

Interactivity, Distributed Workflows, and Thick Provenance: A Review of Challenges Confronting Digital Humanities Research Objects

Katrina Fenlon

College of Information Studies
University of Maryland
College Park, MD, United States
kfenlon@umd.edu

Abstract—Despite the rapid growth of digital scholarship in the humanities, most existing humanities research infrastructures lack adequate support for the creation, management, sharing, maintenance, and preservation of complex, networked digital objects. Research Objects (ROs) have mainly been applied to scientific research workflows, but the RO model and parallel approaches have gained enough uptake in the humanities to suggest their potential to undergird sustainable, networked humanities research infrastructure. This paper reviews several compelling applications in the humanities of RO and closely related models in platforms for data sharing, computational text analysis, collaborative annotation, digital and semantic publishing, and in domain repositories. The paper identifies challenges confronting the broad application of ROs in the humanities—which challenges will confront any emergent model for humanities data- or workflow-packaging and publication—and suggests implications for implementations in humanities cyberinfrastructure.

Index Terms—digital humanities, research objects, workflows, data models

I. INTRODUCTION

While Research Objects (ROs) have primarily been oriented toward scientific research workflows, the RO model and parallel approaches have gained enough uptake in the humanities to suggest their potential to undergird sustainable, networked humanities research infrastructures. Digital scholarship in the humanities takes a wide variety of forms, incorporating narratives, media, datasets, interactive components, etc.—any of which may be physically dispersed as well as dynamic and evolving over time. Despite the rapid growth of humanities digital scholarship, most existing humanities research infrastructures lack support for the creation, management, sharing, maintenance, and preservation of complex, networked digital objects and datasets. ROs, and the community and tools that are growing around ROs, offer a potential, partial solution.

The concept of the RO has seen significantly more uptake in the humanities than has the formal data model [1], [2]; nonetheless, several compelling applications of the concept suggest the time is ripe for considering broader integration of the model into distributed infrastructures. Such applications include platforms for data sharing, computational text analysis, and collaborative scholarship; platforms for digital and semantic publishing; and digital repositories in several domains.

This paper reviews existing applications of the RO model to identify challenges confronting the application of ROs to humanities digital scholarship. This paper builds on [3], which investigated the application of the RO model to digital humanities collections, and which identified three promising strengths of the model for humanities scholarly communication: (1) ROs readily perform the most essential function of a collection: to aggregate related resources in order to support scholarly objectives; (2) ROs have the capacity for explicit, semantic descriptions of interrelationships among components that are often “hidden” in digital humanities collections, and therefore most vulnerable to dissolution; and (3) the RO model accommodates aggregations of linked data, offering researchers the opportunity to create and annotate virtual, fully referential collections. This paper builds on the prior analysis by reviewing the literature on ROs in the humanities and examining a range of applications of the RO and similar models within humanities and cultural heritage domains.

The review is framed around three main challenges and their implications for future implementations of ROs, or indeed any packaging and publication models for data and workflows, to support digital research in the humanities. First, digital humanities scholarship requires highly specialized, interactive uses within communities, so realizing the advantages of ROs for the humanities will depend on implementations that create platforms for community-centered experimentation and development. Second, the idiosyncratic workflows employed in the construction of networked humanities scholarship suggest that workflow-oriented ROs will not gain significant uptake in the humanities unless they can capture distributed, sociotechnical workflows in meaningful ways. Third, humanities ROs will require documenting provenance at a level of nuance that may exceed the capacity of automated or systematic capture; humanities scholarship requires “thick,” multilayered, context-rich provenance descriptions to accommodate conflicting assertions and formalize the expression of uncertainty.

II. ESSENTIAL INTERACTIVITY FOR SPECIALIZED USE

Much of humanities digital scholarship is *essentially* interactive. New modes of production and publication in the hu-

manities are intended for user interaction through dynamic, responsive, and even openly participatory representations based on research context. Digital collections and archives, digital editions, maps, models, simulations, and other modes of digital scholarship all rely on interactive components to express their interpretive contributions, or to enact their scholarly purposes. The interactive and dynamic components of digital scholarship include things like customized browsing and searching facilities that take advantage of extensive, rich scholarly encodings and annotations; platforms for collaborative annotation; dynamic maps and visualizations; etc. Such components are intended to do multiple things at once: to express arguments, manifest interpretive stances, enable knowledge transfer, and simultaneously serve as platforms for ongoing interpretation and research [4], [5].

