# Lessons Learned in Building Linked Data for the American Art Collaborative

Craig Knoblock

University of Southern California

Information Sciences Institute

Pedro Szekely, Eleanor Fink, David Newbury, Robert Sanderson,
Duane Degler, Kate Blanch, Sara Snyder, Nilay Chheda, Nimesh Jain,
Ravi Raju Krishna, Nikhila Begur Sreekanth, and Yixiang Yao

# Project Goals

- Launch the American Art Collaborative
  - Consortium of 14 American art museums
  - Explore the use of Linked Data to make their data available for research, education, and outreach
- Build 5* Linked Data for the museums
  - Map the data about artwork and artists to a common ontology
  - Link the data to other resources
  - Create/extend tools to support the construction of Linked Data
  - Create applications using the data

# Outline

- Mapping the data
- Linking the entities
- Using the Linked Data
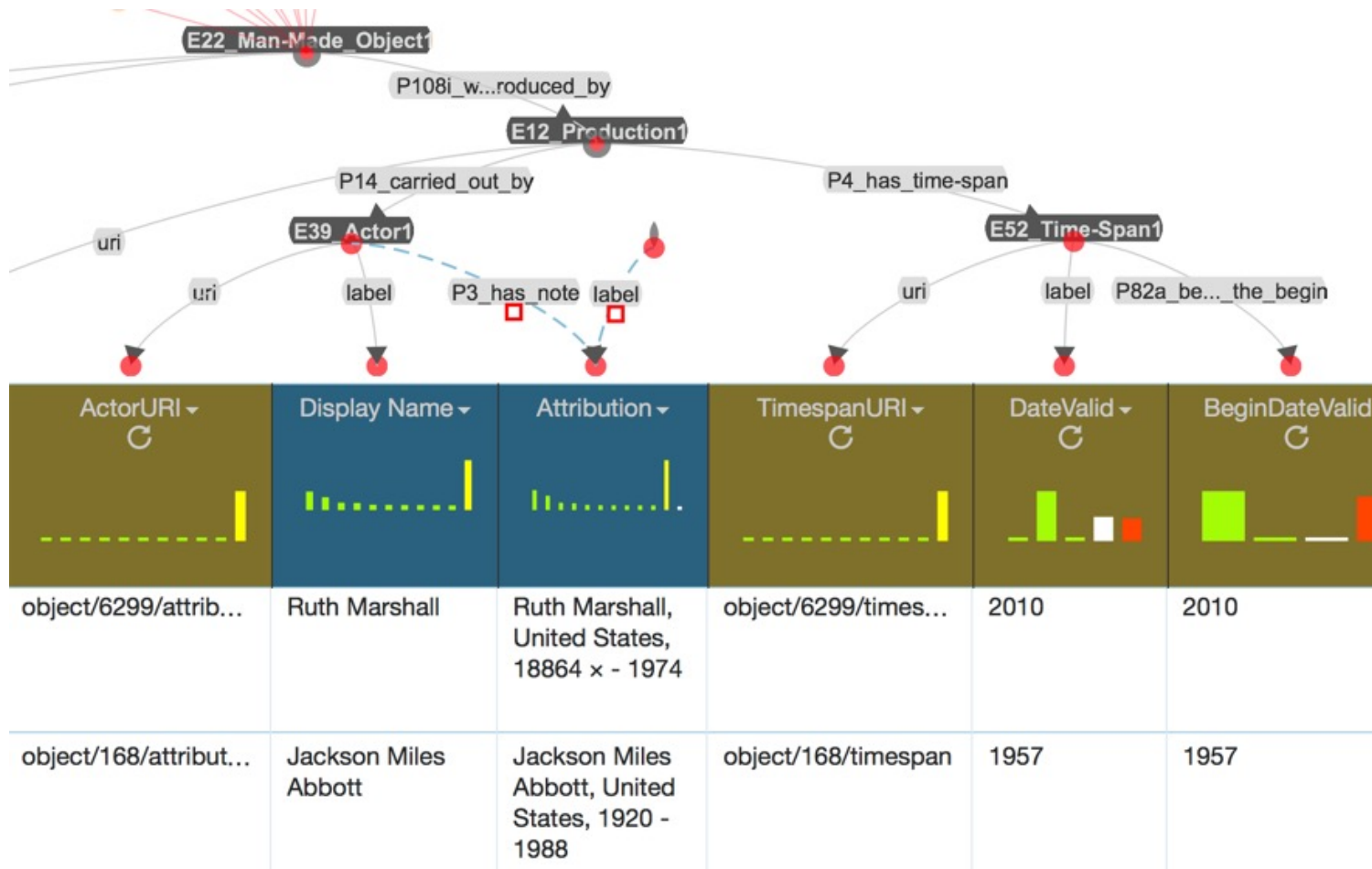- Related Work & Discussion

