**Collections as Data in Texas Digital Library Repositories:**

**Practical Recommendations for Use**

Karla Roig Blay, Anne Morgan, and Katy Tuck

School of Information, University of Texas at Austin

**Collections as Data in Texas Digital Library Repositories:**

**Practical Recommendations for Use**

This document provides practical recommendations for the use of Collections as Data in Texas Digital Library (TDL) repositories. The document presents recommendations for how to undertake a Collections as Data project using the [Texas Conference on Digital Libraries (TCDL) Proceedings](#) as a case study. The recommendations were written by three graduate students at the University of Texas at Austin's School of Information as part of a semester-long group project with TDL. The purpose of this document is to increase understanding and improve methods of using Collections as Data for TDL members.

**What is Collections as Data?**

In "On a Collections as Data Imperative," Thomas Padilla (2017, p. 1) describes Collections as Data as "reframing all digital objects as data," defined as "ordered information, stored digitally, that is inherently amenable to computation." In a sense, then, Collections as Data is a "mindset" (Wittmann et al., 2019, p. 51): viewing collections of digital objects in terms of discrete pieces of information that can be analyzed and used in meaningful ways. On a more practical level, Collections as Data involves the application of computational methods in order to analyze and use collections of digital objects meaningfully. These computational methods include text mining, data visualization, mapping, image analysis, audio analysis, and network analysis (Padilla et al., 2019, May 20b, p. 2).

The *[Always Already Computational: Collections as Data](#)* project, which began in 2016, produced guidance for institutions interested in working with Collections as Data. The project's deliverables include the [Santa Barbara Statement on Collections as Data](#), which outlines foundational Collections as Data principles, such as commitments to shared documentation, open access, and transparency. *Always Already Computational* also produced a set of [50 recommendations](#) intended to provide practitioners with initial steps to take in order to begin

supporting Collections as Data at cultural heritage institutions. In 2018, *Always Already Computational* was succeeded by a new project, [Collections as Data: Part to Whole](#), which has focused on developing models for the implementation and use of Collections as Data.

Using collections as data is an example of a digital humanities project. The field of digital humanities (DH) broadly refers to a certain methodological approach involving the application of technology tools to digital (or digitized) scholarship in order to analyze, visualize, interpret, and present this information in new ways. It is characterized by an emphasis on collaboration, interpretation, openness, interoperability, and innovation and can be applied to multiple disciplinary fields toward a variety of purposes. Examples of DH methods and uses include "text analysis, data mining, visualization, modeling and simulation, geospatial analysis and mapping, multimedia storytelling, information design, network analysis, interface design, and mark-up…[and applying] these tools to humanistic questions." ("About Us - Digital Humanities," n.d.). Many of these methods overlap with the CAD methods of inquiry mentioned above so using collections as data falls squarely into this category of digital humanities.

## What are the Benefits and Uses of Using Collections as Data?

The following are potential uses of Collections as Data:

- Discover patterns in data

- Present data in a new, interesting way

- Attract new users and drive site traffic

- Create compelling visualizations of datasets to tell a story about your collection(s)

These potential uses point to two main benefits of using Collections as Data: increased useability and visibility of collections. Collections as Data moves "beyond traditional use," allowing "more flexible access" to collections (Wittmann et al., 2019, p. 49).

In addition, Padilla discusses how a Collections as Data approach allows for new ways of understanding and deriving meaning from digital objects, including the various contexts in which they are embedded (2017, p. 1-2). Padilla also emphasizes Collections as Data's active engagement with ethical issues, such as inclusivity and bias (2016). According to Elizabeth Russey Roke, Collections as Data approaches apply the key archival principles of transparency, documentation, and provenance, which enhance understanding of digital objects (2019, July 16).

## TDL Repositories and the TCDL Proceedings

Texas Digital Library (TDL) and its member repositories use DSpace, an open source repository application. The Texas Conference on Digital Libraries (TCDL) proceedings are stored in the TDL DSpace repository. This collection includes presentations, posters and other materials from the annual Texas Conference on Digital Libraries. The proceedings were selected for use as a case study because the project contact at TDL indicated an interest in having the group work with this collection.

