# Teaching Computational Archival Science (CAS):
## *Context, Pedagogy, and Future Directions in Metadata*

Richard **MARCIANO**
University of Maryland
marciano@umd.edu

**A**dvanced **I**nformation **C**ollaboratory
https://ai-collaboratory.net/cas/

October 21, 2024

**Dublin Core™ Metadata Initiative (DCMI) 2024**:
*Metadata Innovation: Trust, Transformation, and Humanity*

# This talk builds on an international collaboration: US/Canada

## *Teaching Computational Archival Science: Context, Pedagogy & Future Directions*

Victoria L. Lemieux (U. British Columbia, CANADA -- School of Information)

Richard Marciano (U. Maryland, USA -- College of Information)

Submitted to **iConference 2025**

*https://ai-collaboratory.net/wp-content/uploads/2024/09/Teaching-Computational-Archival-Science-Context-Pedagogy-and-Future-Directions-2024-submitted.pdf*

Professor
MAS Program Chair
Blockchain@UBC Cluster Lead

Where:  We illustrate the introduction of CAS into graduate archival training: *case studies*
         a. University of Maryland
         b. University of British Columbia
         c. Building and Sustaining Educator Networks

We propose CAS graduate competencies

## Computational Training Initiatives:

## Rest of Talk:
Metadata to support GenAI & LLMs

**U.S.:**
- I. TALENT Network (2018-2026):
  - 10 INFO Schools
- II. LEADING Network (2018-2025):
  - Trained 89 Fellows
  - Relied on 26 Mentors
- III. DCIP Certificate Program (2019-2024):
  - 60+ Students
- IV. DCIC Center (2015-2020):
  - 300 Students

**INTERNATIONAL:**
- V. Advanced Information Collaboratory
  - AIC (2020-present):
    - 50+ Partners

- **OAIS & Computational Archival Processing**
- **Computational Archival Science** (CAS)
- **GraphRAG** (Retrieval Augmented Generation with Graphs)

# I. TALENT Network

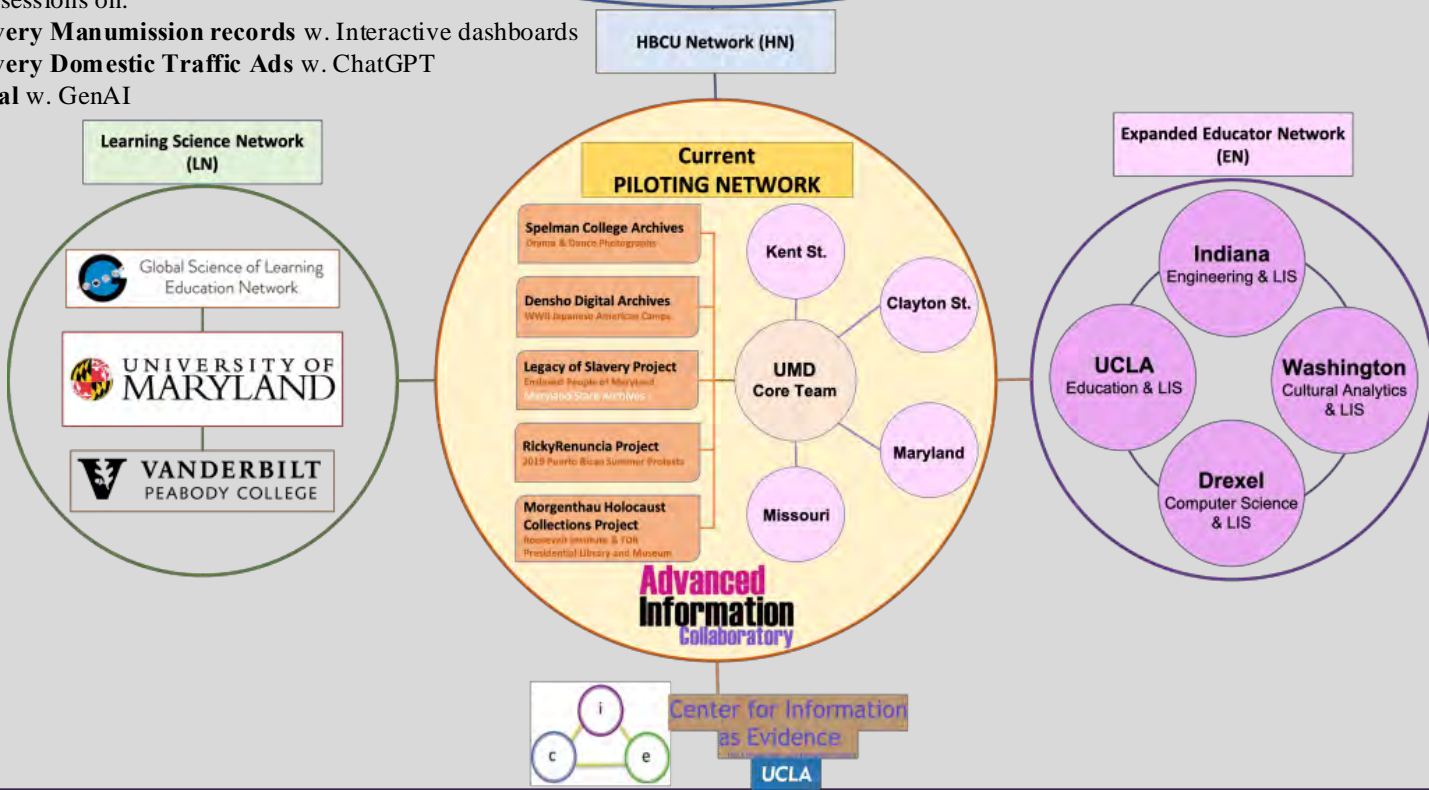*Training of Archival & Library Educators w. iNnovative Technologies*

"Modernizing **MLIS** Teaching"

**Spelman College Fall DATATHON** (Nov. 1, 2024)
Computational Treatments of the Legacies of Black History
Featuring 3 one-hour sessions on:
- **Legacy of Slavery Manumission records** w. Interactive dashboards
- **Legacy of Slavery Domestic Traffic Ads** w. ChatGPT
- **Urban Renewal** w. GenAI

Spelman College
A Choice to Change the World

Georgia Tech Library

CLARK ATLANTA UNIVERSITY

INSTITUTE of Museum and Library SERVICES

HBCU Network (HN)

**Learning Science Network (LN)**

Global Science of Learning Education Network

UNIVERSITY OF MARYLAND

VANDERBILT PEABODY COLLEGE

**Current PILOTING NETWORK**

Spelman College Archives
Drama & Dance Photographs

Densho Digital Archives
WWII Japanese American Camps

Legacy of Slavery Project
Enslaved People of Maryland
Maryland State Archives

RickyRenuncia Project
2019 Puerto Rican Summer Protests

Morgenthau Holocaust Collections Project
Roosevelt Institute & FDR Presidential Library and Museum

Kent St.

Clayton St.

**UMD Core Team**

Maryland

Missouri

**Advanced Information Collaboratory**

**Expanded Educator Network (EN)**

Indiana
Engineering & LIS

UCLA
Education & LIS

Washington
Cultural Analytics & LIS

Drexel
Computer Science & LIS

i c e

Center for Information as Evidence

UCLA

# 1. IMLS 2018-2020: CT-LASER

**Developing a Computational Framework for Library and Archival Education** ($100K)
**Planning grant.**
https://ai-collaboratory.net/wp-content/uploads/2020/11/Final_Report_r.pdf

# 2. IMLS 2020-2024: PILOTING Network

**Pilot Study with 4 US iSchools & 5 Archives** ($300K)
**Piloting grant.**
https://ai-collaboratory.net/projects/piloting-network/

# 3. IMLS 2022-2025: TALENT Network

**Promoting the Training of Archival & Library Educators w. iNnovative Technologies** ($400K)
**National implementation grant.**
https://ai-collaboratory.net/projects/talent-network/

# 4. IMLS 2024-2026: GenAI-4-Arch

Harnessing **Generative AI** to Support Exploration and Discovery in Library and **Archival Collections** ($194K)
https://www.imls.gov/grants/awarded/lg-256565-ols-24



**HERITAGE-AI** Initiative:

**H**arnessing **E**nhanced **R**esearch & **I**nstructional **T**echnologies for **A**rchival **G**enerative **E**xploration - using **AI**

https://heritage-ai.org/

INSTITUTE *of*
**Museum** and **Library**
SERVICES

**II. LEADING**: LIS Education And Data Science Integrated Network Group *[IMLS Drexel: 2020-2025]* **($887K)**
**LEADS**: Library Education And Data Science *[IMLS Drexel: 2017-2020]* **($313K)**

- 14 core team members Drexel iSchool faculty and partners leads UCSD, OCLC, and Univ. of New Mexico
- 26 mentors from leading libraries, archives, data/research centers and agencies across the U.S.

https://mrc.cci.drexel.edu/leading/

"Enhancing **Doctoral** and **Early Career** Professional Training in Data Science"

