissemination, serving as a conduit between the curators and the data producers. The  
341 project manager works with the archive director and the curation supervisor to determine which curation activities to  
342 apply to the data and how to prioritize the data relative to other studies in the queue. User support personnel bridge  
343 between data reusers and either project managers or curators as questions about the data arise. Bibliographers track  
344 use of the all ICPSR's curated datasets and maintain an extensive bibliography of that use.  
345  
346  
347

348 Curation work is accomplished primarily by a dedicated team of curators ( $n = 32$ ) and their curation supervisors  
349 ( $n = 5$ ). We note that data curators are typically entry level employees; this is not always the case at data archives.  
350 ICPSR requires curators to have experience with statistical software (e.g., SPSS, Stata), data preparation, and social  
351 science research methods. ICPSR actively curates data to ensure that they comply with the FAIR principles (i.e., are  
352 findable, accessible, interoperable, and reusable) [Wilkinson et al. 2016]. Generally, curators review data for sensitivity  
353 and re-identification risk, generate metadata [Vardigan et al. 2008], identify missing values, index variables for future  
354 search and discovery, link question text to variables, apply subject terms to the study, and generate multiple formats  
355 (e.g., SPSS, Stata, plain text) of the data files. Curators also pass along citations to publications that use the data to  
356 the bibliographers for inclusion in the ICPSR Bibliography of Data-Related Literature. These tasks are completed by  
357 individual curators and reviewed by their supervisor or a senior curator before disseminating the data for reuse. On an  
358 ongoing basis, the bibliographers also search for additional citations to studies archived at ICPSR.  
359  
360

361 The interviewees were selected using purposive sampling [Miles et al. 2014]; we requested interviews with all  
362 personnel working in the specific roles identified. The only criteria used to filter out potential respondents was for  
363  
364

Table 1. Number of interviews by stakeholder category

Stakeholder Category	Number of Interviews	Interview Codes
Archive Director	7	AD002-AD006, AD008, AD011
Project Manager	9	PM025-PM033
Curation Supervisor	7	CS001, CS007, CS009, CS010, CS012-CS014
Curator	10	CU015-CU024
Bibliography Team	3	BT034-BT036
User Support	1	US037

the curators themselves: due to the time required to become familiar with curatorial work, we limited our interview requests to those curators who had been working for at least one year so that they had built up some expertise in curation work and could better reflect on the processes. The interviews focused on understanding how the different stakeholder groups measured the value and impact of curation work (see Appendix A for our full interview protocol).

We conducted 37 semi-structured interviews with archive staff that enabled us to probe further into the responses and ask questions that were specific to each role [Dearnley 2005; Hesse-Biber and Leavy 2005; Rubin and Rubin 2012]. Interviews were conducted in 2019 and 2020; 12 were face-to-face interviews conducted before the COVID-19 pandemic led our institution (anonymized for review) to transition to remote work, and the remaining 25 were conducted remotely using the Zoom and Google Meet platforms. Table 1 details our interview participants. The interviews were recorded, and then transcribed by the REV transcription service and verified. We anonymized our interview transcripts by assigning all participants identifiers. Our study was reviewed by our university’s Institutional Review Board and found to be exempt from on-going oversight IRB information anonymized for review].

We began analysis using a deductive approach, using a "start list" of codes derived from our research questions, interview questions, and our knowledge of prior literature. In this first round of codes, we paid particular attention to identifying curatorial actions at ICPSR. Codes were iteratively expanded and refined through subsequent rounds of inductive coding [LeCompte and Schensul 2012; Miles et al. 2014]. Analysis of the interviews was completed using the qualitative data analysis program NVivo. Because multiple team members were conducting the coding, we established inter-rater reliability (IRR) to endure coherence. As we began establishing IRR between two members of the interview team, we realized that the interviews between the different stakeholders were divergent enough that IRR would need to be established within each stakeholder group. With the exception of the single User Support interview (59.2%) and the three Bibliography Team transcripts (69.5%), IRR was repeated within each stakeholder group until at least 70% was achieved using Scott’s pi [Scott 1955]. One member of the interview team coded transcripts across all stakeholder groups, and two different members established IRR with her on specific sets of transcripts.

After each round of coding was completed, team members reviewed coded data, then met as a group to discuss emergent themes. After our first round, we identified the role of craft and coordination as being key to data curation work, and decided to conduct a secondary round of axial coding to deepen our analyses. Again, IRR was established between team members until at least 70% was achieved using Scott’s pi. The authors reviewed coded data and again met to discuss the codes as a group. Finally, after reviewer feedback, we conducted a third round of coding, this time diving deeper only into codes related to craft in data curation to deepen our analysis. We met again as a group to discuss emergent themes.



417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442

Jira Dashboards Projects Issues Boards Timetracker Create Search

IRIS / IRIS-2595  
TransPop, United States, 2016-2018

Edit Comment Assign Log work More On Hold - ICPSR On Hold - PI Workflow

Details  
Type: Curation Request Status: CURATION COMPL... (View Workflow)  
Priority: Highest Resolution: Done  
Labels: RCMD

Managers Curators  
Project: DSDR  
Short Code(s): 000236  
Curation Level: Level 3  
Curation Tasks: Initial Review and Plan, DRR, Curation (ph file work), Metadata, Qtext, SDA, Documentation, Make DDI public for SSVd, ... (2)  
Required:  
Project Type: Add Study  
Access Level: Combination  
Deposits: d39468 d39791  
Number of Datasets: 3  
Number of Variables: 1,681

Description  
TransPop is the study title. I attempted to add proper dates and geography. Feel free to adjust as necessary to fit standards and conventions. I put United States in brackets because it is not part of the original title.  
1,681 variables is the total across all three datasets.  
TransPop == 566  
Cisgender == 503  
Combined == 612  
Cannot think of other relevant information at this time. I know I am missing something. Please you borrow from pages 4-5 (Data Sources) of the survey methods document and explicitly mention the data sources. Also, This whole document can be released as part study documentation, but page 39 has a section on applying the sample weights. Please include these instructions in the weights field.

People  
Assignee:  
Reporter:  
Reviewers:  
Votes:  
Watchers:

Dates  
Created:  
Updated:  
Resolved:  
Date Assigned:  
Date Assigned for 1QC:  
Date Assigned for 2QC:  
Date of Processing Plan Approval by Supervisor:  
Date Processing Plan Submitted to Senior Curator:  
Date of Processing Plan Approval by Senior Curator:

Time Tracking

Fig. 1. Jira ticket for a single study

443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460

3.1.2 *Triangulating with Jira tickets.* We triangulated findings from interviews with documentation created through the curation process, again looking for descriptions of curatorial actions. Curation work at ICPSR is coordinated and documented across three main sets of documents: processing plans, Jira tickets, and processing history (PH) files. Jira tickets are the richest and most specific record of data curators' work. Jira is type of project management software that organizes work through the creation of "tickets" that describe the work that needs doing, and that users can update with progress over time. When data are deposited through the ICPSR deposit system, staff review the data for fit and priority, and a data project manager or assistant generates a Jira ticket (see Figure 1; we removed identifying information from the fields on the right but leave their titles to show what information tickets contain). They provide a study title, the priority of the study, the funder or sponsoring archive, a description of the work curation will need to do, and the level of curation (and any additional tasks) required. Before curation begins work on a Jira ticket, metadata unit staff review the ticket and study metadata. After data project staff and metadata staff approve the ticket, it gets sent to curation for assignment. While curation works on the study, they provide details about their work and progress in the "worklog" section of the ticket (see Figure 2). The worklogs offer insights into aggregate time spent on different kinds of curatorial actions at ICPSR.

461  
462  
463  
464  
465

To classify the parts of worklog descriptions (e.g., "Began curation" and "Metadata and proc plan" in the example in Figure 2), we developed a set of eight high-level curatorial actions that describe curation work: initial review and planning; data transformation; metadata; documentation; quality checks; communication; non-curation; and other activities (see Figure 3). These categories mirrored the codes used in our qualitative analysis.

466  
467  
468

We manually coded a randomly selected proportional sample of Jira ticket worklog entries stratified by curation level. These were coded in *brat* software [Stenetorp et al. 2012] to create labeled training data to facilitate the automatic

Activity				
All	Comments	Work Log	History	Activity
Curator		logged work -		Date/time
Time Spent:		3 hours		Began curation
Curator		logged work -		Date/time
Time Spent:		30 minutes		Metadata
Curator		logged work -		Date/time
Time Spent:		4 hours, 30 minutes		Metadata
Curator		logged work -		Date/time
Time Spent:		1 day		Metadata and proc plan

Fig. 2. Worklog excerpt from a Jira ticket

classification of the Jira ticket worklogs (discussed more fully in (anonymized for review)). We trained a computational model with 0.75 accuracy to assign each worklog entry one of the eight categories of curatorial actions (summarized in Figure 3). For example, a worklog entry “Discussed curation standards with supervisor (2 hours)” is classified as an instance of “Communication” while an entry describing “Recording dataset limitations in processing notes (10 hours)” is classified as “Documentation”. We then aggregated each class of action to analyze the relative amount of time spent on each.

There are several limitations to this study. First, it documents curation work at one repository. Second, ICPSR is a mature repository with well-articulated policies and procedures. Finally, we did not do direct observation of the curatorial process but rely on direct reports from curation staff and other stakeholders and indirect observation through the Jira tickets.