## Process

The task presented to our group was the following: collect metadata from the TCDL proceedings, clean the data, and analyze it computationally. We approached our Collections as Data project by dividing it into a three-part process: metadata extraction, metadata cleaning, and metadata visualization. We divided the work equally and worked together in weekly group meetings over Zoom to walk through each step of the process so we could learn together. We also broke down the written work into sections so each member could contribute equitably. Karla Roig has previous experience working with the LLILAS Benson collections on a Digital Humanities project and was able to contribute her knowledge of this subject to enhance our

understanding of steps we needed to take. In addition, Karla's sister, Aleshka Blay, provided the group with information on how to clean the data in Excel.

Note: We were able to extract metadata from the OAI-PMH website of the TDL's Dspace repository but this process is not universal. If an archive or institution does not use this protocol, or uses another Application Programming Interface (API) that can be queried to find metadata, there are many different metadata scraping tools available (please review our list of resources and tools below).

## Data Collection

The data collection process is divided into steps denoting the process we undertook to extract the desired metadata from the OAI PMH DSpace data portal. The process for extracting and cleaning the data was broken down into three main stages: downloading the XML metadata files, importing the uncleaned tabular data in Excel, and cleaning and editing the tabular data in Excel. This section shows a series of screenshots as a guide for the metadata extraction process. There are more detailed instructions in our "Step by Step Process for Extracting and Cleaning metadata" document.

Step 1: Select desired collection from the DSpace OAI-PMH portal, then select the "records" for this collection. In the List of Records page, right click on the website and select "View Page Source" to see the XML HTML source.

Step 2: Save this page onto your local drive by right clicking on it and selecting "Save as." This page only shows the first one hundred items of the collection, since the OAI-PMH portal divides it as such, so we repeated the process until we captured all the data.



Step 3: Open xml file with Notepad++, then copy and paste it into Microsoft Excel using the Text Import Wizard.

Step 4: Create a formula to add a ";" each time an element is repeated, and a "{" each time a new element begins. Once applied, this formula divides the elements into the corresponding columns. Apply this formula to all remaining cells on Column B then select the entire column and paste it into Notepad++.



Step 5: Move all individual rows into a single one by deleting all the "invisible enters," or line breaks. Then create line breaks by pressing "enter" before each iteration of "<metadata>" in order to divide the data into 100 rows, one for each complete metadata record.

Step 6: Import the data from Notepad++ into Excel once more using the Text Import Wizard, delimiting the data by the ";" that was added earlier.
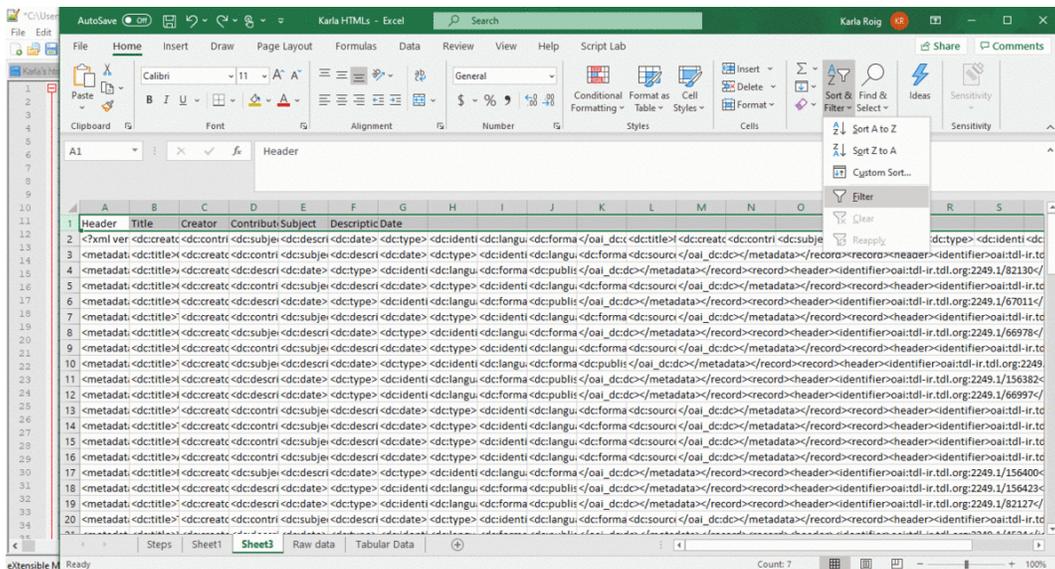


**Data Cleaning**

This section shows a series of screenshots as a guide for the metadata cleaning process. There are more detailed instructions in our "Step by Step Process for Extracting and Cleaning metadata" document.