# Mapping the Data

- Challenges
  - Museums have the data in wildly different formats and use different schemas
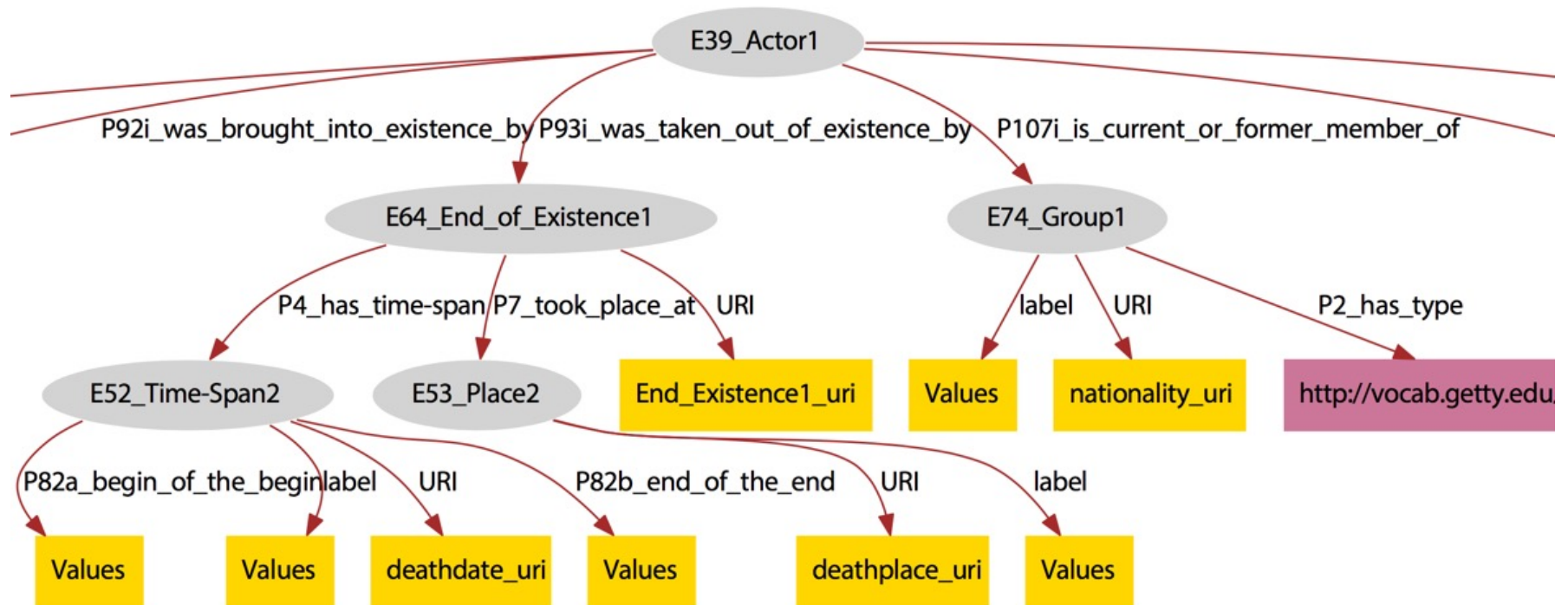  - The CIDOC-CRM ontology is a large and complicated ontology

- Approach
  - Use Github to organize all of the data, mappings, and resulting RDF
  - Use Karma to create the mappings of each dataset
  - Trained a team of USC students to apply the tools to the datasets
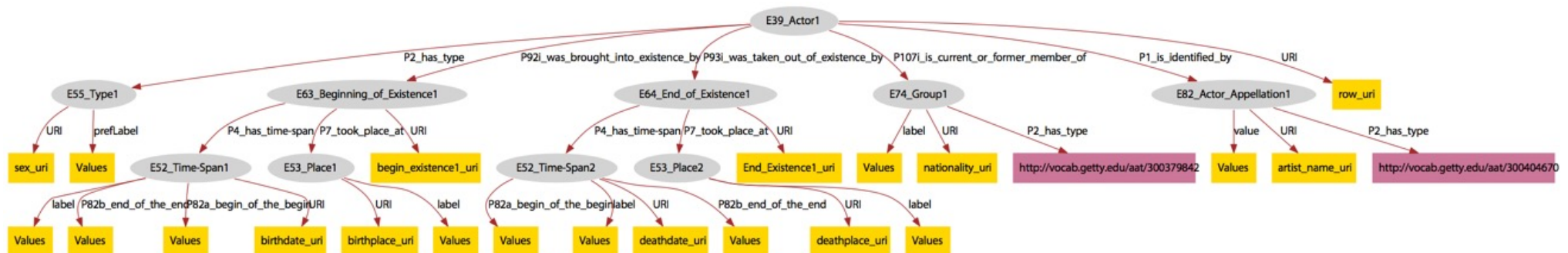
# Use Karma to Map the Data to the Ontology