→ 89 Fellows

**Jane Greenberg**
Professor @ Drexel U.

**Director** @ Metadata
Research Center

**Co-Chair** @ DCMI 2024
Invited Panels & Talks

# CURRENT / RECENT PROJECTS

**EDUCATION:** $2.2M
- Developing a Computational Framework for Library and Archival Education *[IMLS: 2018-2020]*
- Piloting a Collaborative Network for Integrating CT into Library and Archival Ed. and Practice *[IMLS: 2020-2024]*
- Training of Archival & Library Educators w. iNnovative Technologies *[IMLS: 2022-2023]*
- LIS Education And Data Science Integrated Network Group, LEADING *[IMLS Drexel: 2020-2025]* ($887K)
- Library Education And Data Science, LEADS *[IMLS Drexel: 2017-2020]* ($313K)

**INFRASTRUCTURE:** $14M
- WIN: a Window Into Neuroregulation *[NSF Convergence: 2019-2024]*
- Developing a Digital Asset Management System for Additive Manufacturing *[ARL: 2020-2024]*
- Developing a Digital Asset Management System for the Mary McLeod Bethune Historic Site *[NPS: 2019-2022]*
- Improving Fedora 4 to Work with Web-Scale Storage and Services, DRASTIC *[IMLS: 2017-2020]*
- Brown Dog: "Making Sense of Billion-Record Archives" with the NCSA) *[NSF: 2013-2018]*

**SOCIAL JUSTICE, HUMAN RIGHTS, CULTURAL HERITAGE:** $3M
- Using AI and ML to Optimize Information Discovery in Under-utilized, Holocaust-related Records *[Kurtz Foundation]*
- Harnessing Generative AI to Support Exploration and Discovery in Library and Archival Collections *[IMLS proposal]*
- International Research Portal for Holocaust-Era Cultural Property *[Kurtz Foundation]*
- Measuring the Impact of Urban Renewal *[NSF]*
- Computational Thinking to Unlock the Japanese American WWII Camp Experience *[UMD-FIA]*
- Computational Treatments to re-member the Legacy of Slavery (CT-LOS) *[Kurtz Foundation]*
- Testbed for the Redlining Archives of California's Exclusionary Spaces (T-RACES) *[IMLS]*
- Mapping Inequality – Redlining in New Deal America *[U. Richmond Mellon]*

# III. DCIP: Digital Curation for **Information Professionals**

**Certificate Program**

| | |
|---|---|
| January 8 – February 16, 2024 | Introduction to Digital Curation (6 weeks) |
| February 26 – May 17, 2024 | Tools and Software for Digital Curation (12 weeks) |
| May 27 – August 16, 2024 | Implementing Digital Curation in the Workplace (12 weeks) |

INSTITUTE *of* **Museum**and**Library** SERVICES

**Over 60 projects**

- https://ischool.umd.edu/academics/certificates-non-degree-study/digital-curation-for-information-professionals-certificate/

**Examples of Student Projects:**

- 2024: https://ai-collaboratory.net/2024/09/03/2024-dcip-cohort-presents-their-capstone-projects/

- 2023: https://ai-collaboratory.net/2023/08/14/2023-dcip-cohort-presents-their-capstone-projects/
    - *Synthetic data and generative AI: an interactive learning experience*
    - *Digital preservation of legacy file formats*
    - *Buried While Black – Payne Cemetery: The disinterment of a historically black cemetery in Washington, DC*
    - *Exploring Indexing Methods for Handwritten Text*
    - *Columbia MD Archives digital curation manual*

- 2022: https://ai-collaboratory.net/2022/09/07/2022-dcip-cohort-presents-their-capstone-projects/

- 2021: https://ai-collaboratory.net/2021/08/23/dcip-cohort-2021-presents-capstone-projects/

# CAS
## Computational Archival Science

The DCIC is pioneering advances in *computational treatments of archival and cultural content.*

See our CAS portal for the latest developments:
http://dcicblog.umd.edu/cas/

## What is CAS?

An interdisciplinary field concerned with the application of computational methods and resources to large-scale rec-ords /archives processing, analysis, stor-age, long-term preservation, and ac-cess, with the aim of improving efficien-cy, productivity and precision in support of appraisal, arrangement and descrip-tion, preservation and access decisions, and engaging and undertaking re-search with archival materials.

## CAS Founding Partners:

**Richard Marciano,** U. Maryland
**Mark Hedges,** King's College London (UK)
**Vicki Lemieux,** U. British Columbia (Canada)
**Maria Esteva,** Texas Advanced Computing Center (TACC)
**Michael Kurtz,** U. Maryland
**Bill Underwood,** U. Maryland
**Greg Jansen,** U. Maryland
**Mark Conrad,** National Archives and Records Administration (NARA)

## curatelab
**Hornbake South 4110**

Digital lab for group learning, collaborative design, and hands-on digital curation project development (23 seats, 3 interactive screens, 12 workstations with 12TB of storage).

## digitizationlab
**Hornbake South 4110D**

Document scanning, image manipulation, and archival ingestion facility for group projects.

## serverfarm
**Atlantic Building**

On-campus virtual machine farm for research data processing, storage, and hosting (15TB storage, 2 Dell servers, VMWare-powered).

## cloudlab
**Amazon Cloud**

Dashboard-enabled virtual computing lab in the cloud for creating Windows/Ubuntu instances using Amazon Web Services (AWS).

## dataCave
**UMD Cyberinfrastructure Center at the Rivertech Bldg**

## DRAS·TIC

**D**igital **R**epository **A**t **S**cale **T**hat **I**nvites **C**omputa-tion (**T**o **I**mprove **C**ollections): a petascale archival storage and preservation repository (based on the **DRAS-TIC** open-source software [NoSQL Cassandra database] and computational infrastructure (4 Dell nodes).

# dcic
## digital curation
## innovation center

http://dcic.umd.edu

## Mission:

Be a leader in the digital curation research and edu-cational fields, and foster interdisciplinary collabora-tions using Big Records and Archival Analytics with public / industry / government partnerships.

## Goals:

Sponsor interdisciplinary projects that explore the integration of archival research data, user-contributed data, and technology to generate new forms of analysis and historical research engagement, particularly in the arenas of social justice, human rights, and cultural heritage.

## Motto:

*"Integrating Education and Research"*

COLLEGE OF **INFORMATION** STUDIES

Computational archival science is a blend of: (1) computational & (2) archival thinking.

**DP**

**MS**

**CPS**

**ST**

## Data Practices

## Modeling & Simulation Practices

## Computational Problem Solving Practices

## Systems Thinking Practices

| Data Practices | Modeling & Simulation Practices | Computational Problem Solving Practices | Systems Thinking Practices |
|---|---|---|---|
| Collecting Data | Using Computational Models to Understand a Concept | Preparing Problems for Computational Solutions | Investigating a Complex System as a Whole |
| Creating Data | Using Computational Models to Find and Test Solutions | Programming | Understanding the Relationships within a System |
| Manipulating Data | Assessing Computational Models | Choosing Effective Computational Tools | Thinking in Levels |
| Analyzing Data | Designing Computational Models | Assessing Different Approaches/Solutions to a Problem | Communicating Information about a System |
| Visualizing Data | Constructing Computational Models | Developing Modular Computational Solutions | Defining Systems and Managing Complexity |
| | | Creating Computational Abstractions | |
| | | Troubleshooting and Debugging | |

**Example 1**: Fall 2019

# Experiential, Interdisciplinary & Team-based Learning:
## Computational Thinking in Archives

**"Reframing Digital Curation Practices through a Computational Thinking Framework"**

Richard Marciano, et al., 2019 IEEE International Conference on Big Data, 4th CAS Workshop, Dec. 11, 2019, Los Angeles, CA.

https://ai-collaboratory.net/wp-content/uploads/2020/04/ReframingDC-UsingCT_final.pdf

# October 28-29, 2019: "Datathon" at the Maryland State Archives

**https://ai-collaboratory.net/projects/legacy-of-slavery/student-led-datathon-at-the-maryland-state-archives/**



[Runaway Slave Ads]

[Certificates of Freedom]

[Manumissions]

[Cemetery Records]

[LEADS Fellows]

- **Paper**: Gnanasekaran, R.K. and Marciano, R., (2021). *Piloting Data Science Learning Platforms through the Development of Cloud-based interactive Digital Computational Notebooks.* https://ai-collaboratory.net/wp-content/uploads/2021/10/ISGC2021_Gnanasekaran_Marciano.pdf.