Step 7: To begin cleaning the data, create a new row at the top of the Excel document to add the titles to each column. Apply filters to column headings in order to sort and edit multiple cells containing specific elements at the same time and start shifting all elements to the right that don't correspond to the column subject heading until every element is in the correct place.



Step 8: Continue cleaning the data by deleting the header and footer of each element (in this case each element is wrapped in Dublin Core XML) using the find and replace tool. Use this same tool to continue cleaning up additional errors in the data.

**Data Visualization: Using Voyant Tools and Tableau Public**

Data Visualization Using Voyant Tools

The following steps describe the data visualization process using the cleaned "Subject" metadata from the TCDL proceedings.

1. Go to the Voyant Tools home page.

2. Open an existing corpus, upload a file from a computer, or copy/paste text into the text box. In this case, the data was copy/pasted. Click the "Reveal" button.



3. Multiple text analyses, using multiple tools, should be displayed on the screen in panels. Click on the "?" buttons in each panel for more information about each text analysis tool. Hover your mouse over the "?" and click the windowpane icon to choose different tools for each panel.

4.  Voyant Tools' "Getting Started" page has more information about how to work with Voyant, as well as a list of text analysis tools.

<u>Data Visualization Using Tableau Public</u>

1.  Open the Tableau Public application on your desktop.



2.  Select the type of file format of your dataset and you will be prompted to select a file from your computer.

3. Navigate to the file containing your dataset on your hard drive and select "Open." This action will open the data within a Tableau Public workbook.



4. Select the worksheet you want to use on the left-hand side under the "Sheets" header and drag and drop it into the "Drag tables here" area.

5.  Open a new worksheet by selecting "Sheet 1" in the bottom left-hand corner.

6. You can create visualizations by dragging and dropping the fields on the left-hand side onto the canvas (where it says "Drag fields here") or into the "Rows" or "Columns" areas.



7. For this example, we want to create a visualization of how many different item types are in the TCDL Proceedings Collection (2007 to 2020), so we drag the "Type" field on to the canvas.

8.  Next, we entered "424" to the number of columns field because that is the amount of individual items that are in the TCDL Proceedings Collection.



9.  Once you hit "enter," Tableau Public will generate a visualization of the proportion of each item "type" relative to the whole of the collections.

10. On the right-hand side, you can select what kind of visualization you would like to use for your data. Tableau Public automatically indicates which graphical representations would be best suited for your visualization. We picked a simple "packed bubbles" representation to show the breakdown of item types within the collection. As you can see below, the majority of items in the collection are presentations, with posters as the next largest type group.



11. Login to your Tableau Public account (creating one if you need to) to access your data visualization (or "viz" as they are called by Tableau Public).

12. Download your data visualization by pressing the download icon in the bottom right-hand corner of the image.



13. Choose how you would like to download your visualization. For this example, we chose "image."

14. Now you can save the downloaded png image on your desktop for later reuse on a website or for any other use!

**Initial Steps: Starting a Collections as Data Project**

- Choose a collection you would like to explore.

- Ask the question: Why do you want to analyze this collection computationally?

- Weigh the potential costs and benefits of undertaking the project.

- Identify the format of your data (XML, TIFF, PDF, Excel, etc.) and the level of technical expertise required to work with the format. For example, working with XML files may have a steeper learning curve than working with Excel files.

- Determine whether your data needs to be cleaned.

- Select the appropriate tools to analyze your data. Different stages of analysis (collection, cleaning, etc.) may require different tools. Open access tools are preferred. Consult this document's List of Tools.

- Once the data is cleaned, you can select and download a data visualization software (for the purposes of this project, we chose to work with two popular free applications: Voyant Tools and Tableau Public).

**Recommendations**

- Start where you are and break the process into small, manageable tasks so that the work doesn't feel overwhelming.

- Preserve the contexts of data. This includes documenting, and making available to others, data's origins, as well as how data is analyzed and used over time. For example, a README plain text file records basic information about data and can help preserve its context. Recommendation no. 34, 50 Things --- Always Already Computational (Padilla et al., 2019, May 20a, p. 4).