#### 4 FINDINGS: CURATORIAL WORK AT ICPSR

In the interviews and Jira tickets, we found a consistent, overlapping vocabulary of actions describing typical curation work. We also found insights into the ordering and time spent on curation tasks (see Figure 3 for examples of each type of action). Table 2 summarizes the amount of time curators logged for each type of curatorial action. While there were some typical sequences in which actions are performed, curators describe considerable variability in their own day-to-day work, rooted in their specific preferences, expertise, and craft knowledge.

In the subsections that follow, we first describe the core high-level curatorial actions that are undertaken at ICPSR. These expand prior accounts of the work of data curation, particularly in CSCW and data science, where it’s described as one small part of a process. Then, we describe how curators rely on their craft knowledge to navigate the “workflow” dictated by these actions. In doing so, *we show how best practices and craft practices are deeply intertwined.*

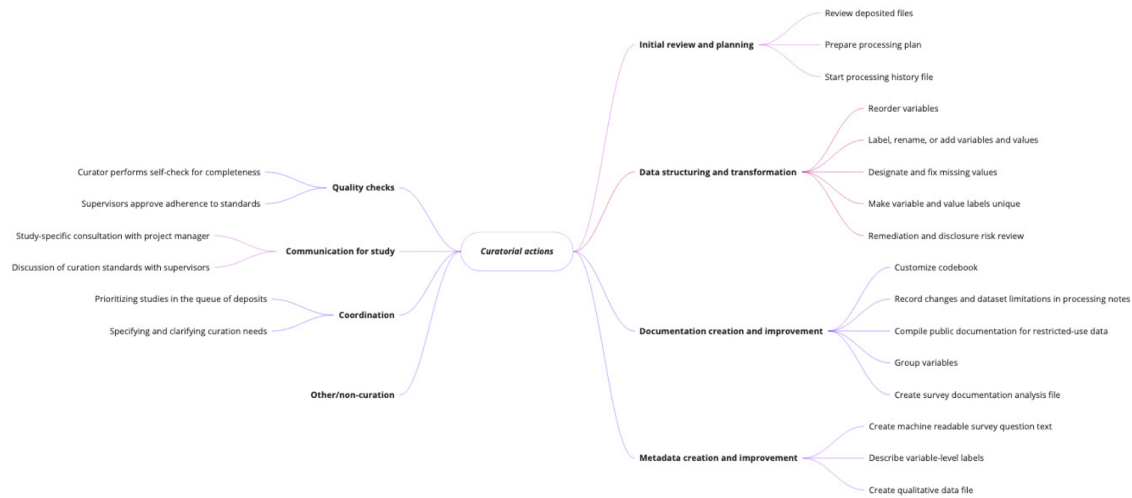


Fig. 3. High-level curatorial actions that occur throughout the curation process

Table 2. Time spent on curation actions (Feb. 2017 - Dec. 2019)

Action	Total hours logged	Percentage
Communication	3,249	7%
Data transformation	12,363	26%
Documentation	3,094	6%
Initial review and planning	5,778	12%
Metadata for study	2,669	6%
Non-curation	6,641	14%
Other	1,157	2%
Quality checks	13,075	27%

#### 4.1 High-level curatorial actions

*Initial review and planning.* Data curators at ICPSR typically begin their curation of a deposit by reviewing deposited files and metadata and *developing a processing plan* – an outline of planned curation tasks, depending on a dataset’s designated curation level. More detail about curation levels at ICPSR can be found in Section 3.1.2. These plans are developed by curators and reviewed by curation supervisors, who answer questions, troubleshoot, and generally advise along the way. In recent years, more initial review and planning actions have also been recorded for higher levels of curation (Curation Levels 2 and 3), suggesting that relatively more attention may be dedicated to developing curation plans at higher levels of curation. Initial review and planning accounted for 12% of curation time over our study period.

Early curation work also includes *disclosure risk review (DRR)*, in which curators evaluate the risk that publishing a dataset might pose to research participants and identify appropriate mitigation steps. Curators described DRR as a critical way in which they add value to a dataset – both because of the anonymization it provides, and for the thorough oversight the DRR represented.

573 *Data structuring and transformation.* This category of curatorial work includes the most “technical piece” of curation  
574 [CU015]: the direct work with the dataset itself to make it easier to use, share and archive. Data transformation tasks  
575 include designating missing values (e.g., assigning metadata to values such as “no response” and “not asked”); adding  
576 question text (inserting the survey questions verbatim); transforming curated SPSS data files into other statistical  
577 packages (e.g., R, SAS, etc); and creating documentation (PDF codebooks and XML metadata files). Datasets sometimes  
578 are split into multiple, more usable parts (for instance, into smaller file sizes, or commonly used file formats), or merged  
579 into single files from multiple sources. This data structuring entails more than just mechanical reformatting; as one  
580 curator describes, “a lot of times, especially in the larger datasets, there’s a lot of pieces to put together that I think  
581 when we make those connections it makes it easier for users to use the data.” [CU023] Considerable expertise and  
582 judgement are needed to structure and transform data well. Data transformation was comparatively time consuming,  
583 taking up 26% of curation time over our study period.

587 *Metadata creation and improvement.* Curatorial work includes the creation of records that will be queried by users  
588 within ICPSR’s online repository. Metadata development is seen as a distinct task; where the focus of data structuring  
589 is on making the dataset usable in and of itself, the focus on creating metadata is on supporting search and retrieval  
590 of datasets. Curators saw metadata as particularly important because it’s “the first line” of access, [CU017] the first  
591 thing users see. Metadata improvements include drafting or revising a dataset’s description; copying and refining  
592 metadata from the initial data provider, such as data collection dates; creating question text (i.e., writing out the full  
593 list of questions in the survey instrument that generated the data); and defining variable-level labels (i.e., creating a  
594 data dictionary that spells out what each data variable represents). Metadata work accounted for 6% of curation time.  
595 This work is both qualitative and technical; curators must have skill manipulating metadata standards to create these  
596 records, and they need to have the experience to understand what context is necessary and helpful for data reusers to  
597 include in metadata records.

600 *Documentation creation and improvement.* In addition to creating metadata records, curators also develop other  
601 forms of documentation about the datasets. This includes creating processing history files, codebooks, which include  
602 information for each variable in a dataset, documentation of major changes made to the data, and compilation of any  
603 additional documentation archived by the data producer. Documentation accounted for 6% of curation time logged.  
604 More instances of documentation were recorded for curation activities in topical archives than in the general archive.

607 *Quality checks (QC).* These include checking data and metadata files for completeness, confirming that the work done  
608 to a dataset aligned with the Jira request, comparing the work done to the processing plan, and confirming adherence to  
609 ICPSR’s guidelines and protocols for curation. The vast majority of studies include quality checks, which was the major  
610 category of curation action we detected in our analysis of Jira ticket worklogs, accounting for over 27% of curation time  
611 logged. These quality checks are performed by a second curator to provide an extra level of review.

613 Beyond designated quality checks, stakeholders discussed the value that curation provided in ensuring that the  
614 data was of high quality overall. Project-related communication is one mechanism for ensuring high quality curation,  
615 which accounted for about 8% of curation time logged in Jira tickets. This includes catching issues and addressing  
616 complicated challenges with data that data producers did not, providing high quality documentation about data and the  
617 data curation process, setting and meeting goals for data release that match depositor expectations and deadlines, and  
618 being a source of consistent, vetted data. A curator described the last item in this way:

621 Our work, I think, is pretty impactful and benefits the community because the work we put in [will] rule  
622 out all the troubleshooting. We look at the data, compare it to the documentation, and then do these things  
623

625 to make sure everything's consistent. If there's any problems, we either resolve with the PI or we have  
626 our own solution for it, so once you get your hands on the data, there's nothing really in question for the  
627 most part. [CU018]  
628

## 629 4.2 Using craft to work the curatorial workflow

630 The previous section outlined the tasks involved in the technical work of data curation. Actually accomplishing  
631 that work, however involves craft. Specifically, curators organize their work by first developing a gestalt, abstract  
632 mental representation of the data to envision what the final released dataset will entail; they then use their judgement  
633 and expertise to interpret standards, creatively come up with solutions, and thereby achieve a standard outcome in  
634 unstandardized ways. Paraphrasing one curatorial supervisor, "I'm the curator: I do anything needed to make the data  
635 archivable" [CU015]. We describe these two aspects of curatorial craft practices in the following subsections.  
636  
637  
638

639 *4.2.1 From abstract representations to fit-for-use.* Curators approach a new dataset by getting the gestalt of the dataset:  
640 understanding the whole of the dataset as beyond the sum of its parts, as well as how these fit together in order to  
641 assess the feasibility of the processing plan and to envision the archivable and disseminated dataset. Several curators  
642 described getting the gestalt:  
643

644 "I don't ... I do find the plan useful and it has to be done and there are things that it walks you through  
645 that can help you find missing things or problematic things. But I also rely heavily on just running  
646 the frequency output on the data and just scrolling through it. ... But I have found, many times I have  
647 found things there that I wouldn't have seen otherwise. So I find that very important. And then I just  
648 ... yeah, between the plan and this check on my own, I find out what I need to do with the data and the  
649 documentation ... I'm not always very linear in how I work with this. I tend to do a lot of poking around..."  
650 (CU017)  
651

652 "So myself personally, I'll try to work with the data first, make some data manipulations or changes after  
653 I've gone through and read the documentation that's been provided, and get a good sense of what's going  
654 on with the study. But I won't... So I'll read through everything and I won't fill out the metadata at that  
655 point because I still like to be able to go through and actually work with the data before I fill out the  
656 metadata that explains the collection in some more detail." [CU023]  
657  
658  
659