- **Video**: https://www.youtube.com/watch?v=cNBc0AY-r-k

- **Jupyter Notebook**: https://cases.umd.edu/github/cases-umd/Legacy-of-Slavery/blob/master/index.ipynb

# (Jupyter) Digital Notebooks

Educators are rapidly adopting
Jupyter Notebooks for:

    * teaching

    * use in the classroom

    * developing teaching materials

    * creating computational stories

See: https://cases.umd.edu

# Historical Lab Notebooks



Paper-based Lab Notebooks:
- Used in science research
- Represent a record of:
  - observations
  - experiments
  - ideas
  - notes
  - formulas
  - data

Electronic Lab Notebooks:
- patient medical records

# 1. Developing Name Registries



Creating a Names Registry

| Data Practices | Modeling & Simulation Practices | Computational Problem Solving Practices | Systems Thinking Practices |
|---|---|---|---|
| Collecting Data | Using Computational Models to Understand a Concept | Preparing Problems for Computational Solutions | Investigating a Complex System as a Whole |
| Creating Data | Using Computational Models to Find and Test Solutions | Programming | Understanding the Relationships within a System |
| Manipulating Data | Assessing Computational Models | Choosing Effective Computational Tools | Thinking in Levels |
| Analyzing Data | Designing Computational Models | Assessing Different Approaches/Solutions to a Problem | Communicating Information about a System |
| Visualizing Data | Constructing Computational Models | Developing Modular Computational Solutions | Defining Systems and Managing Complexity |
| | | Creating Computational Abstractions | |
| | | Troubleshooting and Debugging | |

# 2. Integrating Vital Records



Consolidating and Analyzing Vital Records

| Data Practices | Modeling & Simulation Practices | Computational Problem Solving Practices | Systems Thinking Practices |
|---|---|---|---|
| Collecting Data | Using Computational Models to Understand a Concept | Preparing Problems for Computational Solutions | Investigating a Complex System as a Whole |
| Creating Data | Using Computational Models to Find and Test Solutions | Programming | Understanding the Relationships within a System |
| Manipulating Data | Assessing Computational Models | Choosing Effective Computational Tools | Thinking in Levels |
| Analyzing Data | Designing Computational Models | Assessing Different Approaches/Solutions to a Problem | Communicating Information about a System |
| Visualizing Data | Constructing Computational Models | Developing Modular Computational Solutions | Defining Systems and Managing Complexity |
| | | Creating Computational Abstractions | |
| | | Troubleshooting and Debugging | |

# 3. Designing Controlled Vocabularies



Seeing the Forest and the Trees: Creating Controlled Vocabularies for Historical Collections

| Data Practices | Modeling & Simulation Practices | Computational Problem Solving Practices | Systems Thinking Practices |
|---|---|---|---|
| Collecting Data | Using Computational Models to Understand a Concept | Preparing Problems for Computational Solutions | Investigating a Complex System as a Whole |
| Creating Data | Using Computational Models to Find and Test Solutions | Programming | Understanding the Relationships within a System |
| Manipulating Data | Assessing Computational Models | Choosing Effective Computational Tools | Thinking in Levels |
| Analyzing Data | Designing Computational Models | Assessing Different Approaches/Solutions to a Problem | Communicating Information about a System |
| Visualizing Data | Constructing Computational Models | Developing Modular Computational Solutions | Defining Systems and Managing Complexity |
| | | Creating Computational Abstractions | |
| | | Troubleshooting and Debugging | |

# 4. Mapping Events and People



Mapped Narratives of Resistance in Tule Lake

| Data Practices | Modeling & Simulation Practices | Computational Problem Solving Practices | Systems Thinking Practices |
|---|---|---|---|
| Collecting Data | Using Computational Models to Understand a Concept | Preparing Problems for Computational Solutions | Investigating a Complex System as a Whole |
| Creating Data | Using Computational Models to Find and Test Solutions | Programming | Understanding the Relationships within a System |
| Manipulating Data | Assessing Computational Models | Choosing Effective Computational Tools | Thinking in Levels |
| Analyzing Data | Designing Computational Models | Assessing Different Approaches/Solutions to a Problem | Communicating Information about a System |
| Visualizing Data | Constructing Computational Models | Developing Modular Computational Solutions | Defining Systems and Managing Complexity |
| | | Creating Computational Abstractions | |
| | | Troubleshooting and Debugging | |

# 5. Connecting Events and People



Resistance at Tule Lake: Connecting Events and People through Networks to Understand Japanese American Incarceration During World War II

| Data Practices | Modeling & Simulation Practices | Computational Problem Solving Practices | Systems Thinking Practices |
|---|---|---|---|
| Collecting Data | Using Computational Models to Understand a Concept | Preparing Problems for Computational Solutions | Investigating a Complex System as a Whole |
| Creating Data | Using Computational Models to Find and Test Solutions | Programming | Understanding the Relationships within a System |
| Manipulating Data | Assessing Computational Models | Choosing Effective Computational Tools | Thinking in Levels |
| Analyzing Data | Designing Computational Models | Assessing Different Approaches/Solutions to a Problem | Communicating Information about a System |
| Visualizing Data | Constructing Computational Models | Developing Modular Computational Solutions | Defining Systems and Managing Complexity |
| | | Creating Computational Abstractions | |
| | | Troubleshooting and Debugging | |

DEVELOPING NAME REGISTRIES

Andy Jose SILVA / Emery PATTERSON, Mary McKINLEY
2 MLIS

INTEGRATING VITAL RECORDS

James SANTOS, Genevieve KOCIENDA / Kanishka JAIN
HdLS        MLIS        MIM

DESIGNING CONTROLLED VOCABULARIES

Tahira TUBABI / Margaret Rose HUNT, Hannah FRISCH, Hilary Iris Yin SHUE
InfoSci / Japanese        2 MLIS

MAPPING EVENTS & PEOPLE

Connor MULLANE, Brittan SCHAMS, Marielle VANNELLI / Chenxi LIU, Jade XU
2 MLIS        2 BCIM

CONNECTING EVENTS & PEOPLE

Sarah AGARWAL / Hannah KRAUSS / Danish MIR, Mrinlok SURI / Debindran PRADHAN
CS / Hindi        MLIS        2 MIM        BCIM

**EVENT:** Resistance at Tule Lake: A Conversation with the Filmmaker and iSchool Digital Curators (and Film Viewing)
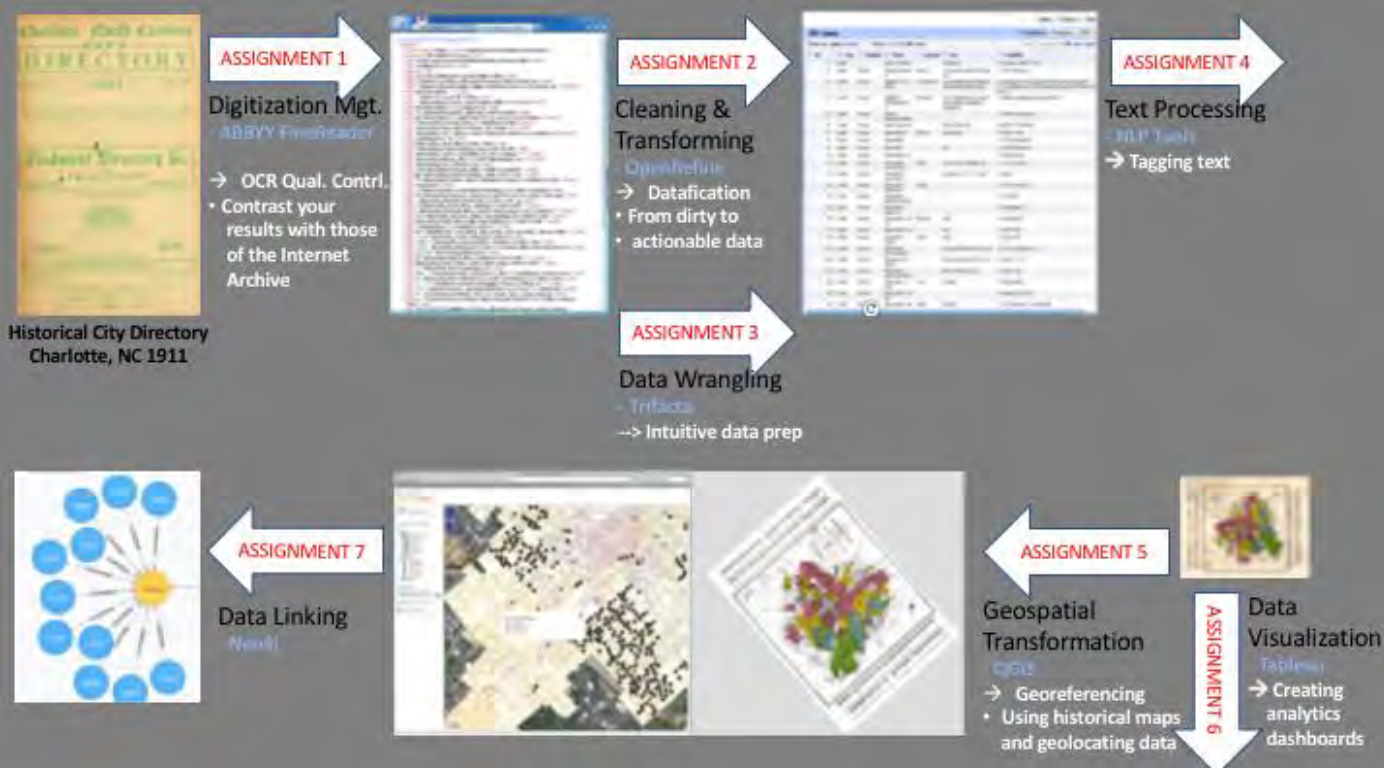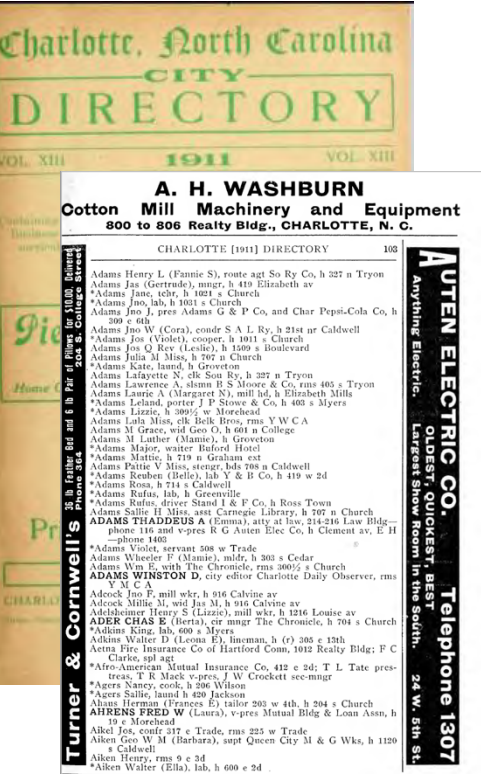
**PROJECT:**