- Provide open access to data and documentation. For example, data and documentation can be made available through a public Github repository.

- Convert your data into a format that doesn't require a lot of technical expertise to work with. This increases accessibility, long-term preservability, and institutional collaboration.

- Seek out digital humanities and digital scholarship practitioners in the field and reach out to IT professionals within your institution for collaboration. "Network with people who work with data and have the skills or knowledge you need to get your work done." Recommendation no. 18, 50 Things --- Always Already Computational (Padilla et al., 2019, May 20a, p. 3).

## Conclusion

The practical recommendations for use provided in this document are intended to be a model for TDL members interested in working with Collections as Data. The methods utilized in this project are not universally applicable. Data collection, cleaning, and visualization methods will vary according to the needs of each Collections as Data project.

## Glossary

- **Digital Humanities:**

  The field of digital humanities (DH) broadly refers to a certain methodological approach involving the application of digital technology tools to scholarship in order to analyze, visualize, interpret, and present this information in new ways. It is characterized by an emphasis on collaboration, interpretation, openness, interoperability, and innovation and can be applied to multiple disciplinary fields toward a variety of purposes. Examples of DH methods and uses include "text analysis, data mining, visualization, modeling and simulation, geospatial analysis and mapping, multimedia storytelling, information design, network analysis, interface design, and mark-up…[and applying] these tools to humanistic questions" ("About Us - Digital Humanities," n.d.).

- **Texas Digital Repository (TDR)**

"The **Texas Data Repository** is a platform for publishing and archiving datasets (and

other data products) created by faculty, staff, and students at Texas higher education

institutions. The repository is built in an open-source application called Dataverse,

developed and used by Harvard University.

The repository is hosted by the Texas Digital Library, a consortium of academic libraries

in Texas with a proven history of providing shared technology services that support

secure, reliable access to digital collections of research and scholarship. For a list of

TDL participating institutions, please visit:  http://tdl.org/member." (Mumma, C.C., n.d.).

- **Dataverse**

"Dataverse is an open source web application to share, preserve, cite, explore, and

analyze research data. It facilitates making data available to others, and allows you to

replicate others' work more easily...A Dataverse repository is the software installation,

which then hosts multiple virtual archives called Dataverses. Each dataverse contains

datasets, and each dataset contains descriptive metadata and data files (including

documentation and code that accompany the data). As an organizing method,

dataverses may also contain other dataverses" (King, G., n.d.).

- **Datasets**

A dataset is a structured collection of related data that can be analyzed computationally

(example: collection metadata in a tabular spreadsheet) ("What are the differences

between data, a dataset, and a database?" n.d.).

- **Data Visualization**

The process of creating graphical representations of data, including maps, charts, and

graphs, which can be used to visually illustrate, identify, and analyze data (Camm, 2017,

p.473).

- **DSpace**

"DSpace is an open source repository application that allows you to capture, store, index, preserve and distribute your digital material including text, video, audio and data" (Hollister, V., 2011).

- **Tableau Public**

  A free software platform for creating visual representations useful for data storytelling and analysis ("Digital Tools - Digital Humanities," 2019).

- **Voyant Tools**

  Voyant Tools is an open-source application for text editing, analysis, and visualization.

- **OpenRefine**

  An open-source application for cleaning and refining data and data sets.

## List of Tools

### Data Extraction, Editing, and Cleaning Tools

- **Notepad++**

  An open-source text and source code editor tool for MS Windows.

- **Atom**

  An open-source text and source code editor tool for Mac OS, MS Windows, and Linux.

- **Microsoft Excel**

  Spreadsheet program for organizing, analyzing, and visualizing data in tabular format. Can use to clean and sort datasets using filters, find and replace, and formula functionality.

- **OpenRefine**

  An open-source application for cleaning "messy" data and transforming large datasets into other formats (https://openrefine.org).

- **Python**

An open-source programming language that can be used for data management, analysis, and visualization. Scripted and reproducible workflows can be developed to run on large datasets (https://guides.lib.utexas.edu/data-and-donuts/schedule).

- **Beautiful Soup**

  This tool is used for scraping data found on the web. It is based on Python and it pulls the information from HTML and XML encoded data.