# Example Model of Actor for Amon Carter

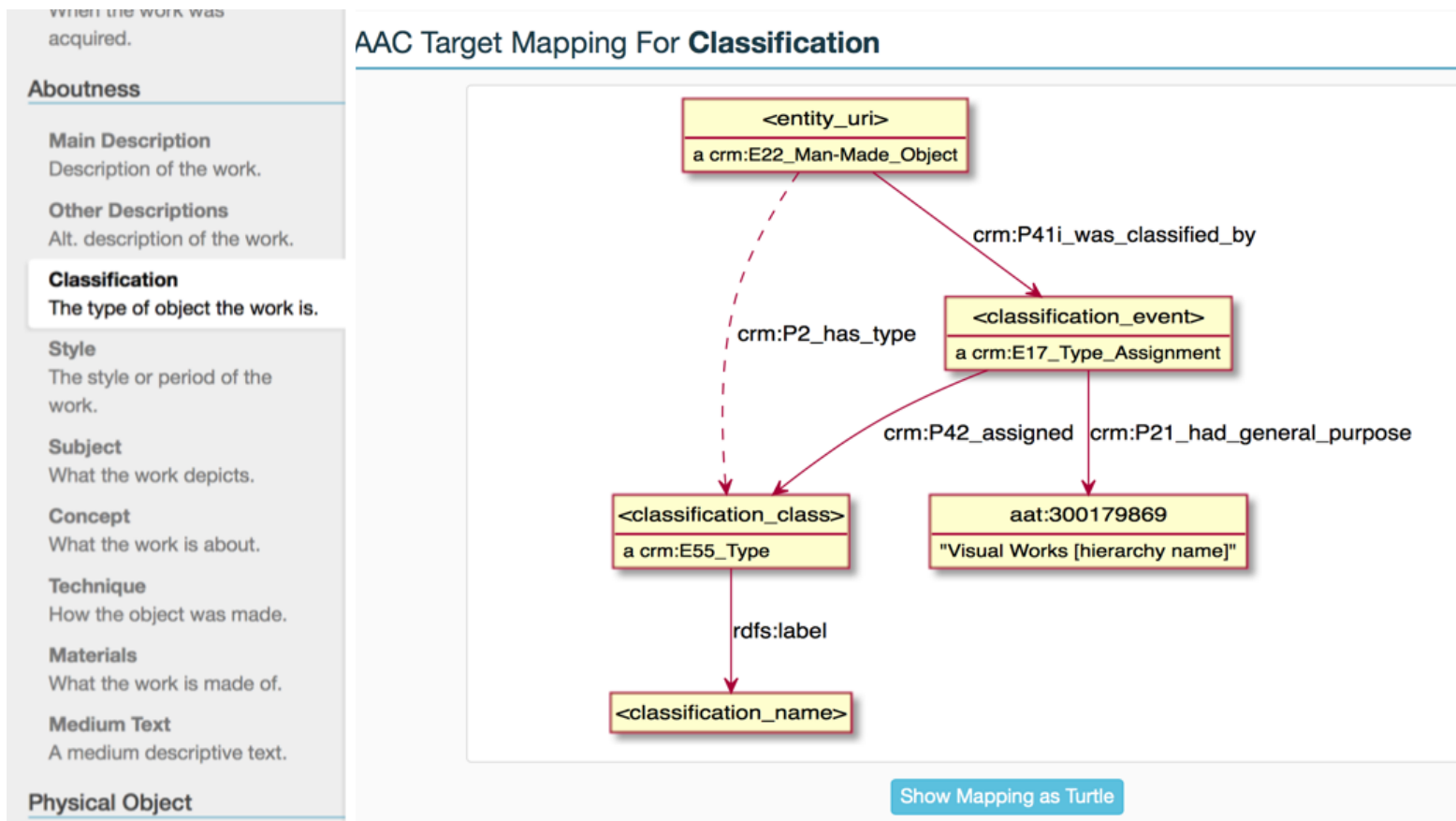# Complete Model of Actor for Amon Carter

