

Provenance visualization: Tracing people, processes, and practices through a data-driven approach to provenance

Tomas Vancisin ^{1*}, Loraine Clarke¹, Mary Orr¹, Uta Hinrichs²

¹University of St Andrews, St Andrews, UK

²University of Edinburgh, Edinburgh, UK

*Correspondence: Tomas Vancisin, School of Computer Science, University of St Andrews, Jack Cole Building, North Haugh, St Andrews KY16 9SX, UK. E-mail: tv8@st-andrews.ac.uk

Abstract

Provenance disclosure—the documentation of an artifact’s origin and how it was produced—is an important aspect to consider when working with historical records which undergo multiple transformations in preparation for and during digitization. Provenance in this context is commonly communicated through explanatory text or static diagrams. However, the methodological and curatorial decisions that have influenced the records’ data are easily overlooked, in particular when exploring the records through visualization as a result of digitization processes. We propose a data-driven approach to provenance disclosure which (1) traces provenance back to when the records were created, (2) documents and categorizes the records’ transformations (transcriptions, content modifications, changes in organization, and representational form), and (3) uses data visualization to disclose provenance in interactive ways. We reflect on how this approach can be practically applied in the context of historical record collections, and we present findings from a qualitative study we conducted to investigate the merits and limitations of provenance-driven visualization. Our findings suggest that data-driven provenance disclosure has the potential to (1) promote transparency and deeper interpretations of historical records, (2) provide rigor in researching historical document collections and underlying production processes, and (3) encourage ethical considerations by making visible labor and implicit bias that influence the production and curation of historical records.

1 Introduction

The digitization of historical documents such as manuscripts, national and university records, letters, or books, and, as part of this, the transformation of such records into computer-readable formats, has given rise to the visualization of such records, to facilitate their interrogation from multiple perspectives (Windhager *et al.*, 2019). Visualization can enable the identification of higher-level patterns across historical record collections and facilitate explorations that are not possible through close-reading techniques (Jänicke *et al.*, 2017; Windhager *et al.*, 2019). In the past few years, visualization has become an important means in various indicative research projects for investigating and interpreting digitized historical record collections (Betti *et al.*, 2014; Edelstein *et al.*, 2017; Hinrichs *et al.*, 2015; Mäkelä *et al.* 2012).

What is less clear in such visualizations of digitized historical record collections, however, is their provenance

disclosure, and why it matters quantitatively and qualitatively for users with regard to how they then interrogate the visualized data. Historical records undergo a number of transformation steps, including transcription and (re-)structuring (e.g. via tagging), before interactive visualizations can be developed. These transformations variously change the data’s content, organization, and artifactual form, in turn influencing how the records can be interrogated (e.g. what research questions are asked) and interpreted—individually and as a collection. Traditionally, provenance information is either described in textual form (e.g. Betti *et al.*, 2014; Hyvönen *et al.*, 2017) or illustrated through process diagrams (e.g. Capodiceci *et al.*, 2015; Hinrichs *et al.*, 2015). Both approaches can only address provenance at a high level, and treat provenance as secondary information that is easily overlooked when exploring the records in question through digital search interfaces or visualization.




Figure 3. Provenance-driven visualization; Selecting ‘M’ surnames in the third layer shows their form in the ‘Record View’ to the right and highlights the structure and organization of those records in the other layers




Figure 4. Provenance-driven visualization’s different record views

by Smart. The Records View (Fig. 3.6) therefore shows the content of individual records after this transformation in a layout, font, and style that resembles BRUSA as published by Smart (Smart, 2004). A selection of one of the squares in Layer 4 would show the corresponding record in the XML: TEI form (see Fig. 4.4).

This example of a provenance-driven visualization illustrates past forms of the student records and inherent changes to their content and structure in a visual way. The records’ content and provenance information are visually linked. Following a data-driven approach to provenance disclosure, we focused on capturing and disclosing aspects of the full spectrum of identified transformation processes, mostly drawing on provenance information that is already implicit in the different artificial forms of the records and that can be directly visualized. For example, we highlight → transcription through the choice of different stroke types (see Fig. 3.1 and 2). We make the → structural modifications introduced by Smart, Crawford and by ourselves explicit through the spatial distribution of elements representing the records (see Fig. 3.3, 4, and 5). We show → content modifications at the record-level and in aggregated ways by modifying the size of visual elements (see Figs 3.3 and 4). Modifications of → artifactual form are illustrated in the Record View directly capturing the original record forms and through typographical representation (see Figs 3.6 and 4).




Figure 5. Traditional visualization

5 Studying the impact of provenance-driven visualization

The notion of a data-driven approach to provenance and provenance-driven visualization illustrated above is a departure from the dominant approach of visualizing the *content* of historical records, to visualizing the context in which these have been collected and modified. This approach also departs from traditional ways of representing provenance (text or diagrams) in that it is informed by data we gathered systematically about each individual record's transformations.

In order to explore the potential of our data-driven approach to provenance disclosure, we conducted a qualitative study. We were particularly interested in the following questions: What types of insights and interpretations do participants gather from our provenance-driven visualization?, Does provenance-driven visualization promote transparency?, and Can it raise awareness of the labor and different layers of interpretation that are inherent in historical document collections?

5.1 Study approach

Evaluating how visualizations inform user insights and interpretations is a complex challenge, since analysis processes and outputs are difficult to capture (Lam *et al.*, 2012), especially in a study context that is typically constrained by time. In order to start addressing

our research questions above, we designed a qualitative study that exposed participants to two independent visualizations that sit at opposite ends of a content—provenance continuum. Our study should not be understood as a comparative appraisal where one condition is tested against the other. Instead, inspired by previous work in the context of personal visualization (Thudt *et al.*, 2016), we consider the two visualizations as probes to trigger situated reflections on two different aspects of historical data that can be made visible and, in consequence, on the impacts, these may have on potential explorations and interpretations of the records.

The provenance-driven visualization we showed participants is the one described in the previous section (see Fig. 3). The, more traditional, content-driven visualization focused on the geo-temporal aspects of the historical student records (see Fig. 5).

Geo-temporal visualizations are commonly used to provide an overview of historical document collections (Brizzi, 2013; Jenkins *et al.*, 2013; Edelstein *et al.*, 2017; Schwinges, 2018; Conroy, 2021). The map view shows the geographical distribution of the students' birth locations (see Fig. 5.1) while the timeline view depicts the number of graduates by year (see Fig. 5.2). Individual student records are shown in the 'Record View' (see Fig. 5.3). All views are linked; hovering over the circles on the map acts as a filter on the timeline