- **Metadata Extraction Tool**

  Open-source tool developed by National Library of New Zealand that can be used for extracting metadata from a wide range of file formats (http://meta-extractor.sourceforge.net)

- **EMET (Embedded Metadata Extraction Tool)**

  Tool for extracting metadata embedded in JPEG and TIFF files.

- **ExifTool**

  Tool for extracting and editing metadata in image, audio, and video files.

- **Metadata++**

  Freeware software for extracting and editing metadata from image, audio, video, text, and other file formats.

- **DigiPres Commons - Tools**

  A comprehensive list of links to tools for a variety of functions that can be used for your institution's CAD project undertaking on the DigiPres Commons website (see resource and annotated bibliography list below).

<div align="center">Data Visualization Tools</div>

- **Voyant Tools**

  A web-based, open-source application for text editing and analysis; useful for digital humanities scholars. (https://voyant-tools.org/docs/#!/guide/about)

- **Tableau Public**

  Free data visualization software application that enables a user to create interactive visualizations for data storytelling. (https://public.tableau.com/en-us/s/resources).

- **Gephi**

  Used for creating network visualizations of data.

- **R**

  This is an open source programming language used for visualizing data and performing statistical analysis.

- **ArcGIS**

  This is an online visualization tool for geographic information that maps, shares and analyzes data.

- **StoryMapJS**

  This is a free tool that allows you to visualize data using maps, images and textual information.

- **TimelineJS**

  This tool is also a free online website that lets you create timelines based on visual and textual data.

**Annotated Bibliography**

About Us - Digital Humanities. (2019, February 25). Retrieved November 21, 2020, from

https://dh.ucla.edu/about/

UCLA has a Digital Humanities Program and their "About Us" page offers a helpful and

condensed definition of DH practices and methods.

"All tools in detail." (n.d.). DigiPres Commons. Retrieved November 21, 2020, from

http://www.digipres.org/tools/all-tools/

List of links to a comprehensive list of community-owned digital preservation tools that

can be sorted by function and used for many of the steps needed for embarking on a

Collections as Data project (including tools for metadata extraction and scraping, data

cleaning, and data visualization, among others).

Camm, J., Fry, M., & Shaffer, J. (2017). A practitioner's guide to best practices in data

visualization. *Interfaces (Providence)*, *47*(6), 473–488.

https://doi.org/10.1287/inte.2017.0916

This resource is a tutorial about best practices for those interested in using data

visualization and introduces some of the key concepts.

Candela, G., Sáez, M. D., Escobar Esteban, Mp., & Marco-Such, M. (2020). Reusing digital

collections from GLAM institutions. *Journal of Information Science*.

https://doi.org/10.1177/0165551520950246

This article provides a guide to creating Collections as Data from different datasets of

cultural heritage institutions. The authors discuss the best practices to make data

discoverable and accessible by looking at several examples.

Clement, T., & Carter, D. (2017). Connecting theory and practice in digital humanities

information work. *Journal of the Association for Information Science and Technology*,

*68*(6), 1385–1396. https://doi.org/10.1002/asi.23732

A study of digital humanities practices and current trends in methodology, goals, and research which includes interviews and surveys of graduate programs currently using DH practices and explores their motivations for doing so.

DigiPres Commons. (n.d.). Retrieved November 21, 2020, from https://www.digipres.org/

A website containing community-owned digital preservation resources. Compiles tips for getting started, a place to submit and look up answers to digital preservation questions, tools, a guide to understanding file formats, and preservation requirements and solutions.

Digital Tools - Digital Humanities. (2019, May 15). Retrieved November 21, 2020, from https://dh.ucla.edu/digital-tools/

A list of Digital Humanities tools on the UCLA website containing descriptions and links to the websites for the tools. The page includes short descriptions and links for the following tools: OpenRefine, Salty, Vectr, Carto, Tableau, Balsamiq, and Voyant.

Hendrigan, H. (2019). Mixing digital humanities and applied science librarianship: Using Voyant Tools to reveal word patterns in faculty research. *Issues in Science and Technology Librarianship*, (91). https://doi.org/10.29173/istl3.

This article represents another case study on how Voyant Tools can be a useful tool for data visualization and further research on collections. This paper explores the use of this tool with regards to subject librarians' work.