Prior empirical work on applying the RO model to digital humanities collections found the main limitation of the model to be that functional components, designed for ongoing end-user interaction, are not usefully captured in a basic RO model and instead fall to the implementations built on top of research-object management systems [3]. ROs can, of course, accommodate as flat code objects that must be functional and interactive to serve their purposes. ROs have been employed for this purpose to support data migration and archiving (e.g., the RO BagIt profile). However, for ROs to be useful in this domain will require implementations that support platforms for flexible, participatory development.

In a conceptual sense, the RO model has demonstrated value for this kind of platform approach in the humanities. The Perseids project offers a platform for sharing and peer-review of research data in the Classics, including transcriptions, annotations, and analyses. The Perseids architecture is built around the concept of *data publications*, modeled as collections of related data objects. The Perseids team explicitly relates the data publication model to the RO model [6]. Like ROs, Perseids *data publications* rely on several domain standards (including the TEI Epidoc schema, W3C Web Annotation, and others) to undergird an infrastructure that supports domain-specific requirements: transcription, fine-grained annotation, collaborative editing (with versioning), a research environment that facilitates data-type-specific extensions, and tailored workflows for peer review [6]. Similarly, the Community Enhanced Repository for Engaged Scholarship (CERES) toolkit, created by the Northeastern University Libraries Digital Scholarship Group, explicitly draws on the concept of the RO in a system for supporting networked humanities scholarship and publishing. CERES allows digital humanities creators to build custom publications pulling objects through API from preservation-oriented repositories (including the Northeastern University Libraries' Digital Repository Service and the Digital Public Library of America) [7].

It is unclear how the RO model may fit into the broader, diversified landscape of linked data and the Semantic Web in humanities and cultural institutions, but the conceptual fit within digital scholarly communication is established. ROs and similar models have substantial potential to underpin

systems that support a variety of implementations. Realizing the advantages of ROs for the humanities will depend on implementations that create platforms for experimentation and collaborative development within distributed communities [3], including participation by the addition of linked data, annotations, and enrichments, including linking among ROs and the concepts and entities represented within ROs.

III. DISTRIBUTED AND IDIOSYNCRATIC WORKFLOWS OF NETWORKED HUMANITIES SCHOLARSHIP

Humanities digital scholarship is increasingly networked, in the sense of being heavily interconnected with and dependent on external resources for functionality and meaning. Many digital humanities forms—monographs, multimedia productions, exhibits, collections—draw on, reference, embed, and patch together distributed resources called from other collections, often via API. For example, a collection may center on a set of high-resolution images of primary sources called from the IIF image server of an independent digital library. Some of the longest-running, large-scale cultural heritage digital libraries (including *Europeana* and the *Digital Public Library of America*) are aggregations of descriptive surrogates, which link to original content hosted externally. Externally maintained schemas, authorities, and utilities undergird digital editions. Visualization and mapping projects generate content using external services. And with the growth of linked data in cultural collections, projects increasingly leverage external data sources as primary content, to which scholars then add layers of interpretive narrative, annotations, context, and interconnection.

Humanities workflows rarely happen in self-contained or end-to-end research infrastructures, thwarting the possibility of sufficiently rich, *automatic* workflow capture. Indeed, efforts to build a workflow-oriented, unified cyberinfrastructure for supporting humanities scholarship tend to founder (e.g. [8]). However, niche, task- or domain-specific infrastructures can capture constrained workflows. For example, in the domain of musicology, Page et al. [9] observe how digital editions and annotations of encoded works are “manifestations of workflows deployed in musicological scholarship,” and offer a compelling framework for representing musical ROs, which include images, text, audio, and encoded music [9], [10]. Fully computational workflows are readily captured within constrained humanities research environments, and ROs have come into play for this purpose. For example, the HathiTrust Research Center Data Capsule environment is moving toward systematic provenance-capture for computational text analysis workflows. These workflows take as inputs *workssets* [11], which are conceptually and technically akin to ROs: aggregate digital objects that implement addressability for and relational expressivity among components using domain ontologies. Unlike ROs, *workssets* are envisioned as the inputs of workflows in the current model of the HathiTrust Data Capsule environment, rather than encompassing whole research workflows [12]. But workflow-oriented ROs will not gain significant uptake in humanities contexts unless they can also capture and make

useful more complex, distributed, sociotechnical workflows in meaningful ways.