660 Much of the process of getting a sense of the data is done with the user in mind. Curators think about how the  
661 dataset would need to be structured or documented in order to be fit-for-use for a range of users. Multiple curators  
662 describe customizing their work, or making decisions with the goal of supporting users' access, essentially envisioning  
663 themselves as the "first user" of a dataset, and "trying to figure out everything that a potential user would want to know  
664 and to make sure the archive version that we release is as complete as possible" [CS009]. The curators' goals are to  
665 answer any questions future users might have about the data: "We obviously can't anticipate everything, but we try to  
666 say, okay, if I was just picking this up, what would I need to know about it that maybe I don't have in a quick glance?"  
667 [CS010]. They also want to let users know about known issues "so they don't have to dig through it themselves to  
668 figure it out" [CS010].  
669

670 Curators anticipated other types of users and user questions. For example, one described tailoring datasets to a range  
671 of users: "Our data needs to be easy enough to use for the most novice user, but sophisticated for the more advanced  
672 user as well. And I think that that can happen by doing the details, making it easier" [CU017]. A curatorial supervisor  
673 considered how the dataset could be represented through crafting good metadata in supporting of use:  
674  
675  
676

677 "And then even in our metadata. So I recently was talking to someone because I could tell them that the  
678 metadata ... Like the way they have worded was too internal-facing. I said, "That's not going to mean  
679 anything to users." If they don't understand, they're not going to hear about it. So we need to make it in a  
680 way that it's going to be something that makes sense to them. And is useful for them." [CS014]  
681

682 And a third curator specifically linked the craft and subjectivity of curation with being able to conceive of how different  
683 users might perceive a dataset.  
684

685 "I look at curation as, you know those technical aspects, there's do's and there's don't's, right and wrong.  
686 There's also some, I like to say artistry, subjectiveness do it and how I might perceive a group of people  
687 wanting to view the data. Another person might see it differently." [CU016]  
688

689 The gestalt techniques used by curators manifested itself differently in different cases. CU017 (above) expressed  
690 creativity in the information they chose to highlight for users. Several curators described differences in the order in  
691 which they approached curation tasks for a data study, or the time or level of detail they devoted to certain activities  
692 over others such as developing the processing plan. This ability to assess the present data, conceive of the path to a  
693 future state and conceive of a future representation is one aspect of craft exhibited by the curators.  
694  
695

696 *4.2.2 Achieving standardization through judgement.* Standards play a large role in data curation: at ICPSR, these include  
697 internal "house" standards set by ICPSR itself, and external standards developed by the broader community, such as  
698 metadata standards or preservation best practices. However, the application of standards is far from rote. Curators use  
699 their expertise and make judgements about when, how, and why to apply standards throughout the curation process.  
700 Several participants said that because the data that ICPSR receives is just too diverse for strict standards to be feasible,  
701 and curators must rely on their craft knowledge and expertise to navigate "gray areas" [CS012]. This can happen with  
702 "unusual data sets" in unique formats or with idiosyncratic structures [PM028], or for particular archives with distinct  
703 user communities, or for instances where a PI has requested what one participant called "a la carte" curation, where they  
704 do everything from one level, plus one task from another [CS009]. One supervisor said there are multiple workflows  
705 that stem from agreements with PIs, and therefore, one singular workflow isn't possible:  
706  
707  
708

709 "I do think that we will always have more than one workflow just because sometimes the way proposals  
710 have to be written ... Sometimes project officers have a certain thing in mind. So, I think I mentioned  
711 earlier some PIs want a lot of involvement in disclosure review or the changes that we're making. So,  
712 for the demography archive, for pretty much most, if not all studies, we basically list out the examples  
713 of things that we're changing, and send it to the PI for approval. So, that's a different work flow than  
714 normally we just kind of can proceed. And then there's an archive within the criminal justice archive  
715 where the analysts want to use the data to do their analysis and then publish reports before we release the  
716 data. So, we have a process there were we do the first quality check then send it to them for review. They  
717 might send changes back to us, they might spend the next six months analyzing data. And so, we may not  
718 return to that study to make changes and/or do the second quality check and release it until months, or  
719 even a year later." [CS001]  
720  
721  
722

723 Some curators expressed frustration with standards. Long time curators particularly viewed the standards as living,  
724 malleable tools that change over time, rather than as unbreakable rules. Some went so far as to say that standards could  
725 be an obstacle to their work, because they interfered with their preferred way or working. One curator felt that "... some  
726 of the standards that we have get in the way of it when they're supposed to help it." [023] A long-time curator observed:  
727  
728



729 They're getting better. Previously, there was just a lot of unanswered questions in the document. A lot  
730 of ambiguity. ... Sometimes I'm a little outspoken just because I've been around and... We call them  
731 standards." [CU017]  
732

733 This curator went on to say:

734 "I don't necessarily agree that they're standards. They're somebody's opinion that got put down and then  
735 set as a standard. It's somebody's preference then they made it a standard, especially on sentence case  
736 for variable labels that one's like... Those aren't standards. That's somebody's preference and I don't like  
737 it, because I would just do it differently. Not that my way is right and they're wrong. It's just different  
738 ways of thinking. If I was the one creating that standard it'd be different because my preference and my  
739 thought process is different." [CU017]  
740  
741

742 Several curators discussed applying standards creatively and flexibly in service to users, extending the curator's  
743 focus on the user discussed previously. One curation supervisor acknowledged the importance of creativity to sidestep  
744 standards when the user was not served:  
745

746 "So yeah, keeping those standards in mind, it's just sometimes you have to be creative. If there's something  
747 that you know users need to know who, put it in the summary field, don't put it in the collection notes,  
748 make it more visible. Maybe your tools don't allow us to make it as visible as we'd like. But there's always  
749 a way." [007]  
750

751 Curators also were aware that the applying the standards involved judgement calls and could be self-reflective on  
752 their comfort level in making some types of these calls:  
753

754 "... so there's a lot of judgement calls involved despite all of our efforts to write up standards and follow  
755 them, again it's human produced data and documentation. And there're still judgement calls to be made  
756 from grammar and spelling and capitalization to more serious matters. So there's a number of judgement  
757 calls that I feel comfortable making. such as those involving labels. But then on the other hand, if it's a  
758 really sort of hairy disclosure risk scenario, I would definitely check in with my supervisor on those."  
759 [MICA 24]  
760  
761

762 In deciding to apply or not apply standards, curators use expert judgement, creativity, and skill to achieve a standard  
763 outcome. One curation supervisor described their work as helping maintain the standards among curators but not setting  
764 or enforcing them. [012] Curators were also self-reflective about the standards and considered the data themselves as  
765 well as potential users when arriving at solutions that fell outside a 'normal' application of the standards. Whether the  
766 curators were more respectful or skeptical of the standards, many discussed instances where the standards fell short of  
767 achieving a dataset fit-for-use.  
768  
769

770 4.2.3 *Organizing curatorial actions.* Though there is a *common* sequence of curatorial actions at ICPSR, there is not a  
771 strict workflow; processing plans outline the work that's needed at a high level, but not how it should be carried out.  
772 Curators use craft knowledge to sequence their technical work, and to customize their work practices to the dataset at  
773 hand or their own preferences. As one supervisor described,  
774

775 "We always say that curation isn't a linear process. There are a set of tasks that it makes sense to do in a  
776 particular order sometimes... I mean, we try to leave it up to the curator to what it works best for them  
777 because everyone has different ways of curating so... some people like to do metadata first, some people  
778 do that last, some people want to make all these data edits right away, some people want to focus on  
779  
780

781 peripheral stuff. We try to get the processing plan done as soon as possible just because that helps expose  
782 all of the other issues that we might need to go to the project manager or the PI about. And that gives us  
783 more information about if we need to prioritize something particular." [CS013]  
784

785 The curators themselves described considerable variability in how they ordered their worked:  
786

787 "I'm very collective, and it's not always the case, but just as I would prefer to be working on multiple  
788 curation projects at a time and instead of just focusing on one all day, every day until it's done, I like to  
789 have two or three going, if possible, just to break it up, break up my day, break up my focus. I also have  
790 that same approach for the tasks. So I might jump between different things." [CU016]  
791

792 "Well the prioritization is to complete the plan and the disclosure risk worksheet. So that is where I start.  
793 And in completing those, I set the agenda, so to speak, for where it goes next." [CU024]  
794

795 "I tend to start with the data, I tend to leave metadata to the end because I often find ... it could go either  
796 way, right? You could do the metadata to inform how you approach the data, but I often find that going  
797 through the data, going through the questionnaire lets me fill in the metadata better. Yeah. So my first step  
798 is the plan and the worksheet, because that is the first part of the process, and there's checks involved  
799 with other people. And from there I usually tackle the data first. [CU024]  
800

801 "I definitely jump around." [CU015]  
802

803 The common thread throughout these different approaches is that curators draw on their own expertise to structure  
804 their work and days: they know what works best for them and how best to hone their attention for detailed, technical,  
805 and sometimes tedious work.  
806  
807

### 808 **4.3 Coordination in service of curation** 809

810 Above, we described how curators use craft practices to gain a gestalt understanding of a dataset's structure, and  
811 then to organize their own work. This work is not done in a vacuum, however; curators must also coordinate with  
812 other stakeholders at ICPSR to proceed with this work, and to clarify priorities. Because ICPSR's workflows resist  
813 standardization, curators and curation supervisors consult archive directors, project managers, data producers, and  
814 each other to ensure there is a consensus (if not agreement) on the best way to approach curating a study. This occurs  
815 throughout the curation process. For example, coordination occurs early on to determine where a study is placed in the  
816 curation queue and identify the level of curation. It can also occur later in the curation process if issues emerge requiring  
817 a decision about additional curation activities which are required to make the study fit for use. Acts of communication  
818 are captured in the Jira ticket worklogs, but the content is often vague. Our interviews elucidated the frequency and  
819 critical place of coordination in the curation process. These include the following: prioritizing studies in the queue;  
820 specifying how data will be curated; and monitoring progress and alerts.  
821  
822  
823