**Computational Thinking to Unlock the Japanese American WWII Camp Experience (w. Densho.org)**

https://ai-collaboratory.net/projects/ct-ja_ww2_camps/

**Example 2**: Spring 2023 Graduate Course in the MLIS Program – Implementing Digital Curation

- *Teaching and Learning with Archival Materials through the Development of Interactive Computational Notebooks"*,
  P. Piety, M. Conrad, R. Marciano, I. Cornfield, E. Dallimore, R. Fettig, E. Hansen, H. Kemp, T. Turabi (2023).
  Chapter Submission for the 2023 Archives and Primary Source Handbook, peer-reviewed open-access NewPrairiePress textbook.
  Link: https://ai-collaboratory.net/wp-content/uploads/2023/10/Piety_Conrad_Marciano_et_al-FINAL.pdf
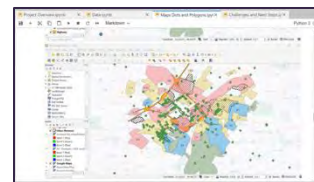
**Table 2 - Comparison of student final projects**

| Student Project Title | CLD | CRD | MD | AD | VD | UaC | F&T | ACM | DCM | CCM | PCS | PRG | CET | ADS | DMS | CCA | TAD | ISW | URS | TIL | CIS | DSC | OCR | Tableau | Excel | Neo4J | Open Refine | QGIS | Pandas | Code Boxes | Dataset Scale | Use of Code |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **A. Expanding the Network: Modeling Relationships with Neo4j** | | | X | | | | | | X | X | | | | | | | X | | X | | | | | | X | X | | | | No | Small | 1 |
| **B. Revisualizing Geographic Disparities: Examining Trends in Racial and Economic Inequality on the Streets w/o GIS** | | | X | X | X | X | | | | | X | | | X | | | | | X | X | | | | | | | X | | | No | Full | 1 |
| **C. Mad or Madam: Investigating an Undefined Data Term** | X | X | X | X | | X | | | | | | | | | | | | | | | | | | X | | | X | | | No | Full | 1 |
| **D. Race, Marriage, and Profession: Data at Scale Test Case** | | | X | X | X | X | | | | | | | | | | X | | | X | X | | | | X | | | X | | | Tableau | Full | 2 |
| **E. Building a Bigger Picture: A Case Study of Combining the General City and Business Directories** | X | | X | X | X | | | | | | | | | | | | X | | X | X | | | X | X | | | | | | Tableau | Full | 2 |
| **F. Religious Life in 1911 Charlotte, NC** | | X | X | X | X | | | | | | | | | | | | X | | | | | | | X | | | X | | | Tableau | Full | 2 |
| **G. Gender, Race, and Archival Silences** | | X | X | X | X | | | | | | | | | | | | | | | | | | | X | | | X | | | Tableau | Full | 2 |
| **H. Conceptualizing Prosperity: A Case Study Analyzing Housing through Job Types** | | | X | X | X | | | | | | | | | | | | | X | X | X | X | | | X | | | X | | | Tableau | Full | 2 |
| **I. Visualizing Neighborhood Demographics** | | | X | X | X | | | | | | | | | X | | | | | X | X | | | | X | | | | X | X | Tableau | Full | 2 |
| **J. Mapping Over Time in Charlotte NC: Population, Redlining, and Urban Renewal** | | | X | X | X | X | | | X | X | X | | X | | | | | X | X | X | X | X | | X | | | X | | X | Python / Pandas Lib. | Full | 3 |

# The Future is Here & *Now*

*- Fundamental changes to the way we acquire, manage, and present cultural collections*
*- COVID-19: impact on access to archives, libraries, museums*

## HERE:

- **US: OMB M-23-07 Update to Transition to Electronic Records**
  - **By June 30, 2024,** Federal agencies must manage all permanent records in an electronic format for eventual transfer and accessioning by NARA.
  - **After June 30, 2024**, Federal agencies must transfer all permanent records in an electronic format and with appropriate metadata.
- **NARA Budget Estimates 622 Years to Process Backlogged FOIA Requests at Just Two Presidential Libraries**
  - A "FOIA backlog of an estimated 183 million pages at the George W. Bush Library and [a] 128-million-page backlog at the Barack Obama Library" alone. At current rates, it would take NARA 622 years to declassify the pending declassification requests just these two presidential libraries.
  - NARA presently holds 13.5 billion pages, **only about 2% of which is in digitized form**. This is just one among many thousands of archival repositories; with more than 25,000 such repositories in just the United States.

## NOW:

- **"An IMPERATIVE for educating the archivists and records managers of the future for the digital world"** *[Mark Conrad]*

# LAUNCH OF THE AIC COLLABORATORY

https://ai-collaboratory.net

V.

**Advanced Information Collaboratory**

The **AIC** was launched at The Alan Turing Institute in London, UK on Jan. 20, 2020. It brings together partners from leading academic and cultural institutions from six continents. Its goals are to:

1.  **EXPLORE** the opportunities and challenges of "disruptive technologies" for archives and records management (digital curation, machine learning, AI, etc.).

2.  **PURSUE** multidisciplinary collaborations to share relevant knowledge across domains.

3.  **LEVERAGE** the latest technologies to unlock the hidden information in massive stores of records.

4.  **TRAIN** current and future generations of information professionals to think computationally and rapidly adapt new technologies to meet their increasingly large and complex workloads.

5.  **PROMOTE** ethical information access and use.

# AIC Founding Partners:



**Dr. Richard Marciano**
Professor
UMD iSchool (US)

**Mark Conrad**
Archives Specialist
NARA - former (US)

**Dr. Eirini Goudourali**
Head of Dig. Research Progs.
TNA (UK)

**Dr. Jane Greenberg**
Professor
Dir. Metadata Research Center
Drexel U. (US)

**Dr. Mark Hedges**
Dep. of Dig. Hum.
King's College London (UK)

**Greg Jansen**
Senior Res. Soft. Architect
UMD iSchool (US)

**Dr. Michael Kurtz**
Asst. Archivist for Rec. Services
NARA - former (US)

**Dr. Victoria Lemieux**
Assoc. Professor
Blockchain@UBC Cluster
lead (Canada)

**Dr. Bill Underwood**
Res. Scientist
GTRI Res. Sci (former)
UMD iSchool (US)

**Dr. Lyneise Williams**
Associate Prof. Art History
Founder VERA Collaborative
UNC Chapel Hill (US)

Home

## Advanced Information Collaboratory

## North America:
- **MEDiAL Lab @U. Maryland:** Dr. Phil Piety
- **NARA (former):** Bruce Ambacher
- **UCLA:** Dr. Anne Gilliland
- **Kent State U.:** Dr. Karen Gracy
- **U. Missouri:** Dr. Sarah Buchanan
- **Clayton State U.:** Dr. Joshua Kitchens
- **The Smithsonian Institutions (NMAH):** Bob Horton
- **Harvard Library:** Ceilyn Boyd
- **UC Santa Barbara:** Marisol Ramos
- **UC San Diego:** Dr. Andrea Chiba
- **US Holocaust Memorial Museum:** Michael Levy
- **Densho.org:** Geoff Froh
- **Maryland State Archives:** Chris Haley & Maya Davis
- **Spelman College:** Holly Smith
- **Puerto Rican Spring Project:** Marison Ramos, Irmarie Fraticelli, Joel Blanco

## South America:
- **U. de Brasilia:** Dr. Cláudio Gottschalg-Duque

## UK:
- **Loughborough U.:** Lise Jaillant
- **The Alan Turing Institute:** David Beavan
- **UK TNA:** Pip Wilcox, Mark Bell, Paul Young, Jenny Bunn & Sonia Ranade
- **Oxford U.:** Dr. David De Roure
- **European Holocaust Research Infrastructure:** Dr. Reto Speck