With their capacity for linked data using domain vocabularies, ROs readily accommodate many of the *artifacts* of networked digital scholarship in the humanities, along with their interrelationships [3]. But can ROs accommodate humanities workflows in useful ways? In their effort to undergird DARIAH (pan-European infrastructure for digital arts and humanities research) through the systematic production of humanities ROs, Blanke and Hedges observed that humanities scholars employ sequential workflows, but found scant evidence of usefully reproducible workflows [13]. While auto-generated, computer-useable workflows may not apply to most humanities research processes, formalized, (semi-) manually captured workflows would be highly useful for review, validation, archiving, reproducibility, reuse, and other purposes. While the RO model has the capacity and flexibility for complex workflow representation, more research is needed to characterize humanities workflows; to identify how such characterizations can be made useful; and to identify model extensions and unique implementation strategies workflows might require in different domains.

IV. THICK PROVENANCE

Drilling down on the problem of workflow capture, digital humanities scholarship places special demands on data provenance—not only on the provenance of digital resources (such as files, compound objects, datasets) or components thereof (such as passages of music, paragraphs of a text, or lines of a poem), but also the provenance of attached, contextual information. Archival artifacts—the evidence of the humanities—often possess simultaneous, multiple and parallel provenances [14], [15]. Documenting the provenance of the evidence itself can be complicated, but beyond that, the provenance *of the provenance* must also be documented. Any assertion made about any artifact (in the form of metadata or annotation), or any contextual and secondary information attached to artifacts in the context of digital scholarship, require provenance. Annotations and metadata are often, in the humanities, products of scholarly, interpretive work. Therefore, each annotation or metadata proposition itself is subject to claims of authorship, competing perspectives, expression of uncertainty, and further annotation—all requiring provenance information.

Because provenance is a multilayered thing in humanities scholarship, different humanities disciplines and sub-disciplines may require domain-specific provenance schemas and standards, which specialize existing standards for the expression of the provenance of different kinds of resources, ranging from digital media files to annotations. Humanities ROs will require thick, multilayered, context-rich provenance descriptions, which can accommodate conflicting assertions and formalize uncertainty. It is unclear whether existing implementations of the RO model can accommodate this level of description, though the model itself has the capacity.

The ResearchSpace environment [16] offers exemplary support for documentation of thick, multifaceted provenance of humanities ROs. ResearchSpace is an open-source platform created by the British Museum to facilitate scholarly data sharing, formal argumentation, and semantic publishing within communities of researchers. ResearchSpace does not directly employ the RO model, though its architecture does rely on aggregates of linked data, taking advantage of related standards including W3C Web Annotation and Linked Data Platform containers. In this environment, provenance and argumentation are expressed using the CIDOC-CRM specialization CRMInf (The Argumentation Model). Scholars can use this vocabulary to build narratives and thick descriptions around digital ROs through annotation and data-linking. These narratives of provenance allow and formalize the expression of uncertainty and competing perspectives, and the environment also serves to document the scholarly work that goes into building these narratives [17].

The reasons for highlighting the ResearchSpace approach to provenance in this review of humanities ROs are to (1) exemplify the unique demands of formalizing humanities provenance, and (2) exemplify the highly distinctive, domain-specific implementation requirements that confront the RO and other domain-independent data models. Describing humanities provenance will require vocabularies to express argument and belief, as Oldman et al. [18] observe. Beyond the RO model's use of Prov and Web Annotation, humanities provenance will demand domain-specific argumentation extensions such as CRMInf. It is clear that ROs can theoretically accommodate thick provenance description, just as they can theoretically accommodate the representation of highly complex workflows, but can they usefully undergird implementations that are centered in humanities research needs? The ResearchSpace interface is tailored toward knowledge work, toward the collaborative construction of multifaceted provenance descriptions, without requiring users to code or gain expert-level knowledge of domain ontologies. Tools for the authorship of humanities ROs, or tools that implement ROs behind the scenes, may benefit from taking the same approach.