Hollister, V. (2011, December 02). What is DSpace? - DSpace KnowledgeBase. Retrieved November 21, 2020, from https://wiki.lyrasis.org/pages/viewpage.action?pageId=25467341

A wiki page written by Valorie Hollister on the Lyrasis website introducing DSpace with information about its uses, benefits, how it works, and how to install it, as well as software documentation and other helpful information.

King, G. (n.d.). About. Retrieved November 21, 2020, from https://dataverse.org/about

The "About" page for the Dataverse Project explains what a Dataverse is, the mission and goals of the project, history, and collaborators (institutions and individuals).

Kirschenbaum, M. (2012). What is digital humanities and what's it doing in English departments? In M. K. Gold (Ed.), *Debates in the Digital Humanities* (pp. 1-2). Minneapolis, MN: University of Minnesota Press. https://dhdebates.gc.cuny.edu/read/40de72d8-f153-43fa-836b-a41d241e949c/section/f5640d43-b8eb-4d49-bc4b-eb31a16f3d06#ch01

Mumma, C. C. (n.d.). About - Texas Data Repository user documentation. Retrieved November 21, 2020, from https://texasdigitallibrary.atlassian.net/wiki/spaces/TDRUD/pages/289144946/About. TDL wiki page written by Courtney Mumma outlining reasons to deposit data into TDR, what can be deposited and what cannot be deposited in the repository, and an explanation with visuals of how dataverses work.

Neely, L., Luther, A., and Weinard, C. Cultural collections as data: Aiming for digital data literacy and tool development. *MW19: MW 2019*. Published February 1, 2019. Consulted November 18, 2020. https://mw19.mwconf.org/paper/cultural-collections-as-data-aiming-for-digital-data-literacy-and-tool-development/ This article discusses Collections as Data in a museum and cultural institutions context. It also mentions the importance of open data in collections and how having data accessible is a crucial part of working with Collections as Data.

Oglesby, N. (n.d.). Digital humanities tools and resources: Home. Retrieved November 21, 2020, from https://guides.lib.utexas.edu/digitalhumanities University of Texas Libraries' LibGuide for digital humanities tools and resources with links to useful resources.

Padilla, T. (2016). "Collections as Data: Conditions of possibility." In Collections as Data:

Stewardship and Use Models to Enhance Access.

https://www.thomaspadilla.org/2016/09/29/possibility/

Padilla's closing keynote for the 2016 Library of Congress conference *Collections as*

*Data: Stewardship and Use Models to Enhance Access* describes the ways in which

Collections as Data promotes agency, empowerment, and ethics in working with data.

Padilla, T. (2017). On a collections as data imperative. UC Santa Barbara.

https://escholarship.org/uc/item/9881c8sv

In this article, Padilla describes Collections as Data's potential to fundamentally change

the way libraries and users engage with collections. Padilla also provides a concise

definition of data in a Collections as Data context: "ordered information, stored digitally,

that are amenable to computation."

Padilla, T. (2020). Always Already Computational - Collections as Data. Retrieved November 21,

2020, from https://collectionsasdata.github.io/

Website of Always Already Computational project (which ran from 2016 to 2018) with

links to the final report and findings. The project researched approaches to applying

computational methods to collections from cultural heritage institutions.

Padilla, T., Allen, L., Frost, H., Potvin, S., Russey Roke, E., & Varner, S. (2019, May 20a). 50

Things --- Always already computational: Collections as data. Zenodo.

http://doi.org/10.5281/zenodo.3066237

This document is intended to provide practitioners with initial steps to take in order to

begin supporting Collections as Data at cultural heritage institutions. The

recommendations provide a starting point for institutions such as TDL to develop their

own recommendations for working with Collections as Data.

Padilla, T., Allen, L., Frost, H., Potvin, S., Russey Roke, E., & Varner, S. (2019, May 20b). Santa

      Barbara statement on collections as data --- Always already computational: Collections

      as data (Version 2). Zenodo. http://doi.org/10.5281/zenodo.3066209

      The Santa Barbara statement is one of the deliverables of the Always Already

      Computational project. The document outlines ten foundational Collections as Data

      principles, including commitments to shared documentation, open access, and

      transparency.