824 *4.3.1 Prioritizing studies in the queue of deposits.* Curation supervisors manage a large queue of studies waiting to be  
825 processed and assess how, when, and to whom to assign them based on the priorities and funding available to the various  
826 topical archives at ICPSR. This assessment includes tight coordination with archive directors, data producers, project  
827 sponsors, and project managers. Curation supervisors factor in a project's budget, promised deliverables, relevant  
828 external deadlines, and the potential impact of the study's release to determine placement in the queue. Two archive  
829 directors described this balance of considerations:  
830  
831  
832

833 ...our funder really decides what to archive. I work with our project manager to [...] ensure that the  
834 curation team is prioritizing our data the way we want it prioritized. [...] And the project manager ensures  
835 that those [...] priorities get communicated to the curation team. [AD004]  
836

837 ...I would say that feedback from the funders influences both the curation levels and the priorities that  
838 we give to studies. So we do coordinate with our program officer. [...] If there are certain studies that  
839 are a high priority and that is something that we would then incorporate into Jira and into the curation  
840 requests so that we can adjust priorities and make sure that the highest priority work gets prioritized  
841 accordingly [AD029].  
842

843 Project managers also communicate priorities to the curation unit. They do this as a matter of routine through  
844 multiple, reinforcing channels: project managers enter deadlines and rank relative priority in Jira tickets (e.g. Highest,  
845 High, Medium, Low), and they hold quarterly meetings with curation supervisors to “talk about the queue for a  
846 particular quarter” (PM025). However, as several curation and project management staff noted, priorities change, and  
847 the communication often involves significant back and forth:  
848  
849

850 There’s a lot of back and forth in terms of what their priorities are versus what we feel we can reasonably  
851 accomplish in a given timeframe. And so sometimes that can get a little tricky. So in terms of like, if they  
852 say we have this and this, we’re going to ask them which one is more important to them? I’m not going to  
853 try to figure that out, if I can assign them both I will, if I can’t I’ll make sure it’s their highest priority.  
854 [CS010]  
855  
856

857 *4.3.2 Negotiating levels of curation.* Though the project managers initially define the work expected by choosing a  
858 curation level (1, 2, or 3), curators sometimes find that a given dataset needs more or different curation than originally  
859 planned. When this happens, they (and curation supervisors) must negotiate up and down the organizational chart to  
860 come to an agreement about how curation will proceed. Curators and curation supervisors negotiate with  
861 project managers, who coordinate with archive directors, who sometimes coordinate with PIs. A curation supervisor  
862 described the negotiation process that can be involved in curation level clarification:  
863  
864

865 [...] We have curation levels and the project manager reviews those levels and says, “Okay, I want this level  
866 of curation.” And then we would review it to see if that’s accurate. So that would be like a collaboration  
867 between the supervisor and the curator. So when [the curators] do their processing plan, they may identify,  
868 “Hey, they’re asking me to do something that isn’t in this level.” Or “they’re asking me to do things but it  
869 should be like a level up or down.” And then we also do a review of the plan and then we assess as well.  
870 [CS014]  
871  
872

873 We note here that the project managers may not necessarily get the same gestalt view of the data as the curators, they  
874 trust the curators’ view in further structuring work. As curation proceeds, and curators get into the data, they often  
875 discover things that suggest several possible courses of action that can prompt discussion with curation supervisors,  
876 project managers, archive directors, and the data providers themselves. One curator described working with their  
877 supervisor to make final decisions. Though the curator ultimately defers to the curation supervisor, there’s still a  
878 conversation about potential options:  
879

880 If we’ve identified an issue or something, [...] I might give [the supervisor] some options and then we  
881 talk about it a minute, and ultimately [...] I let her decide as the supervisor, especially when it comes to  
882 things on how to address confidentiality things. Those are definitely things that supervisors would like to  
883  
884

885 have their approval on before things get out. It lessens the responsibility in a way on us by having that  
886 supervisor, or someone who's in charge, being able to make the final decisions [...] It's good to have, I  
887 think, other people's opinions on those things. Yeah, for my supervisor, it's definitely a conversation of,  
888 "Here's some possibilities of what we could do." [CU015]  
889

890 Thus, there is some tension between respect for the curators' expertise and deep knowledge of their data, ICPSR's  
891 standards and decision-making hierarchy, and the overall budget for a project.  
892

893 *4.3.3 Monitoring progress via alerts, and navigating varying degrees of visibility.* One mildly controversial method  
894 of facilitating coordination is through the use of Jira tickets to monitor and record progress on a project. Curators,  
895 supervisors, and project managers communicate about the study via Jira ticket comments. Curators receive an alert  
896 every time tickets are updated, and the tickets act as a running log of the work performed on the study. Though project  
897 managers found Jira to be generally helpful, some curators characterized Jira as annoying or overwhelming. The deluge  
898 of alerts and documentation also made some curators feel micromanaged, and as if Jira was keeping a running log of  
899 their work for their supervisors to review at any moment. As with any representation of work, however, the Jira tickets  
900 can be more or less precise. As one curator noted:  
901  
902

903 "So I have trouble ... occasionally I have trouble keeping up with the ticket, the JIRA ticket, where we're  
904 meant to tick off things as we go because I'm kind of doing a little bit of everything at once because every  
905 part has information you need that affects other parts. I often find myself quite close to the end and I'm  
906 like, "Oh shoot, I have to go update the ticket." [CU024]  
907  
908

909 The curators' feeling of sometimes being overly visible is mirrored by other concerns about curatorial work being  
910 under recognized, or that curatorial work was insufficiently visible. For instance, one project manager said that they  
911 felt curation skills were invisible to those who are more removed from it:  
912

913 I think that generally a lot of project managers and directors think that it's simple, simple syntax being  
914 applied but some of the challenges that come up while curating data can be quite complicated. It can take  
915 a certain level of skill. [PM026]  
916

917 The centralization of the curation function has taken the curators out of the individual archives, and the implemen-  
918 tation of a single, shared standard has altered their practice to align with the organization, rather than with a single  
919 archive within ICPSR. In this new arrangement, ICPSR staff interact with curators at discrete points in the curation  
920 process, but few interact throughout the process. Therefore, there is less opportunity to see how curators use complex  
921 representations to envision the data as fit-for use and how they use judgement and creativity to achieve standardized  
922 outcomes for data. More coordination via intermediaries – whether Jira or project managers – becomes necessary to  
923 support curatorial work.  
924  
925

## 926 **5 DISCUSSION**

### 927 **5.1 Understanding data work: more than just technical**

928 One of the motivations of this study was to create a finer-grained understanding of data work, specifically focusing on  
929 curatorial actions. We developed a rich description of data curation work at ICPSR – one that goes beyond the technical,  
930 procedural work with data and metadata to include the expertise-driven decision-making involved in crafting data,  
931 and the coordination required to develop a consensus around curatorial priorities and activities. Our participants have  
932 shown that data curation is neither "magic" nor "janitorial" work [D'Ignazio and Klein 2020; Owens 2018; Rawson  
933  
934  
935  
936

937 and Muñoz 2016], but rather, is the result of technical skill enacted through craft practices. Indeed, we find that staff  
938 members in all roles bristle against characterizations of curation as something rote or mechanical. Curators do what  
939 what needs to be done to achieve the outcomes of a standard – even when not necessarily *following* a standardized  
940 workflow. This requires significant collaboration with other stakeholders in the data science workflow. Thus, the  
941 workflow is achieved but much of the actual work that made that happen disappears.

942 Our work makes two main contributions to understanding data curation, and thereby data work. First, the description  
943 of “hands-on” technical tasks we provide in Section 4.1 expands an existing body of literature describing data curation  
944 practices in different contexts. Understanding different data (curation) practices is critical for building infrastructure,  
945 software tools, and ontologies that capture disciplinary contexts, and for educating curators. For instance, Chao et al.  
946 [2015] developed the Data Practices and Curation Vocabulary, which describes how a community (in that case, earth  
947 scientists) defines data curation. Comparison of our two frameworks reveals that ICPSR has much more detailed quality  
948 check protocols, and ICPSR’s curators spend considerable time on tasks like “adding question text” that simply are not  
949 needed in the earth science fields. The diversity of curatorial actions shown in just these two papers highlights the need  
950 for further research into the specific curatorial workflows and communication regimes in different scholarly settings.  
951 It is well understood from research on data practices that there are significant domain differences in curation needs  
952 [Akers and Doty 2013; Cragin et al. 2010; Faniel et al. 2019; Witt et al. 2009]. Yet models of data curation rarely account  
953 for this diversity of practice, or provide guidance in how to navigate them.

954 Second, our work shows the vital role that craft practices play in successfully organizing curatorial work and applying  
955 and navigating standards. *In data work, we see craft manifesting as the ability to develop an abstract, gestalt representation*  
956 *of a data product and then envision how to make changes to that data product so that it is more fit-for-use. This work*  
957 *involves following best practices and creating a standardized product, but not necessarily following a standardized workflow.*  
958 Furthermore, the kind of data curation carried out at ICPSR requires significant collaboration and consultation with  
959 other stakeholders. This extends prior work on craft in technical settings in the CSCW literature, most recently discussed  
960 by Muller et al. [2019] in their summary of craft in the context of data science, as well as a more recent focus on craft in  
961 the LIS literature by Owens [2018]. Muller and co-authors summarized key themes regarding craft in CSCW, including  
962 “Conversation with materials: Through the conversation with materials, there is often a sense of intimacy with materials  
963 and media” and “Control: Craft-workers labor at an intersection of control and unpredictability.”