## Europe:
- **Hamburg U. Archives:** Francesco Gelati
- **University of Amsterdam:** Dr. Tobias Blanke
- **INESC-ID Portugal:** Dr. Diogo Proença

## Africa:
- **U. South Africa:** Dr. Shadrack Katuu

## Asia:
- **Central U. of Gujarat (India):** Dr. Bhakti Gala
- **Centre for Dev. of Advanced Computation (India):** Dr. D. Katre
- **Indian Inst. of Management:** Dr. H. Anil Kumar
- **Kyushu U. (Japan):** Dr. Yoichi Tomiura & Dr. Emi Ishita

## Australia:
- **U. Canberra:** Dr. Tim Sherratt

|  | I. | II. | III. | IV. | V. |
|---|---|---|---|---|---|
| **LTAR PROJECTS** | WWII Japanese American Incarceration | Legacy of Slavery | American Responses to the Holocaust | Urban Renewal | Redlining |
| **PARTNERS** | DENSHŌ densho.org | LEGACY OF SLAVERY | FRANKLIN D. ROOSEVELT PRESIDENTIAL LIBRARY | Remapping Southside | T-RACES |
|  | 10 years | 7 years | 4 years | 11 years | 25 years |

**Example:**

- **SUMMER 2024: MLIS Students Engage w. Innovative Technologies to Explore the Future Processing of Archival Collections through *spatial, graph & genAI* Techniques:** https://ai-collaboratory.net/2024/07/09/summer2024/

- **Victoria Lemieux @ UBC: "**Archival Competencies Framework for Training in AI/ML"
  (using concepts of trustworthiness and authenticity of records)

# How Do Computational Processes Touch Upon Archival Work?

**Preparing Archivists in Computational Thinking & Innovative Technologies**
https://www.youtube.com/watch?v=ZdLJHQLbR4k

**Anne Gilliland (UCLA)** & Richard Marciano (UMD),

International Conference of **"Technology, Society, Humanities:
Digital Intelligence Empowers the Modernization of Archival Work"**

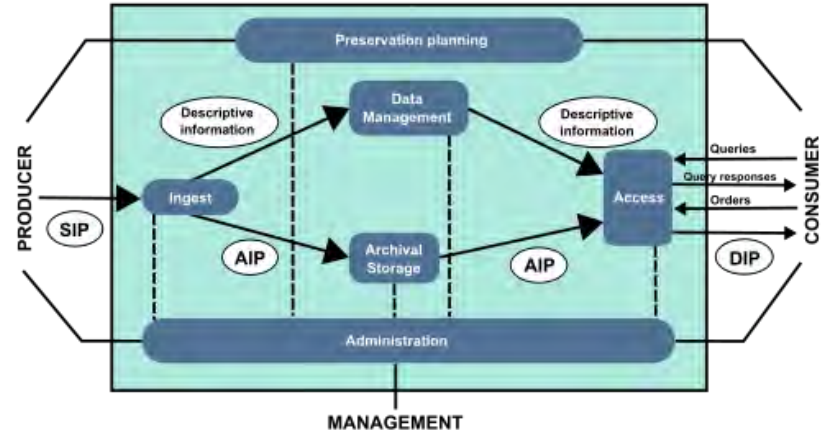Shandong U., China. Oct. 26, 2023.

# Open Archival Information System (OAIS)

# On the producer or creator side:

Computational processing is used by bureaucracies and research:

- to create/gather and analyze data
- in ways that should but may not always generate records for evidentiary, accountability and transparency purposes
- to manage active data and records

Archivists and records managers may be responsible for ensuring that:

- accountable and transparent records are created and are archivable or disposable
- the algorithms used to produce, process and dispose of the records are also accountable, transparent, accurate and ethical

# On the consumer or user side:

- Computational processing is being used in the digital humanities and STEM fields to collate, cross-compile and analyze digitized and born-digital archival materials in new ways in order to gain new insights

- Archivists have to prepare digital archival content for end-user computational processing e.g.,
  - by ensuring archival content can be compiled, manipulated and curated at a very granular level
  - by adding descriptive metadata that supports new kinds of disciplinary research questions

# During archival appraisal and ingest:

- ○ Human appraisal of digitally-born records is becoming increasingly impossible due to massive volume, complexity and contingency, e.g., digital communications such as email and social media; networked and Cloud-based recordkeeping

- ○ Archivists will need to employ computational analysis to identify, ingest, and secure relationships between records and their components
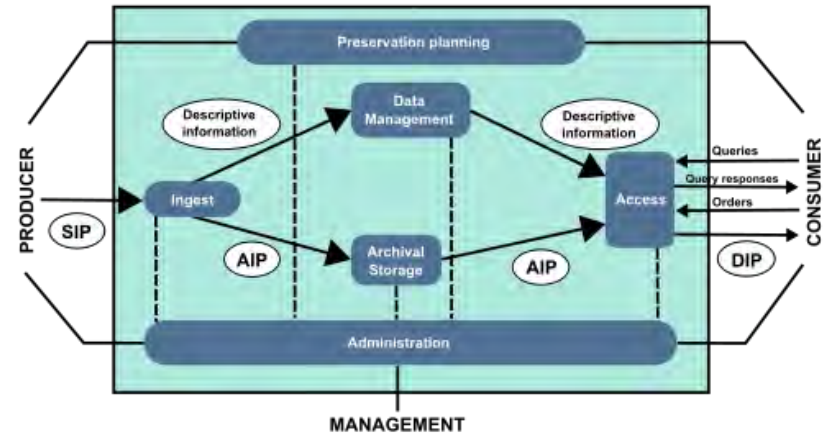
# During archival preservation:

- Manual preservation of increasing volumes of digital archival content does not scale

- Archivists will need to employ computational approaches in activities such as integrity checking, regular migration processes and allocating and tracking storage
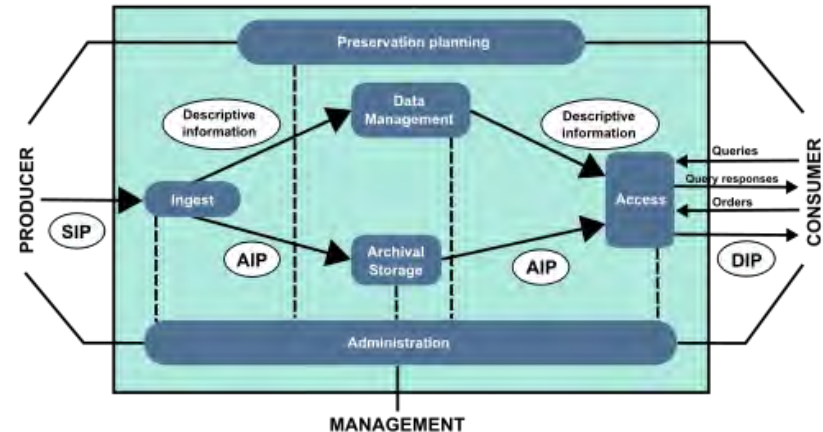
# During archival processing:

Archivists will use computational processes to:
- automate aspects of description, maintain administrative metadata, surface previously hidden aspects of collections, identify and explain anomalies, and curate collections to support specific researcher/consumer needs

# At the access interface:

- Archivists will use computational processes to:
    - establish consumer identity and privileges regarding collection access
    - search for, retrieve and package content that matches consumer queries and specifications
    - ensure privacy and security conditions governing dissemination and access of collections are met

# Foundational Paper on Computational Archival Science (CAS): Apr. 2018, Marciano et al.
## *"Archival records and training in the Age of Big Data"*
In *Re-Envisioning the MLS: Perspectives on the Future of Library and Information Science Education.*

**Eight Case Studies** w. examples of interdisciplinary efforts to address the changing context of recordkeeping and character of records:

1. **Evolutionary prototyping** and **computational linguistics**,
2. **Graph analytics, digital humanities** and **archival representation**,
3. **Computational finding aids,**
4. **Digital curation**,
5. **Public engagement** with (archival) content,
6. **Authenticity,**
7. Confluences between archival theory and computational methods: **cyberinfrastructure and the Records Continuum**,
8. **Spatial and temporal analytics**.

Each of the case studies concludes with a *"Takeaways for CAS/MLS Education"* statement.

| BROAD ARCHIVAL FUNCTION | |
|---|---|
| | Appraisal |
| | Accessioning |
| | Arrangement |
| | Description |
| | Preservation |
| | Access |
| | Records management |

| CAS TOPICS | |
|---|---|
| | Iterative design<br>Computational thinking<br>NLP |
| | Graph analytics |
| | Computational finding aids |
| | Digital curation |
| | Public engagement |
| | Authenticity |
| | Archival theory<br>Computational methods |
| | Spatial analytics<br>Temporal analytics |

# What is CAS?

## Working Definition of Computational Archival Science:

(R. Marciano et al. 2016, amended by N. Payne in 2018)

- A transdisciplinary field grounded in archival, information, and computational science that is concerned with the application of computational methods & resources, design patterns, sociotechnical constructs, and human-technology interaction, to large-scale *(big data)*:
    - records/archives processing, analysis, storage, long-term preservation, and access problems,
- with the aim of improving and optimizing efficiency, authenticity, truthfulness, provenance, productivity, computation, information structure & design, precision & human technology interaction in support of:
    - acquisition, appraisal, arrangement & description, preservation, communication, transmission, analysis &access decisions.