Padilla, T., Allen, L., Frost, H., Potvin, S., Russey Roke, E., & Varner, S. (2019, May 22). Final

      report --- Always already computational: Collections as data (Version 1). Zenodo.

      http://doi.org/10.5281/zenodo.3152935

      The *Always Already Computational: Collections as Data* project's final report

      summarizes the project's Collections as Data framework, impacts, findings, and areas

      for further investigation. The final report's appendices provide the project deliverables,

      including the 50 recommendations document and the Santa Barbara statement on

      Collections as Data.

Padilla, T., Allen, L., Varner, S., Potvin, S., Russey Roke, E., Frost, H., & Heppler, J. (2017,

      March 09). Serendipitous Collections as Data. Retrieved November 21, 2020, from

      https://collectionsasdata.github.io/serendipity/

      "Serendipitous Collections as Data" is a resource from the Already Always

      Computational site. It is a tool that randomly generates ideas (in the form of quotes)

      about collections as data and one can refresh the page to view a new idea. The ideas

      are drawn from partners of the project, the Collections as Data National Forum writings,

      and collections as data position statements from the site. This tool was coded by Jason

      Heppler.

Paradise, L. (2015). When you find out what digital humanities is, will you tell me? *The Serials*

      *Librarian 69*(2), 194-203. https://doi.org/10.1080/0361526X.2015.1036198

An essay that acknowledges the nebulous definitions of DH as a field and the confusion that can cause, while attempting to define what it means in terms of historical research. The essay covers digital archives for born digital content, the digital presentation of scholarly content, and working with historical materials in a digital environment (as well as the controversy that can cause).

Rethinking Collections as Data. (2019). In *Open a GLAM Lab* (1st ed.).

https://doi.org/10.21428/16ac48ec.f54af6ae

This website provides an introduction on how to assess new collections, from thinking about copyright, accessibility, and reusability. It also explains some useful concepts regarding the description of collections.

Russey Roke, E. (2019, July 16). Collections as data. *bloggERS*.

https://saaers.wordpress.com/2019/07/16/collections-as-data/

Russey Roke is a co-investigator of the *Collections as Data: Always Already Computational* project. In this blog post, Russey Roke provides a brief overview of Collections as Data. She states that Collections as Data approaches are archival and enhance access to archival material, and advises practitioners to "just start" working with their collections as data.

Shensky, M. (2020). Data & Donuts: Schedule. Retrieved November 21, 2020, from

https://guides.lib.utexas.edu/data-and-donuts/schedule

University of Texas Libraries LibGuide page for the Data & Donuts series, which features links to recordings of webinars on data-related topics including using Python for data visualization, mapping data with R, and GIS operations. Compiled by Michael Shensky.

Sinclair, S., & Rockwell, G. (2016). Voyant Tools. Retrieved November 21, 2020, from

https://voyant-tools.org/docs/

Open-source software for text editing, analysis, and visualization that can be used for DH data projects.

Women Techmakers (Directors). (2018, October 24). What are APIs and how to use them in 60

    seconds! [Video file]. Retrieved November 21, 2020, from

    https://www.youtube.com/watch?v=BRAMRUxfiXY

    Video from the Women Techmakers group introducing the concept of application

    programming interfaces (APIs) and how they are used.

Welcome! (n.d.). Retrieved November 21, 2020, from https://openrefine.org/

    Official website for OpenRefine tool for cleaning and transforming data with videos

    covering the following topics: "Explore Data," "Clean and Transform Data," and

    "Reconcile and Match Data." Website also includes the software documentation.

What are the differences between data, a dataset, and a database? (n.d.). Retrieved November

    21, 2020, from

    https://www.usgs.gov/faqs/what-are-differences-between-data-a-dataset-and-a-database

    ?qt-news_science_products=0

    United States Geological Survey (USGS) webpage that contains a helpful description of

    the distinction between data, datasets, and databases.

Wittmann, R., Neatrour, A., Cummings, R., & Myntti, J. (2019). From digital library to open

    datasets: Embracing a "collections as data" framework. *Information Technology and*

    *Libraries*, *38*(4), 49-61. https://doi.org/10.6017/ital.v38i4.11101

    Wittmann et al. speak of a "collections as data mindset," which is a helpful way of

    conceptualizing Collections as Data (p. 51). The authors examine five case studies that

    use a variety of Collections as Data approaches, such as text mining, GIS, and topic

    modeling.