964 At ICPSR, we see clear alignments with some of Muller et al.’s account of craft. Curators repeatedly emphasized the  
965 importance of the “conversation with materials” in their work through repeated descriptions of the contingency of their  
966 workflows and specific tasks. Likewise, ICPSR curators exist at the intersection of “control and unpredictability” – they  
967 are constrained into somewhat narrow roles by ICPSR’s organizational structure, yet must navigate unpredictable and  
968 unique curation challenges for each dataset with stakeholders throughout and outside of ICPSR. Our research further  
969 shows how craft practices “fit” into best practices and other standards for working with data; *in short, we find that craft*  
970 *practices are necessary to enact best practices.* It’s well understood that data standards can vary in their application and  
971 results based on variations in how they are enacted by a group Millerand and Bowker [[n.d.]]. Yet at ICPSR we see a  
972 *standardized result arising from the nonstandardized application of standards via craft practices.* By giving curators the  
973 freedom to rely on their own skill to structure their work and make decisions, ICPSR is able to truly rely on them as the  
974 human-in-the-loop.

975 Accounting for the role of craft and expertise in data work is important in designing effective data workflows,  
976 training data workers, and in better supporting data workers in showing the impacts of their work. We expand on  
977 this further in section 5.3 We argue that this view raises important questions for the practice of science (data science,  
978

989 social science, etc.), such as: How do notions of “craft” complicate the development of data curation pipelines to support  
990 complex data science applications or support repository infrastructures that automate curation? We begin to consider  
991 the latter question in the following section. How does understanding the craft involved in data work support data  
992 workers in gaining credit for their work and it’s impact? We address this question in section 5.3  
993

## 994 5.2 Coordinating work in data curation: complicating “workflow” or “pipeline” views of data science 995 and curation 996

997 One of our primary findings is that data curators must structure their own work within the context of their organization’s  
998 structure and job descriptions and constraints. In this way, they and other stakeholders “work the workflow” and  
999 navigate across standards and up and down the organizational chart; they gain a gestalt view of not just the data at hand  
1000 but also of the organization as a whole. Coordination and communication are key in this. In identifying coordination and  
1001 craft practices as important parts of data curation work, we complicate not just solely technical accounts of data curation,  
1002 but also “workflow” or “pipeline” conceptions of data work. By “workflow” views, we mean conceptualizations of data  
1003 curation as a sequential process, easily represented by a UML diagram or similar technique. These representations are  
1004 quite common in CSCW and the information sciences, where they are used to model curation processes at a higher  
1005 level [Johnston 2014; Kross and Guo 2021; Muller et al. 2019; Zhang et al. 2020], or the plethora of data/digital lifecycle  
1006 models in the digital curation literature [DataONE 2015; Faundeen et al. 2013; Higgins 2008]), or to capture detailed  
1007 change logs and provenance chain of a dataset [Goble et al. 2008, 2010; Thomer et al. 2018b; Zhao et al. 2012]. The  
1008 models are common because of their utility; they represent complex processes in a way that is digestible, and they can  
1009 act as boundary objects that help communication between disparate groups of stakeholders [Dourish 2001].  
1010

1011 However, our work here underscores that data curation is more than the sum of its parts, involving much more than  
1012 the objects that are curated; it is also a process in which distributed knowledge management decisions are made to  
1013 facilitate information reuse [Ackerman and Halverson 1999]. Our research supports the notion that data curation is  
1014 a highly collaborative process occurring across a distributed system over time. While some curation actions tend to  
1015 occur in sequence, important components of curation work, like quality checks, are performed in parallel or iteratively  
1016 throughout the curation process. Project-related communication is also embedded in all other curatorial actions, making  
1017 it difficult to delineate. A closer look at project-related communication reveals the importance of discussion and  
1018 delegation in curatorial work; for example, supervisors and curators often discuss how best to mitigate disclosive  
1019 variables on a case-by-case basis, following risk minimization heuristics rather than hard rules. And though coordination  
1020 strategies such as *Prioritizing studies in the queue of deposits* and *Specifying and clarifying how the data will be curated*  
1021 may seem like they could fit neatly into a workflow diagram, in reality, they require a meta-level understanding of  
1022 the curation workflow itself to proceed. Articulation work is needed to navigate a data science workflow [Neang et al.  
1023 2021; Thomer et al. 2018a], yet this labor can, somewhat ironically, be obscured in workflow-centric views. We want  
1024 to be clear: we are not trying to discourage or dismiss workflow-based explorations of data work. Rather, we want  
1025 to note the importance of continued, rich exploration of what goes on in and around each “box” of the diagram -- lest we  
1026 obscure that which we wish to reveal.  
1027

## 1028 5.3 Revisiting the invisible nature of data curation 1029

1030 In our prior work, we have argued that hiding curation makes it harder to plan, prioritize and fund curatorial activities  
1031 (anonymized for review). Additionally, by rendering their work invisible to outsiders, curation can hide curators’ value  
1032 and impact. Our interviews verified this latter point, in that curators – and even their managers, to a degree – described  
1033



1041 some concern that their work was not truly seen or appreciated. Invisibility can make it harder for these data workers  
1042 to advance in their careers, lobby for salary increases, and participate fully in their fields. Our work here reveals some  
1043 tensions, though, in making curatorial work totally visible. Below we discuss both the visibility and hypervisibility of  
1044 curatorial work at ICPSR.  
1045

1046 The craft and coordination in curatorial work at ICPSR are mostly invisible to data users. The public datasets hide  
1047 the work that went into their creation precisely because they are standardized [Plantin 2019]. Even the documents  
1048 that emerge from curation hide aspects of this work; while the Jira tickets contain descriptions of the high-level tasks,  
1049 they do not provide a full account of the curators' labor and decision making process. The existence of data curation  
1050 standards makes the work seem routine even though all our interviewees recognize, to varying degrees, that curation  
1051 requires technical skill, flexibility, and coordination.  
1052

1053 At the same time, some aspects of curators' work is hypervisible within ICPSR through Jira tickets and other  
1054 documentation. The Jira ticket worklogs and comments, especially, serve first to coordinate work and then to document  
1055 it. And while Jira can document their labor and decisions, making their work visible, it can also open curators to negative  
1056 side-effects such as micromanagement. For instance, Jira tickets make it possible for more powerful colleagues (e.g.,  
1057 archive directors) to monitor curators' work. As Suchman [1995] and Yates [1989] pointed out years ago, technologies  
1058 that help workers coordinate locally can become mechanisms of global control by enabling surveillance and proscription.  
1059 Thus, not all invisible work should be made visible. Bishop [1999] uses Weber's concepts of "status contract" and  
1060 "instrumental contract" to understand the changing relationships between employers and employees. She notes that  
1061 status contracts – those that are about our relations to one another rather than our performance – often rely on the  
1062 trust that results from these relationships, and not from formal articulations of the work. At ICPSR, we see evidence of  
1063 this status contract; by and large, those higher in the organization respect the skill and expertise of their curators. There  
1064 is an understanding that some aspect of data work will always be invisible. The use of Jira tickets to monitor, however,  
1065 threatens to replace this status contract with an instrumental one, in which the worker is valued for visible products.  
1066

1067 How does viewing curation as a craft impact this (in)visibility? When supervisors, project managers, and archive  
1068 directors view and treat curation as a craft, it supports the status contract between curators and higher management. It  
1069 appreciates this data work as skilled labor, and thereby "affords identity, status and a sense of connection to others  
1070 in the enterprise and to the enterprise itself" (Nardi and Engeström [1999] citing Bishop [1999]). When we as data  
1071 practices researchers, data science educators, and CSCW theorists argue for curation as craft, we, too contribute to  
1072 the support of this status contract. Thus, recognizing curation as craft is important to supporting labor arrangements  
1073 that do not render the worker invisible even when the work is. A well known impact of the invisibility of curation  
1074 work is that outsiders underestimate its costs and value, and, by implication, the value of curators. Work like curation  
1075 that is conducted in the background is often taken for granted. Recent efforts to surface curatorial contributions to  
1076 scholarship via structured metadata [Thessen et al. 2019] or improvement of legacy data records [Bionomia [n.d.]] echo  
1077 prior efforts such as the Nursing Interventions Classification to make work visible in efforts to legitimate both the  
1078 work and workers [Bowker et al. 1996; Star and Strauss 1999]. Here, we are pushing to recognize the labor needed to  
1079 organize, understand, and negotiate the tidy boxes on workflow diagrams – and to recognize it's seeming ineffibility as  
1080 important to preserve and respect.  
1081

#### 1082 5.4 Implications for practice

1083 Better articulating the work and craft of data curation has several implications for practice. First and foremost,  
1084 understanding the complex role that different forms of visibility play in data work may help us design technologies for  
1085

1093 users that move beyond reporting and surveillance [Suchman 1995]. As we described, many ICPSR curators bristled at  
1094 the constant use of Jira because it made them feel hypervisible, monitored, and mildly harangued. We consequently ask:  
1095 given a view of curation as craft rather than rotely mechanical labor, what changes might we imagine for ticketing  
1096 systems like Jira -- ones that might lead the management system to serve the curators as well as their supervisors and  
1097 managers?  
1098

1099 Our work also has several implications for data curation training and education. Within the information sciences,  
1100 considerable effort has been put into designing data curation curriculum for budding information professionals; much of  
1101 this has been highly focused on articulating different versions of data curation workflows, and describing data practices  
1102 in different fields. The range of high-level curatorial actions we identified in section 4.1 contributes to this tradition.  
1103