V. CONCLUSION

ROs make a great deal of sense for modeling cultural information. Skeletons of similar shape—the simple and powerful combination of aggregation and annotation to represent compound digital objects—already structure large-scale cultural data aggregations, e.g., through the *Europeana Data Model* and the *Digital Public Library of America Metadata Application Profile*, which are both founded on *ore:aggregation* plus *oa:annotation*. But the challenges confronting widespread application of the RO model to humanities digital scholarship are significant. This review of existing applications has identified three central challenges, which may also prove resonant in disciplines beyond the humanities:

- (1) Digital humanities scholarship requires specialized interactive use; ROs for the humanities will depend on

implementations that create platforms for experimentation and development within communities.

- (2) The idiosyncratic workflows of humanities scholarship means that workflow-oriented ROs must capture distributed, sociotechnical workflows in meaningful ways.
- (3) Humanities ROs require thick, multilayered, context-rich provenance descriptions to accommodate conflicting assertions and formalize uncertainty, along with implementations that support such provenance documentation.

In particular, the challenge of characterizing and formally expressing diverse humanities workflows, along with the provenance of data and contextual information within those workflows, presents the most urgent challenge and exciting opportunity for the future of humanities cyberinfrastructure. To many stakeholders in humanities cyberinfrastructure, “workflows are the new content” [19]–[21]. While research on workflows is underway on multiple fronts (including [22]), it is clear already that there will be significant semantic differences between conceptual and technical elements in scientific workflows and those in the humanities, which will affect the implementation of ROs and other packaging and publication models for humanities research. Many prior attempts to implement scientific research infrastructures and data models to support humanities scholarship have run aground on basic semantic differences or conceptual gulfs between disciplines or domains within disciplines. The *Linking and Querying Ancient Texts* project, an effort to transfer eScience infrastructure in support of a humanities virtual research environment, observed a fundamental challenge in integrating humanities data from different databases and located the solution to that problem within humanities research communities: “integrating humanities research material...will require researchers to make the connections themselves, including decisions on how they are expressed and how to understand and explore the data more effectively” [23]. Oldman et al. [18], reviewing the state of linked data in the humanities, observe that basic linked data publication for many kinds of humanities sources can be counterproductive, “unless adapted to reflect specific methods and practices, and integrated into the epistemological processes they genuinely belong to.” This qualification resonates with the challenges identified for the adoption of the RO model—or indeed for the importation of any data model, including domain-independent data models—into the humanities. The main challenges to implementing ROs or any packaging or publication model for humanities research also present exciting opportunities for a more sustainable cross-disciplinary infrastructure [3], but implementation strategies must be centered in scholarly communities, and grow out from the practices, needs, and epistemologies of individual scholars and specific areas of study in the humanities and cultural institutions.