# CAS PORTAL: https://ai-collaboratory.net/cas/



#1 IEEE BIG DATA 2016
#2 IEEE BIG DATA 2017
#3 IEEE BIG DATA 2018
#4 IEEE BIG DATA 2019
#5 IEEE BIG DATA 2020
#6 IEEE BIG DATA 2021
#7 IEEE BIG DATA 2022
#8 IEEE BIG DATA 2023
#9 IEEE BIG DATA 2024

\* Workshops:

- 50+ workshops since 2016
- 9 CAS @ IEEE Big Data Conf.
  w. 150+ papers

Lessons learned from:
- **CAS#1**: 2016 in Washington, DC
- **CAS#2**: 2017 in Boston
- **CAS#3**: 2018 in Seattle
- **CAS#4**: 2019 in LA
- **CAS#5**: 2020 in Atlanta
- **CAS#6**: 2021 in Orlando
- **CAS#7**: 2022 in Osaka, Japan
- **CAS#8**: 2023 in Sorrento, Italy
- **CAS#9**: 2024 in Washington, DC

\* Presentations
\* Publications
\* Infrastructure

https://ai-collaboratory.net/cas/cas-workshops/2024-9th-cas-workshop/

*Welcome!*

**2024 IEEE International Conference on Big Data (IEEE BigData 2024)**

Dec 15-18, 2024 @ Washington DC, USA

IEEE Big Data 2024: CAS #9

- **Mon., Nov. 4, 2024 (final):** Due date for full workshop papers submission
- **Fri., Nov. 15, 2024:** Notification of paper acceptance to authors
- **Wed., Nov. 20, 2024 (hard deadline):** Camera-ready of accepted papers
- **Tue., Dec. 17, 2024:** Day-long CAS workshop (in person) in Washington DC, USA

**#9**
**IEEE BIG DATA**
**2024**

**RESEARCH TOPICS COVERED:**

- **Application of analytics to archival material**, including AI, ML, text- & data-mining, sentiment analysis, network analysis.
- **Analytics in support of archival processing**, including e-discovery, identification of personal information, appraisal, arrangement and description.
- **Scalable services for archives**, including identification, preservation, metadata generation, integrity checking, normalization, reconciliation, linked data, entity extraction, anonymization and reduction.
- **New forms of archives**, including Web, social media, audiovisual archives, and blockchain.
- **Cyber-infrastructures for archive-based research** and for development and hosting of collections
- **Big data and archival theory and practice**
- **Digital curation and preservation**
- **Crowd-sourcing** and archives
- **Big data and the construction of memory and identity**
- **Specific big data technologies** (e.g. NoSQL databases) and their applications
- **Corpora and reference collections** of big archival data
- **Linked data** and archives
- **Big data and provenance**
- **Constructing big data research objects** from archives
- **Legal and ethical issues** in big data archives

# Records Management Journal (Emerald Publishing)

**Special 2020 Issue on:**

Disruptive technologies for archives & records management and records professionals

**Editors:**
Julie McLeod, *Northumbria University, UK*
Richard Marciano, *University of Maryland, USA*

Summer 2020: Volume 30, Issue 2 & Issue 3
https://www.emerald.com/insight/content/doi/10.1108/RMJ-07-2020-057/full/html
https://www.emerald.com/insight/publication/issn/0956-5698/vol/30/iss/3

- Algorithm produced records
- Explainable Artificial Intelligence
- Natural Language Processing
- Automated Appraisal
- Internet-of-Things in the Archives
- Managing IoT-data for Gov. Agencies
- Collaboration between AI and Archival Science
- Preserving Virtual Reality
- Record Linking
- Mapping Archival Catalogs from Trees to Networks
- Blockchain and Records Management
- A Code of Ethics for the Digital Age

# ACM Journal on Computing and Cultural Heritage (JOCCH)

**Special 2022 Issue on:**

Computational Archival Science (CAS)

**Guest Editors:**
Mark Hedges, *King's College London, UK*
Eirini Goudarouli, *The National Archives, UK*
Richard Marciano, *University of Maryland, USA*

Vol. 13, Issue 1 (Feb. 2022) – Issue 3 (Sep. 2022)

https://ai-collaboratory.net/2020/05/21/jocch-cas_call_for_papers/
https://dl.acm.org/toc/jocch/2022/15/1
https://dl.acm.org/toc/jocch/2022/15/3

# Compendium of Core Computational Archival Science (CAS) Papers

| # | IEEE Big Data Conf.: Computational Archival Science (CAS) Workshops — Paper Title | Countries | Paper count | BROAD ARCHIVAL FUNCTION | CAS TOPICS | EXTENDED CAS TOPICS | SOCIAL JUSTICE TOPICS | RECORD TYPE |
|---|---|---|---|---|---|---|---|---|
| | **IEEE Big Data 2016 – Washington D.C., USA** | | 10 | | | | | |
| 16-1 | Exploring archives with probabilistic models: Topic Modelling for the valorisation of digitised archives of the European Commission | Belgium, Germany | | Description Access | NLP | | | Textual Records |
| 16-2 | Traces through Time: A Probabilistic Approach to Connected Archival Data | UK | | Description Access | NLP | | | Textual Records |
| 16-3 | Opening Up Dark Digital Archives Through The Use of Analytics to Identify Sensitive Content | USA | | Description Access | Graph analytics | AI/ML | | Textual Records |
| 16-4 | Computational Provenance in DataONE: Implications for Cultural Heritage Institutions | USA | | Description Preservation Access | Computational finding aids | | | Textual Records |
| 16-5 | Content-based Comparison for Collections Identification | USA | | Records management | Computational methods | | | Data Files |
| 16-6 | Breaking Down the Invisible Wall to Enrich Archival Science and Practice | USA | | Description | NLP Graph analytics | | | ALL |
| 16-7 | Mind the explanatory gap: Quality from Quantity | USA | | Records management | Computational thinking | | | ALL |
| 16-8 | Understanding Computational Web Archives Research Methods Using Research Objects | UK | | Description Access | Digital curation | | | Web Pages |
| 16-9 | Appraising Digital Archives with Archivematica | Canada | | Appraisal Preservation | Digital curation | | | ALL |
| 16-10 | Mining and Analysing One Billion Requests to Linguistic Services | Germany | | Description Access | NLP | | | ALL |
| | **IEEE Big Data 2017 – Boston, USA** | | 14 | | | | | |
| 17-1 | Building new knowledge from distributed scientific corpus: HERBADROP & EUROPEANA, Two concrete case studies for exploring big archival data | France, NL | | Description | Computational methods | | | Photographs & other Graphic Materials |
| 17-2 | An Infrastructure and Application of Computational Archival Science to Enrich and Integrate Big Digital Archival Data: Using Taiwan Indigenous Peoples Open Research Data (TIPD) as an Example | Taiwan | | Description | Spatial analytics | | Taiwan Indigenous People | Textual Records Maps & Charts |
| 17-3 | Computational Curation of a Digitized Record Series of WWII Japanese-American Internment | USA | | Description Access | NLP Graph analytics Spatial analytics | | WWII Japanese American Incarceration | Textual Records Maps & Charts |
| 17-4 | The Cybernetics Thought Collective Project: Using Computational Methods to Reveal Intellectual Context in Archival Material | USA | | Access | NLP | AI/ML | | Textual Records |
| 17-5 | Towards Automated Quality Curation of Video Collections from a Realistic Perspective | USA | | Appraisal Preservation | Computational methods | AI/ML | | Moving Images |
| 17-6 | Line Detection in Binary Document Scans: A Case Study with the International Tracing Service Archives | USA | | Description Access | Computational methods | Computer vision (CV) | Holocaust | Textual Records |
| 17-7 | Auto-Categorization Methods for Digital Archives | Canada, USA | | Description Records management | Computational methods | AI/ML | | Data Files |
| 17-8 | Heuristics for Assessing Computational Archival Science (CAS) Research: The Case of the Human Face of Big Data Project | USA | | Description Access | Iterative design | | Urban renewal | Textual Records Maps & Charts |
| 17-9 | What Can a Knowledge Complexity Approach Reveal About Big Data and Archival Practice? | NL | | Access | Computational thinking | | | ALL |