1104 However, the greater contribution of our work is the importance of training data curators as craftspeople and not  
1105 just technicians. As we quoted from Owens [2018] previously, a craftful approach to curation is one that stays engaged  
1106 with the “unresolved” and contingent aspects of curatorial work, and one that sees the “inherent messiness” of data  
1107 work as a feature, rather than a bug. Our work helps more specifically identify the strategies ICPSR curators use to  
1108 navigate this unresolved messiness, particularly in how they use gestalt approaches to see the dataset as more than the  
1109 sum of its parts. Though more work would be needed to better understand this process, we believe it is a promising  
1110 direction for further curriculum develop – whether in data curation classes at the master’s level, or online lessons in  
1111 the vein of Data Carpentry (<https://datacarpentry.org/>).  
1112

1113 One of the limitations of this study is ICPSR’s unusual size and scope; they simply have a much larger and more well  
1114 organized data curation team that many other peer institutions and archives. It is possible that lessons learned here  
1115 will not translate well to smaller contexts or teams. However – we believe this view of data work as grounded in craft  
1116 practices could be important to explore elsewhere. How do craft practices differ in smaller organizations, or in teams  
1117 where there are not dedicated curators? For those that think of data curation as a “wrangling” or “munging” process,  
1118 how could adopting a craft perspective help guide this work and make it more reproducible?  
1119  
1120  
1121

1122  
1123

## 1124 6 CONCLUSION

1125

1126 Data curation is a critical component of data science, and an important aspect of data work. Obscuring the work of data  
1127 curation not only renders the labor and contributions of the data curators invisible; it also makes it harder to tease out  
1128 the impact curators’ work has on the later usability, reliability, and reproducibility of data. In this paper we have made  
1129 curatorial work visible through a case study of data curation at ICPSR, a large social science data repository. We have  
1130 contributed a rich description of curatorial work at this site, including a range of technical curatorial actions, and the  
1131 craft and coordination needed to successfully enact those actions. We echo prior work calling for a craftful view of  
1132 work with data: curation requires not just a rote following of standards and protocols, but rather, a creative, on-going  
1133 conversation with the data, with one’s colleagues, and with one’s community. Our work complicates “workflow” based  
1134 views of data curation, in that we find ICPSR curators do considerable work that can’t be easily visualized with a  
1135 UML diagram, and indeed, rely on craft practices to work their workflow. We also find that ICPSR curators sit at an  
1136 intersection between visibility and invisibility: their work is highly documented (and even monitored, to a degree),  
1137 yet when they do their jobs well, it is invisible. Finding ways of selectively making curatorial work visible in service  
1138 of curators will be key in supporting their work and professional development, as well as the development of data  
1139 curation tools.  
1140  
1141  
1142  
1143

## REFERENCES

- 1145  
1146 M S Ackerman and C Halverson. 1999. Organizational memory: processes, boundary objects, and trajectories. In *Proceedings of the 32nd Annual Hawaii*  
1147 *International Conference on Systems Sciences. 1999. HICSS-32. Abstracts and CD-ROM of Full Papers*, Vol. Track1. IEEE, Maui, HI, USA.
- 1148 Katherine G Akers and Jennifer Doty. 2013. Disciplinary differences in faculty research data management practices and perspectives. *International Journal*  
1149 *of Digital Curation* 8, 2 (2013), 5–26.
- 1150 Karen Baker and Florence Millerand. 2007. Articulation Work Supporting Information Infrastructure Design: Coordination, Categorization, and Assessment  
1151 in Practice. In *2007 40th Annual Hawaii International Conference on System Sciences (HICSS'07). 2007 40th Annual Hawaii International Conference on*  
1152 *System Sciences (HICSS'07)*, 242a–242a.
- 1153 Stephen R Barley and Julian E Orr. 1997. Introduction: The neglected workforce. In *Between Craft and Science*. Cornell University Press, Ithaca, NY, 1–20.
- 1154 William C Barley, Jeffrey W Treem, and Paul M Leonardi. 2020. Experts at coordination: Examining the performance, production, and value of process  
1155 expertise. *Journal of Communication* 70, 1 (2020), 60–89.
- 1156 Bionomia [n.d.]. *Bionomia*. Accessed: 2021-7-14.
- 1157 Libby Bishop. 1999. Visible and Invisible Work: The Emerging Post-Industrial Employment Relation. *Computer Supported Cooperative Work* 8, 1-2 (March  
1158 1999), 115–126.
- 1159 Christine L Borgman. 2015. *Big Data, Little Data, No Data: Scholarship in the Networked World*. MIT Press.
- 1160 Christine L Borgman, Andrea Schamhorst, and Milena S Golshan. 2019. Digital data archives as knowledge infrastructures: Mediating data sharing and  
1161 reuse. *J. Assoc. Inf. Sci. Technol.* 70, 8 (Aug. 2019), 888–904.
- 1162 Geoffrey C Bowker, Stefan Timmermans, and Susan Leigh Star. 1996. Infrastructure and Organizational Transformation: Classifying Nurses' Work. In  
1163 *Information Technology and Changes in Organizational Work: Proceedings of the IFIP WG8.2 working conference on information technology and changes in*  
1164 *organizational work, December 1995*, Wanda J Orlikowski, Geoff Walsham, Matthew R Jones, and Janice I Degross (Eds.). Springer US, Boston, MA,  
344–370.
- 1165 Tiffany C Chao, Melissa H Cragin, and Carole L Palmer. 2015. Data Practices and Curation Vocabulary (DPCVocab): An empirically derived framework of  
1166 scientific data practices and curatorial processes: Data Practices and Curation Vocabulary (DPCVocab). *J. Assoc. Inf. Sci. Technol.* 66, 3 (March 2015),  
616–633.
- 1167 Melissa H Cragin, Carole L Palmer, Jacob R Carlson, and Michael Witt. 2010. Data sharing, small science and institutional repositories. *Philos. Trans. A*  
1168 *Math. Phys. Eng. Sci.* 368, 1926 (Sept. 2010), 4023–4038.
- 1169 Peter T Darch, Ashley E Sands, Christine L Borgman, and Milena S Golshan. 2020. Library cultures of data curation: Adventures in astronomy. *Journal of*  
1170 *the Association for Information Science and Technology* 71, 12 (2020), 1470–1483.
- 1171 DataONE. 2015. Data Life Cycle. <https://old.dataone.org/data-life-cycle>. Accessed: 2021-7-14.
- 1172 C Dearnley. 2005. A reflection on the use of semi-structured interviews. *Nurse Res.* 13, 1 (2005), 19–28.
- 1173 Catherine D'Ignazio and Lauren F Klein. 2020. *Data Feminism*. MIT Press.
- 1174 Paul Dourish. 2001. Process descriptions as organisational accounting devices: the dual use of workflow technologies. In *Proceedings of the 2001*  
1175 *International ACM SIGGROUP Conference on Supporting Group Work*. ACM, 52–60.
- 1176 Ingrid Erickson and Mohammad Hossein Jarrahi. 2016. Infrastructuring and the challenge of dynamic seams in mobile knowledge work. In *Proceedings of*  
1177 *the 19th ACM conference on Computer-Supported cooperative work & social computing*. ACM, 1323–1336.
- 1178 Ixchel M Faniel, Rebecca D Frank, and Elizabeth Yakel. 2019. Context from the data reuser's point of view. *Journal of Documentation* (2019).
- 1179 Ixchel M Faniel and Ann Zimmerman. 2011. Beyond the Data Deluge: A Research Agenda for Large-Scale Data Sharing and Reuse. *International Journal*  
1180 *of Digital Curation* 6, 1 (March 2011), 58–69.
- 1181 John L Faundeen, Thomas E Burley, Jennifer Carlino, David L Govoni, Heather S Henkel, Sally Holl, Vivian B Hutchison, Elizabeth Martin, Ellyn T  
1182 Montgomery, Cassandra C Ladino, et al. 2013. *The United States geological survey science data lifecycle model*. US Department of the Interior, US  
Geological Survey.
- 1183 Sebastian S. Feger, Paweł W. Wozniak, Lars Lischke, and Albrecht Schmidt. 2020. 'Yes, I Comply!': Motivations and Practices around Research  
1184 Data Management and Reuse across Scientific Fields. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW2, Article 141 (Oct. 2020), 26 pages. <https://doi.org/10.1145/3415212>
- 1185 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2021. Datasheets  
1186 for datasets. *Commun. ACM* 64, 12 (2021), 86–92.
- 1187 Elihu M Gerson and Susan Leigh Star. 1986. Analyzing due process in the workplace. *ACM Transactions on Information Systems (TOIS)* 4, 3 (1986), 257–270.
- 1188 Carole Goble, Robert Stevens, Duncan Hull, Katy Wolstencroft, and Rodrigo Lopez. 2008. Data curation + process curation = data integration + science.  
1189 *Brief. Bioinform.* 9, 6 (Nov. 2008), 506–517.
- 1190 Carole A Goble, Jiten Bhagat, Sergejs Aleksejevs, Don Cruickshank, Danius Michaelides, David Newman, Mark Borkum, Sean Bechhofer, Marco Roos,  
1191 Peter Li, and David De Roure. 2010. myExperiment: a repository and social network for the sharing of bioinformatics workflows. *Nucleic Acids Res.* 38,  
1192 Web Server issue (July 2010), W677–82.
- 1193 Ann G Green and Myron P Gutmann. 2007. Building partnerships among social science researchers, institution-based repositories and domain specific  
1194 data archives. *OCLC Systems & Services: International digital library perspectives* 23, 1 (Feb. 2007), 35–53.
- 1195  
1196