REFERENCES

- [1] S. Bechhofer et al., “Why linked data is not enough for scientists,” *Future Generation Computer Systems*, vol. 29, no. 2, pp. 599–611, Feb. 2013.
- [2] K. Belhajjame et al., “Using a suite of ontologies for preserving workflow-centric research objects,” *Journal of Web Semantics*, vol. 32, pp. 16–42, May 2015.
- [3] K. Fenlon, “Modeling Digital Humanities Collections as Research Objects,” in *Proceedings of the 19th ACM/IEEE Joint Conference on Digital Libraries*, Champaign, IL, USA, 2019.
- [4] C. L. Palmer, L. C. Tefteau, and C. M. Pirmann, “Scholarly Information Practices in the Online Environment: Themes from the Literature and Implications for Library Service Development,” OCLC Research and Programs, Dublin, OH, 2009.
- [5] K. Fenlon, “Thematic research collections: Libraries and the evolution of alternative scholarly publishing in the humanities,” Doctoral dissertation, University of Illinois at Urbana-Champaign, 2017.
- [6] B. Almas, “Perseids: Experimenting with Infrastructure for Creating and Sharing Research Data in the Digital Humanities,” *Data Science Journal*, vol. 16, no. 0, p. 19, Apr. 2017.
- [7] S. J. Sweeney, J. Flanders, and A. Levesque, “Community-Enhanced Repository for Engaged Scholarship: A case study on supporting digital humanities research,” *College Undergraduate Libraries*, vol. 24, no. 2–4, pp. 322–336, Oct. 2017.
- [8] Q. Dombrowski, “What Ever Happened to Project Bamboo?,” *Lit Linguist Computing*, vol. 29, no. 3, pp. 326–339, Sep. 2014.
- [9] K. Page, D. Lewis, and D. Weigl, “Contextual interpretation of digital music notation,” presented at the Digital Humanities (DH2017), Montreal, Canada, 2017, p. 3.
- [10] D. De Roure, G. Klyne, K. Page, J. Pybus, D. M. Weigl, and P. Willcox, “Digital Music Objects: Research Objects for Music,” presented at the Research Object workshop (RO2018) at IEEE eScience Conference 2018, 31-Jul-2018.
- [11] J. Jett, T. W. Cole, and J. S. Downie, “Exploiting graph-based data to realize new functionalities for scholar-built worksets,” *Proceedings of the Association for Information Science and Technology*, vol. 54, no. 1, pp. 716–717, 2017.
- [12] J. Murdock, J. Jett, T. Cole, Y. Ma, J. S. Downie, and B. Plale, “Towards Publishing Secure Capsule-based Analysis,” in *Proceedings of the 17th ACM/IEEE Joint Conference on Digital Libraries*, Piscataway, NJ, USA, 2017, pp. 261–264.
- [13] T. Blanke and M. Hedges, “Scholarly primitives: Building institutional infrastructure for humanities e-Science,” *Future Generation Computer Systems*, vol. 29, no. 2, pp. 654–661, Feb. 2013.
- [14] A. J. Gilliland, *Conceptualizing 21st-Century Archives*. ALA Editions, 2014.
- [15] C. Hurley, “Parallel provenance [Series of parts]: Part 1: What, if anything, is archival description?,” *Archives and Manuscripts*, vol. 33, no. 1, p. 110, May 2005.
- [16] D. Oldman and D. Tanase, “Reshaping the Knowledge Graph by Connecting Researchers, Data and Practices in ResearchSpace,” in *The Semantic Web —ISWC 2018*, 2018, pp. 325–340.
- [17] ResearchSpace Team, British Museum, “Moving from Documentation to Knowledge Building: ResearchSpace Principles and Practices,” Stiftung Preussischer Kulturbesitz (Prussian Cultural Heritage Foundation) Berlin, 10-Dec-2018.
- [18] D. Oldman, M. Doerr, and S. Gradmann, “Zen and the Art of Linked Data,” in *A New Companion to Digital Humanities*, John Wiley Sons, Ltd, 2015, pp. 251–273.
- [19] L. Dempsey, “The Library in the Life of the User: Two Collection Directions,” presented at the The transformation of academic library collecting: A symposium inspired by Dan C. Hazen, Harvard Library, Oct-2016.
- [20] M. A. Baynes, D. Sommer, D. Melley, and T. Lickiss, “Workflow is the new content: Expanding the scope of interaction between publishers and researchers,” presented at the Society for Scholarly Publishing, Apr-2016.
- [21] R. C. Schonfeld and D. Waters, “The turn to research workflow and the strategic implications for the academy,” presented at the Coalition for Networked Information (CNI) Spring Membership Meeting, San Diego, CA, Apr-2018.
- [22] A. Liu, S. Kleinman, J. Douglass, L. Thomas, A. Champagne, and J. Russell, “Open, Shareable, Reproducible Workflows for the Digital Humanities: The Case of the 4Humanities.org ‘WhatEvery1Says’ Project,” presented at the Digital Humanities (DH2017), 2017.
- [23] S. Anderson and T. Blanke, “Taking the Long View: From e-Science Humanities to Humanities Digital Ecosystems,” *Historical Social Research / Historische Sozialforschung*, vol. 37, no. 3 (141), pp. 147–164, 2012.