2  IEEE CAS Workshops (101) ▾   JOCCH (21) ▾   Records Management Journal (16) ▾   Misc (15) ▾

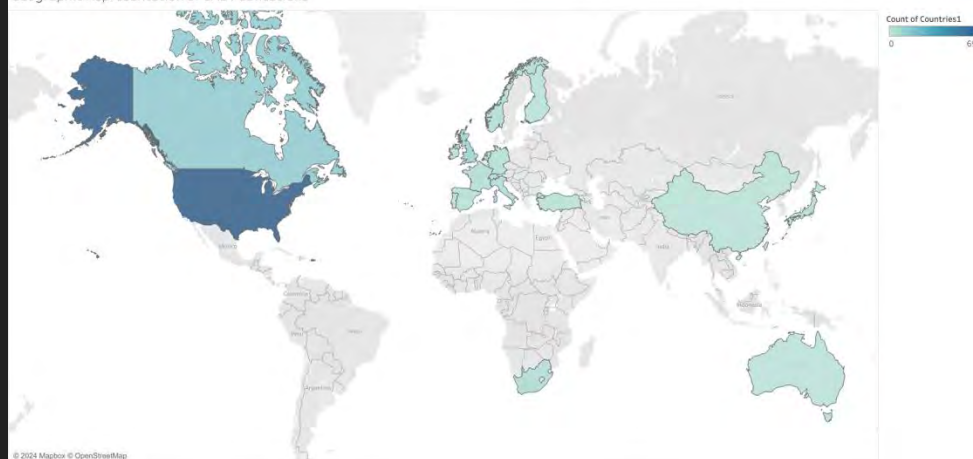| # | Paper Title | Countries | Paper count | BROAD ARCHIVAL FUNCTION | CAS TOPICS | EXTENDED CAS TOPICS | SOCIAL JUSTICE TOPICS | RECORD TYPE |
|---|---|---|---|---|---|---|---|---|
| | **IEEE Big Data Conf.: Computational Archival Science (CAS) Workshops** _tory.net/cas_ | | | | | **ANALYSIS** | | |
| | **IEEE Big Data 2023 – Sorrento, ITALY** | | 12 | | | | | |
| 23-1 | The Sequel: The Development of a Novel Context Capturing Method for the Functional Auto Classification of Records | Canada | | Description Records management | Computational methods | AI/ML | | Data Files |
| 23-2 | Specimen Outlining: A Computational Archival Science Approach | USA | | Description | Computational methods | AI/ML | | Photographs & other Graphic Materials |
| 23-3 | Who's in My Archive? An End-to-End Framework for Automatic Annotation of TV Personalities | Italy | | Description | Computational finding aids | AI/ML | | Moving Images |
| 23-4 | Authenticating Citizen Journalism by Incorporating the View of Archival Diplomatics into the Verification of Open-source Investigators | Canada | | Description Preservation Access | Computational methods | | | Moving Images |
| 23-5 | Will Blockchain Technology Change How Well National Archives Preserve the Trustworthiness of Digital Records?: Preliminary Results of a Survey | Turkey, Canada | | Description Preservation Access | Authenticity | | | ALL |
| 23-6 | Analogous Analogues: Digital Twins and Hardware Tracking in GLAM Collections | Canada | | Preservation Records management | Authenticity | | | Photographs & other Graphic Materials |
| 23-7 | Critical Community-Centeredness: Ethical Considerations for Computational Archival Studies | USA | | Creation | Public engagement | | | ALL |
| 23-8 | Accelerating Precision Research and Resolution Through Computational Archival Science Pedagogy | USA | | Arrangement Description Access | Computational thinking Spatial analytics | | Holocaust | Textual Records |
| 23-9 | The Utility of Standards and Good Practice Guidelines for Records Professionals: Comparing Apples, Oranges, and Other Fruits | South Africa | | Records management | Computational methods | | | ALL |
| 23-10 | Can GPT-4 Think Computationally about Digital Archival Practices? | USA | | Description | Computational thinking | GenAI LLM | WWII Japanese American Incarceration | Textual Records Maps & Charts |
| 23-11 | Exploring the Application of Large Language Models in Detecting and Protecting Personally Identifiable Information in Archival Data: A Comprehensive Study | China | | Description Access | NLP | GenAI LLM | | Textual Records |
| 23-12 | AI-Generated Images as an Emergent Record Format | USA | | Appraisal | Digital curation | AI/ML GenAI Computer vision (CV) | | Photographs & other Graphic Materials |
| | | | 101 | | | | | |

Jennifer Proctor

DublinCore Webinar: Oct. 19, 2023
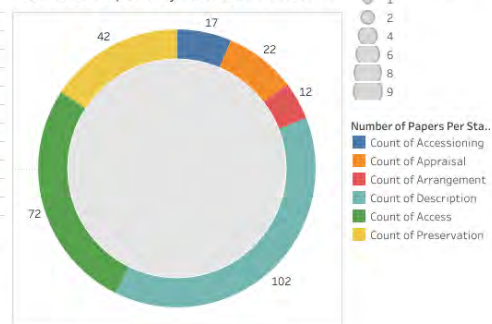
AI & NLP for Open-Source Archival Linked Data Workflows
https://www.dublincore.org/webinars/2023/ai-nlp-archival-linked-data/

GraphRAG with Neo4j

**Neo4j LLM Knowledge Graph Builder**
Allison Cossette
https://www.youtube.com/watch?v=UYJbG3p68NM

# GraphRAG (Retrieval Augmented Generation with Graphs):

Merges Knowledge Graphs "with context, structure & trust".

Instead of relying only on text chunk searches (using Vector DBs), it uses graph queries to pull relevant, connected data. The claim is that it helps with "Explainable AI": every result comes with an audit trail, helping with trust and transparency.

In the Year: "2024", Paper: "M-24-9" was Published in the Journal: "ArchivesHandbook". It connects to the following Steps: "Description, Access", relates to the following Tools: "Computational thinking, Digital curation, Graph analytics, Spatial analytics, Temporal analytics", speaks to the following Topics: "City directories", and links to the following Items: "Textual Records, Maps & Charts".

In the Year: "2024", Paper: "M-24-10" was Published in the Journal: "iConference". It connects to the following Steps: "Appraisal, Accessioning, Arrangement, Description, Preservation, Access", relates to the following Tools: "Iterative design, Computational thinking, NLP, Graph analytics, Computational finding aids, Digital curation, Public engagement, Authenticity, Archival theory, Computational methods, Spatial analytics, Temporal analytics", adds to the following Technologies: "GenAI, LLM", and links to the following Items: "Architectural & Engineering Drawings, Data Files, Maps & Charts, Moving Images, Photographs & other Graphic Materials, Sound Recordings, Textual Records, Web Pages".

In the Year: "2018", Paper: "M-18-11" was Published in the Journal: "DigitalHeritage". It connects to the following Steps: "Description, Access", relates to the following Tools: "Iterative design, Digital curation, Computational methods", speaks to the following Topics: "WWII Japanese American Incarceration", and links to the following Items: "Textual Records, Maps & Charts".

In the Year: "2024", Paper: "M-24-12" was Published in the Journal: "AEOLIAN". It connects to the following Steps: "Accessioning, Arrangement, Description, Access", relates to the following Tools: "Digital curation", adds to the following Technologies: "GenAI, LLM", speaks to the following Topics: "Legacy of Slavery", and links to the following Items: "Textual Records".
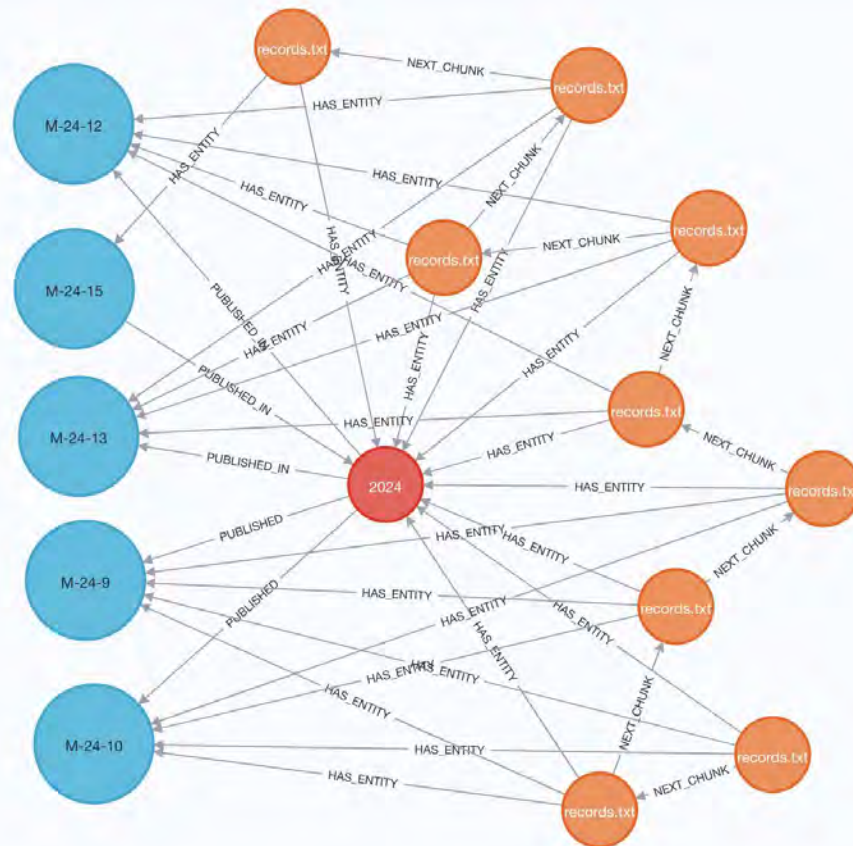
In the Year: "2024", Paper: "M-24-13" was Published in the Journal: "ISGC". It connects to the following Steps: "Accessioning, Arrangement, Description, Access", relates to the following Tools: "Iterative design, Digital curation, Spatial analytics, Temporal analytics", adds to the following Technologies: "GenAI, LLM", speaks to the following Topics: "Legacy of Slavery", and links to the following Items: "Textual Records".