- 1197 Libby Hemphill, Margaret L Hedstrom, and Susan Hautaniemi Leonard. 2021. Saving social media data: Understanding data management practices among  
1198 social media researchers and their implications for archives. *J. Assoc. Inf. Sci. Technol.* 72, 1 (Jan. 2021), 97–109.
- 1199 Sharlene N Hesse-Biber and Patricia Leavy. 2005. *The practice of qualitative research* (third ed.). SAGE Publications.
- 1200 Hey T., Tansley S., Tolle K. (Ed.). 2009. *The Fourth Paradigm: Data-intensive Scientific Discovery*. Vol. 1. Microsoft Research Redmond, WA.
- 1201 Sarah Higgins. 2008. The DCC Curation Lifecycle Model. *International Journal of Digital Curation* 3, 1 (Dec. 2008), 134–140.
- 1202 ICPSR. 2020. ICPSR Curation Levels. <https://www.icpsr.umich.edu/files/datamanagement/icpsr-curation-levels.pdf>
- 1203 Eun Seo Jo and Timnit Gebru. 2020. Lessons from archives: strategies for collecting sociocultural data in machine learning. In *Proceedings of the 2020*  
1204 *Conference on Fairness, Accountability, and Transparency* (Barcelona, Spain) (FAT\* '20). Association for Computing Machinery, New York, NY, USA, 306–316.
- 1205 Lisa R Johnston. 2014. A Workflow Model for Curating Research Data in the University of Minnesota Libraries: Report from the 2013 Data Curation Pilot.
- 1206 Lisa R Johnston, Jacob Carlson, Cynthia Hudson-Vitale, Heidi Imker, Wendy Kozlowski, Robert Olendorf, and Claire Stewart. 2018. How Important Are  
1207 Data Curation Activities to Researchers? Gaps and Opportunities for Academic Libraries. *Journal of Librarianship and Scholarly Communication* 6, 1  
1208 (2018), eP2198.
- 1209 Sean Kandel, Andreas Paepcke, Joseph Hellerstein, and Jeffrey Heer. 2011. Wrangler: interactive visual specification of data transformation scripts. In  
1210 *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 3363–3372.
- 1211 H Karasti and K S Baker. 2004. Infrastructuring for the long-term: ecological information management.
- 1212 Helena Karasti, Karen S Baker, and Eija Halkola. 2006. Enriching the notion of data curation in E-science: Data managing and information infrastructuring  
1213 in the long term ecological research (LTER) network. *Comput. Support. Coop. Work* 15, 4 (Oct. 2006), 321–358.
- 1214 Karina Kervin, Robert B Cook, and William K Michener. 2014. The Backstage Work of Data Sharing. In *Proceedings of the 18th International Conference on*  
1215 *Supporting Group Work* (Sanibel Island, Florida, USA) (GROUP '14). Association for Computing Machinery, New York, NY, USA, 152–156.
- 1216 Sean Kross and Philip J Guo. 2021. Orienting, Framing, Bridging, Magic, and Counseling: How Data Scientists Navigate the Outer Loop of Client  
1217 Collaborations in Industry and Academia. arXiv:2105.05849 [cs.HC]
- 1218 Margaret D LeCompte and Jean J Schensul. 2012. *Analysis and interpretation of ethnographic data: A mixed methods approach* (second ed.). Rowman &  
1219 Littlefield.
- 1220 Helena M Mentis, Ahmed Rahim, and Pierre Theodore. 2016. Crafting the Image in Surgical Telemedicine. In *Proceedings of the 19th ACM Conference on*  
1221 *Computer-Supported Cooperative Work & Social Computing* (San Francisco, California, USA) (CSCW '16). Association for Computing Machinery, New  
1222 York, NY, USA, 744–755.
- 1223 Matthew B Miles, A Michael Huberman, and Johnny Saldaña. 2014. *Qualitative data analysis: A methods sourcebook* (third ed.). SAGE Publications.
- 1224 Florence Millerand and Geoffrey C Bowker. [n.d.]. Trajectories and Enactment in the Life of an Ontology. In *Standards and Their Stories*, Susan Leigh Star  
1225 and Martha Lampland (Eds.). Cornell University Press, 149–165.
- 1226 Michael Muller, Ingrid Lange, Dakuo Wang, David Piorkowski, Jason Tsay, Q Vera Liao, Casey Dugan, and Thomas Erickson. 2019. How Data Science  
1227 Workers Work with Data: Discovery, Capture, Curation, Design, Creation. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing*  
1228 *Systems*. Association for Computing Machinery, New York, NY, USA, 1–15.
- 1229 Tahani Nadim. 2016. Data labours: How the sequence databases GenBank and EMBL-bank make data. *Sci. Cult.* 25, 4 (Oct. 2016), 496–519.
- 1230 Bonnie A Nardi and Yrjö Engeström. 1999. A Web on the Wind: The Structure of Invisible Work. *Comput. Support. Coop. Work* 8, 1 (March 1999), 1–8.
- 1231 National Academies of Sciences, Engineering, and Medicine, Policy and Global Affairs, Committee on Science, Engineering, Medicine, and Public Policy,  
1232 Board on Research Data and Information, Division on Engineering and Physical Sciences, Committee on Applied and Theoretical Statistics, Board on  
1233 Mathematical Sciences and Analytics, Division on Earth and Life Studies, Nuclear and Radiation Studies Board, Division of Behavioral and Social  
1234 Sciences and Education, Committee on National Statistics, Board on Behavioral, Cognitive, and Sensory Sciences, and Committee on Reproducibility  
1235 and Replicability in Science. 2019. *Reproducibility and Replicability in Science*. National Academies Press.
- 1236 Andrew B Neang, Will Sutherland, Michael W Beach, and Charlotte P Lee. 2021. Data Integration as Coordination: The Articulation of Data Work in an  
1237 Ocean Science Collaboration. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW3 (2021), 1–25.
- 1238 Trevor Owens. 2018. *The Theory and Craft of Digital Preservation*. Johns Hopkins University Press.
- 1239 Carole L Palmer. 2006. Weak information work and “doable” problems in interdisciplinary science. *Proceedings of the American Society for Information*  
1240 *Science and Technology* 43, 1 (2006), 1–16.
- 1241 Carole L Palmer, Melissa H Cragin, and Timothy P Hogan. 2007. Weak information work in scientific discovery. *Information processing & management* 43,  
1242 3 (2007), 808–820.
- 1243 Carole L. Palmer, Nicholas M. Weber, Trevor Muñoz, and Allen H. Renear. 2013. Foundations of Data Curation: The Pedagogy and Practice of “Purposeful  
1244 Work” with Research Data. *Archive Journal* (2013).
- 1245 Amandalynne Paullada, Inioluwa Deborah Raji, Emily M Bender, Emily Denton, and Alex Hanna. 2021. Data and its (dis)contents: A survey of dataset  
1246 development and use in machine learning research. *Patterns* 2, 11 (Nov. 2021).
- 1247 Jean-Christophe Plantin. 2019. Data Cleaners for Pristine Datasets: Visibility and Invisibility of Data Processors in Social Science. *Sci. Technol. Human*  
1248 *Values* 44, 1 (Jan. 2019), 52–73.
- 1249 Line Pouchard. 2016. Revisiting the Data Lifecycle with Big Data Curation. *International Journal of Digital Curation* 10, 2 (May 2016), 176–192.
- 1250 Katie Rawson and Trevor Muñoz. 2016. Against cleaning. *Curating Menus* 6 (2016), 1–14.
- 1251 Manuscript submitted to ACM