In the Year: "2021", Paper: "M-21-14" was Published in the Journal: "AI&Society". It connects to the following Steps: "Description, Access", relates to the following Tools: "Iterative design, Temporal analytics", adds to the following Technologies: "AI/ML", speaks to the following Topics: "Holocaust", and links to the following Items: "Textual Records".

In the Year: "2024", Paper: "M-24-15" was Published in the Journal: "MTSR". It connects to the following Steps: "Description, Access", relates to the following Tools: "Spatial analytics, Digital curation", adds to the following Technologies: "AI/ML, Computer vision (CV)", speaks to the following Topics: "WWII Japanese American Incarceration", and links to the following Items: "Textual Records, Maps & Charts".

**Overview**

**Node labels**

* (377)  Chunk (153)  Year (9)  Paper (153)
Journal (14)  Step (7)  Tool (12)  Item (8)
Technology (4)  Topic (14)  _Bloom_Perspective_ (1)
_Bloom_Scene_ (1)  Document (1)

**Relationship types**

* (5376)  PART_OF (153)  NEXT_CHUNK (152)
HAS_ENTITY (3484)  SIMILAR (180)
PUBLISHED_IN_YEAR (12)  PUBLISHED_IN_JOURNAL (12)
CONNECTS_TO_STEP (23)  RELATES_TO_TOOL (15)
LINKS_TO_ITEM (26)  ADDS_TO_TECHNOLOGY (2)
RELATES_TO (243)  PUBLISHED_IN (274)
CONNECTS_TO (274)  LINKS_TO (416)  SPEAKS_TO (42)
ADDS_TO (58)  PUBLISHED (8)
_Bloom_HAS_SCENE_ (1)  FIRST_CHUNK (1)

Displaying 377 nodes, 5,376 relationships.

`neo4j$ MATCH (y:Paper {id: 'M-24-10'})-[l]-(n) RETURN y,l,n`

Node properties

Chunk

| | |
|---|---|
| <elementid> | 4:1290cf3d-f157-4a5b-a025-eba8bbd0b344:84 |
| <id> | 84 |
| content_offset | 45533 |
| embedding | [0.004388340283185244,0.04025658592581749,-0.0 23980682715773582,-0.05043290555477142,0.0300 50864443182945,-0.037488505244255066,-0.09333 04950594902,0.0... Show all] |
| fileName | records.txt |
| id | 0faf2b003836997ece580e18535a48eefff4bdd3 |
| length | 686 |
| position | 148 |
| text | In the Year: 2024. Paper: M-24-10 was Published in the Journal: iConference. It connects to the following Steps: Appraisal, Accessioning, Arrangement, Description, Preservation, Access, relates to the following Tools: Iterative design, Computational thinking, NLP, Graph analytics, Computational finding aids, Digital curation, Public engagement, Authenticity, Archival theory, Computational methods, Spatial analytics, Temporal analytics, adds to the following Technologies: GenAI, LLM. and links to the following items: Architectural & Engineering Drawings, Data Files, Maps & Charts, Moving Images, Photographs & other Graphic Materials, Sound Recordings, Textual Records, Web Pages. |

**GraphRAG** using ChatGPT4.o prompts:

The paper IDs related to the year 2024 are M-24-12, M-24-10, M-24-13, M-24-15, and M-24-9.
10/14/2024 6:03:50 PM

Details

## Retrieval information

To generate this response, the process took **5.63 seconds**, utilizing **11579** tokens with the model **gpt-4o-2024-08-06** in **graph & vector & fulltext** mode.

Sources used    **Top Entities used**    Chunks

**Paper (153)** C-16-1, C-16-2, C-16-3

**Journal (14)** Cas, Jocch, Rmj

**Topic (14)** Taiwan Indigenous People, Wwii Japanese American Incarceration, Hol..

**Tool (12)** Nlp, Graph Analytics, Computational Finding Aids

**Year (9)** 2016, 2022, 2017

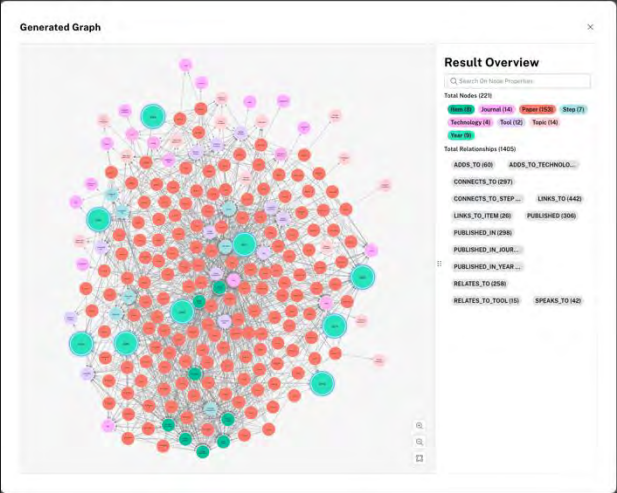**Item (8)** Textual Records, Data Files, Architectural & Engineering Drawings

**Step (7)** Description, Access, Preservation

**Technology (4)** Ai/Ml, Computer Vision (Cv), Genai

Graph Entities used for Answer Generation

### Generated Graph



**Result Overview**

Total Nodes (221)
Item (8)  Journal (14)  Paper (153)  Step (7)
Technology (4)  Tool (12)  Topic (14)
Year (9)

Total Relationships (1405)
ADDS_TO (61)    ADDS_TO_TECHNOLO...
CONNECTS_TO (297)
CONNECTS_TO_STEP_...    LINKS_TO (442)
LINKS_TO_ITEM (26)    PUBLISHED (306)
PUBLISHED_IN (296)
PUBLISHED_IN_JOUR...
PUBLISHED_IN_YEAR_...
RELATES_TO (258)
RELATES_TO_TOOL (15)    SPEAKS_TO (42)

## Retrieval information

To generate this response, the process took **5.63 seconds**, utilizing **11579** tokens with the model **gpt-4o-2024-08-06** in **graph & vector & fulltext** mode.

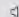Sources used    Top Entities used    **Chunks**

- **Similarity Score: 0.786**
  In the Year: 2020, Paper: M-20-3 was Published in the Journal: SAA. It connects to the following Steps: Description, Access, relates to the following Tools: Iterative design, Computational thinking, speaks to the following Topics: Legacy of Slavery, and links to the following Items: Textual Records.

- **Similarity Score: 0.7583**
  In the Year: 2018, Paper: M-18-11 was Published in the Journal: DigitalHeritage. It connects to the following Steps: Description, Access, relates to the following Tools: Iterative design, Digital curation, Computational methods, speaks to the following Topics: WWII Japanese American Incarceration, and links to the following Items: Textual Records, Maps & Charts.

- **Similarity Score: 0.7571**

The paper IDs mentioned (M-24-10, M-24-12, M-24-13, M-24-15, and M-24-9) link to the following Items:
**M-24-10**: Architectural & Engineering Drawings, Data Files, Maps & Charts, Moving Images, Photographs & Other Graphic Materials, Sound Recordings, Textual Records, Web Pages.
**M-24-12**: Textual Records.
**M-24-13**: Textual Records.
**M-24-15**: Maps & Charts, Textual Records.
**M-24-9**: Maps & Charts, Textual Records.
10/14/2024 6:05:04 PM

Details

From these paper IDs which Steps are there?
10/14/2024 6:05:35 PM
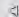
The paper IDs you mentioned are associated with the following Steps:
**M-24-10**: Access, Accessioning, Appraisal, Arrangement, Description, Preservation
**M-24-12**: Access, Accessioning, Arrangement, Description
**M-24-13**: Access, Accessioning, Arrangement, Description
**M-24-15**: Access, Description
**M-24-9**: Access, Description
10/14/2024 6:05:46 PM

Details

From these paper IDs which Tools are there?
10/14/2024 6:06:20 PM

The paper IDs you mentioned are associated with the following Tools:
**M-24-10**: Computational Methods, Computational Thinking, Digital Curation, Graph Analytics
**M-24-12**: Digital Curation
**M-24-13**: Digital Curation
**M-24-15**: Digital Curation
**M-24-9**: Computational Thinking, Digital Curation, Graph Analytics
10/14/2024 6:06:31 PM

# Richard Marciano

## College of INFORMATION

computational-archival-science@googlegroups.com

https://ai-collaboratory.net/

marciano@umd.edu

https://www.linkedin.com/in/richard-marciano/