- 1249 Daniela K Rosner, Samantha Shorey, Brock R Craft, and Helen Remick. 2018. Making core memory: Design inquiry into gendered legacies of engineering  
1250 and craftwork. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. 1–13.
- 1251 Herbert J Rubin and Irene S Rubin. 2012. *Qualitative interviewing: The art of hearing data* (third ed.). SAGE Publications.
- 1252 Morgan Klaus Scheuerman, Alex Hanna, and Emily Denton. 2021. Do Datasets Have Politics? Disciplinary Values in Computer Vision Dataset Development.  
1253 *Proc. ACM Hum.-Comput. Interact.* 5, CSCW2 (Oct. 2021), 1–37.
- 1254 William A Scott. 1955. Reliability of Content Analysis: The Case of Nominal Scale Coding. *Public Opinion Quarterly* 19, 3 (1955), 321–325.
- 1255 Susan Leigh Star and Anselm Strauss. 1999. Layers of Silence, Arenas of Voice: The Ecology of Visible and Invisible Work. *Computer Supported Cooperative  
1256 Work* 8, 1 (March 1999), 9–30.
- 1257 Pontus Stenertorp, Sampo Pyysalo, Goran Topić, Tomoko Ohta, Sophia Ananiadou, and Jun'ichi Tsujii. 2012. brat: a Web-based Tool for NLP-Assisted  
1258 Text Annotation. In *Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics*.  
Association for Computational Linguistics, Avignon, France, 102–107.
- 1259 Lucy Suchman. 1995. Making Work Visible. *Commun. ACM* 38, 9 (Sept. 1995), 56–64.
- 1260 Alex S Taylor, Siân Lindley, Tim Regan, David Sweeney, Vasillis Vlachokyriakos, Lillie Grainger, and Jessica Lingel. 2015. Data-in-Place: Thinking through  
1261 the Relations Between Data and Community. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. Association  
1262 for Computing Machinery, New York, NY, USA, 2863–2872.
- 1263 Anne E Thessen, Matt Woodburn, Dimitrios Koureas, Deborah Paul, Michael Conlon, David P Shorthouse, and Sarah Ramdeen. 2019. Proper Attribution  
1264 for Curation and Maintenance of Research Collections: Metadata Recommendations of the RDA/TDWWG Working Group. *Data Science Journal* 18, 1  
1265 (Nov. 2019), 54.
- 1266 Andrea K Thomer, Michael Bernard Twidale, and Matthew J Yoder. 2018a. Transforming Taxonomic Interfaces: "Arm? s Length" Cooperative Work and  
1267 the Maintenance of a Long-lived Classification System. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (2018), 1–23.
- 1268 Andrea K Thomer, Karen M Wickett, Karen S Baker, Bruce W Fouke, and Carole L Palmer. 2018b. Documenting provenance in noncomputational  
1269 workflows: Research process models based on geobiology fieldwork in Yellowstone National Park. *Journal of the Association for Information Science  
1270 and Technology* 69, 10 (2018), 1234–1245.
- 1271 Mary Vardigan, Pascal Heus, and Wendy Thomas. 2008. Data Documentation Initiative: Toward a Standard for the Social Sciences. *International Journal  
1272 of Digital Curation* 3, 1 (Dec. 2008), 107–113.
- 1273 Jullian C Wallis, Christine L Borgman, Matthew S Mayernik, and Alberto Pepe. 2008. Moving Archival Practices Upstream: An Exploration of the Life  
1274 Cycle of Ecological Sensing Data in Collaborative Field Research. *International Journal of Digital Curation* 3, 1 (Dec. 2008), 114–126.
- 1275 Hadley Wickham. 2014. Tidy Data. *J. Stat. Softw.* 59, 10 (2014), 1–23.
- 1276 Mark D Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten,  
1277 Luiz Bonino da Silva Santos, Philip E Bourne, Jildau Bouwman, Anthony J Brookes, Tim Clark, Mercè Crosas, Ingrid Dillo, Olivier Dumon, Scott  
1278 Edmunds, Chris T Evelo, Richard Finkers, Alejandra Gonzalez-Beltran, Alasdair J G Gray, Paul Groth, Carole Goble, Jeffrey S Grethe, Jaap Heringa,  
1279 Peter A C 't Hoen, Rob Hoof, Tobias Kuhn, Ruben Kok, Joost Kok, Scott J Lusher, Maryann E Martone, Albert Mons, Abel L Packer, Bengt Persson,  
1280 Philippe Rocca-Serra, Marco Roos, Rene van Schaik, Susanna-Assunta Sansone, Erik Schultes, Thierry Sengstag, Ted Slater, George Strawn, Morris A  
1281 Swertz, Mark Thompson, Johan van der Lei, Erik van Mulligen, Jan Velterop, Andra Waagmeester, Peter Wittenburg, Katherine Wolstencroft, Jun  
1282 Zhao, and Barend Mons. 2016. The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* 3 (March 2016), 160018.
- 1283 Michael Witt, Jacob Carlson, D Scott Brandt, and Melissa H Cragin. 2009. Constructing data curation profiles. *International Journal of Digital Curation* 4,  
1284 3 (2009), 93–103.
- 1285 Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2016. Wikipedia Talk Labels: Personal Attacks.
- 1286 Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2017. Ex Machina: Personal Attacks Seen at Scale. In *Proceedings of the 26th International Conference on  
1287 World Wide Web (Perth Australia) (WWW '17, Vol. 11)*. International World Wide Web Conferences Steering Committee, Republic and Canton of  
1288 Geneva, Switzerland, 1391–1399.
- 1289 Elizabeth Yakel. 2007. Digital curation. *OCLC Systems & Services: International digital library perspectives* 23, 4 (Nov. 2007), 335–340.
- 1290 J Yates. 1989. *Control Through Communication: The Rise of System in American Management*. Johns Hopkins University Press, Baltimore, MD.
- 1291 Amy X Zhang, Michael Muller, and Dakuo Wang. 2020. How do Data Science Workers Collaborate? Roles, Workflows, and Tools. *Proc. ACM Hum.-Comput.  
1292 Interact.* 4, CSCW1 (May 2020), 1–23.
- 1293 Ning Zhang, Manohar Paluri, Yaniv Taigman, Rob Fergus, and Lubomir Bourdev. 2015. Beyond frontal faces: Improving person recognition using multiple  
1294 cues. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. cv-foundation.org, 4804–4813.
- 1295 Jun Zhao, Jose Manuel Gomez-Perez, Khalid Belhajjame, Graham Klyne, Esteban Garcia-Cuesta, Aleix Garrido, Kristina Hettne, Marco Roos, David  
1296 De Roure, and Carole Goble. 2012. Why workflows break—Understanding and combating decay in Taverna workflows. In *2012 IEEE 8th international  
1297 conference on e-science*. IEEE, 1–9.

## 1296 A INTERVIEW PROTOCOL

### 1297 Background

- 1298 • How long have you been a [ROLE]?

- 1301 – What’s your academic background?
- 1302 – Do you have any prior experience as a data manager, curator, etc.
- 1303 • I would like to start by hearing more about what your [ROLE] work at ICPSR.
- 1304 • What is your role in the curation process?
- 1305 – Would you describe the chain of command? (e.g., Archive directors, curators)
- 1306 – How do you work together to make decisions?
- 1307 \* How do you interact with project managers?
- 1308 \* Do you feel involved in making judgement calls?
- 1309 \* When a study seems like it falls between two levels of curation, who determines which level to assign it?
- 1310 · What factors are important to this determination?
- 1311 – Can you describe your overall workflow to me?
- 1312 – Is your work specialized to an archive/domain/data type?
- 1313 • How much curation do the datasets you work with typically need?
- 1314 – Do they tend to arrive in the same state?
- 1315 • Would you tell us a bit about the scholarly community that uses your archive?
- 1316 • What type of relationship does the archive have with the scholarly community reusing its data?
- 1317
- 1318
- 1319
- 1320
- 1321
- 1322

### 1323 Curation

- 1324 • [If applicable] What types of interactions do you have with the curation unit?
- 1325 • How involved are you in curation decisions?
- 1326 – When does this occur (grant proposal, initiation of a grant/project planning, before/during data sharing)?
- 1327 – [If applicable] Were you involved in recent decision to make ICPSR curation workflows more systematic?
- 1328 If so, can you tell us what led to that decision?
- 1329 – How have recent changes to curation workflows at ICPSR changed your involvement in curation decisions, if at all?
- 1330 – Is there a formal process?
- 1331 – Are these decisions always easy to implement?
- 1332 • Which curation activities add the most value to your archive/datasets?
- 1333 – Why do you say this?
- 1334 • How do you prioritize different curatorial activities?
- 1335 – How do you know when a dataset is “done” being curated?
- 1336 – How involved are you in making judgement calls (e.g. between levels, between different curatorial actions)?
- 1337 • Are the curation levels well defined?
- 1338 – Do you think they work for most studies?
- 1339 \* Why or why not?
- 1340 • How has your job changed since the curation reorganization? [Tailor to ROLE and BACKGROUND]
- 1341 • Do your data reusers or designated community provide input into curatorial decisions?
- 1342 • How has the curation provided by ICPSR changed the use or impact of your collections?
- 1343 • Do you ever question the amount of curation planned for or being applied to a dataset?
- 1344 • Is there additional/different curation you’d like to see applied to some of your datasets?
- 1345
- 1346
- 1347
- 1348
- 1349
- 1350
- 1351
- 1352



- What metrics would you propose or like to guide the level of curation of data?

#### Impact and metrics

- What type of impact would you like your archive to have?
  - How close is the archive to achieving this goal?
  - Where would you like to see the impact of your archive in 5 years?
- Are you aware of any metrics at ICPSR guiding the curation process?
- What metrics do you currently use to measure the impact of your collection, if any?
- What metrics would you propose to measure the impact of your collection?
- [If applicable] Do you plan or discuss your curation work with anyone else at ICPSR?
  - If yes, how do their comments impact your curatorial decisions?
- [If applicable] Do you consider the potential impact of the dataset during curation?
  - If so, how?
- [If applicable] How does your work add value to the datasets you curate?
  - Probe if not answering specifically: application of standards, metadata
    - \* Which ones?
- How do you see x (e.g., metadata, data cleaning, etc.) having impact on the datasets?
  - [If applicable] In what ways do the JIRA tickets document the curation work you've done?
  - [If applicable] In what ways do the JIRA tickets *not* document the curation work you've done?
  - [If applicable] Do you find JIRA intrusive?
- Is there a dataset that you've worked on that's had substantial impact?
  - If so, could you describe what it was and what impact it made?
  - What contribution did your curatorial work have on this dataset's impact?
- How would your designated community define impact of the collections?
- What impact would the faculty that contribute data to your archive want the collection to have?
  - Do you believe that data sharing and reuse should be considered for promotion and tenure?
  - How broadly is this shared in the scholarly community served by the archive?
- What kind of impact do you want your work to have? Your collections to have?

#### Reuse

- Do you consider data reusers during the curation process?
  - If so, what are the significant characteristics or properties (e.g., information about the data that is important for effective preservation management or reuse) you think are important to capture to enable data reuse?
- Do you interact with data reusers?
  - If yes, do their comments impact your curatorial decisions?
- What are the greatest barriers in re-using collections from your archive?
- What kinds of input or questions do you get from data reusers?
- What do reusers tell you about the value of different datasets?
- What kind of reuse would you like to facilitate in the future?

#### Wrap up

- Do you have any questions for us, or about this project?