# AAC Data Statistics

| Museum | Format | Files | Mappings | People | Commits | Issues |
|---|---|---|---|---|---|---|
| Archives of American Art | xls | 5 | 5 | 5 | 67 | 17 |
| Amon Carter Museum | xml | 2 | 3 | 7 | 195 | 17 |
| Autry Museum | xlsx | 6 | 6 | 9 | 309 | 68 |
| Crystal Bridges Museum | csv | 8 | 14 | 7 | 572 | 76 |
| Colby College Museum of Art | json | 1 | 2 | 7 | 345 | 31 |
| Dallas Museum of Art | csv | 2 | 2 | 3 | 250 | 11 |
| Gilcrease Museum | xlsx | 9 | 12 | 5 | 447 | 24 |
| Indianapolis Museum of Art | json | 3 | 3 | 6 | 214 | 16 |
| National Museum of Wildlife Art | csv | 2 | 3 | 6 | 196 | 9 |
| National Portrait Gallery | xlsx | 11 | 12 | 7 | 334 | 75 |
| Princeton University Art Museum | json | 10 | 11 | 7 | 421 | 53 |
| Smithsonian American Art Museum | csv | 11 | 14 | 4 | 408 | 49 |
| Walters Art Museum | xml | 6 | 12 | 6 | 878 | 28 |
| Total | 4 | 76 | 99 | | 4,636 | 474 |

# AAC Target Mappings

# AAC Mapping Validator

# Statistics on the Mappings

| Museum | Data Trans. | Structure Trans. | Semantic Classes | Semantic Types | Semantic Links |
|---|---|---|---|---|---|
| Archives of American Art | 46 | 0 | 30 | 65 | 43 |
| Amon Carter Museum | 13 | 3 | 13 | 26 | 14 |
| Autry Museum | 76 | 0 | 46 | 87 | 49 |
| Crystal Bridges Museum | 112 | 6 | 74 | 132 | 89 |
| Colby College Museum of Art | 52 | 0 | 36 | 69 | 52 |
| Dallas Museum of Art | 46 | 0 | 27 | 55 | 39 |
| Gilcrease Museum | 105 | 5 | 75 | 132 | 109 |
| Indianapolis Museum of Art | 87 | 2 | 55 | 101 | 75 |
| National Museum of Wildlife Art | 37 | 0 | 24 | 47 | 34 |
| National Portrait Gallery | 112 | 2 | 64 | 118 | 69 |
| Princeton University Art Museum | 116 | 5 | 95 | 153 | 115 |
| Smithsonian American Art Museum | 88 | 4 | 67 | 114 | 95 |
| Walters Art Museum | 78 | 8 | 56 | 99 | 71 |
| Total | 968 | 35 | 662 | 1,198 | 854 |

# Statistics on What Was Mapped

| Museum | Constituents | Objects | Events | Places | Triples |
|---|---|---|---|---|---|
| Archives of American Art | 6,944 | 15,025 | 7,301 | 1,592 | 210,360 |
| Amon Carter Museum | 806 | 6,421 | 13,164 | 532 | 225,528 |
| Autry Museum | 148 | 193 | 558 | 0 | 14,639 |
| Crystal Bridges Museum | 514 | 1,691 | 3,384 | 0 | 96,533 |
| Colby College Museum of Art | 2,210 | 8,217 | 18,905 | 0 | 456,711 |
| Dallas Museum of Art | 1,299 | 2,229 | 5,639 | 0 | 114,184 |
| Gilcrease Museum | 1,578 | 20,904 | 83,603 | 4,159 | 1,851,246 |
| Indianapolis Museum of Art | 2,131 | 22,314 | 34,560 | 432 | 846,952 |
| National Museum of Wildlife Art | 376 | 2,208 | 2,226 | 0 | 83,486 |
| National Portrait Gallery | 12,553 | 16,829 | 54,097 | 5,713 | 1,902,699 |
| Princeton University Art Museum | 2,899 | 13,314 | 43,828 | 881 | 1,253,239 |
| Smithsonian American Art Museum | 20,490 | 43,038 | 106,534 | 3,042 | 2,597,938 |
| Walters Art Museum | 182 | 801 | 1722 | 159 | 60,136 |
| Total | 52,130 | 153,184 | 375,521 | 16,510 | 9,713,651 |

# Mapping Lessons

- Lesson 1: Reproducible Workflows
  - Allow museums to export raw data from their collection management systems

- Lesson 2: Shared Repository
  - Github was invaluable for managing all the data and mapping files

- Lesson 3: Data Cleaning
  - Significant data cleaning was required
  - Integrated as part of the data processing workflow in Karma

- Lesson 4: Mapping Inconsistencies
  - Validation tool was critical to completing a consistent set of mappings

- Lesson 5: Expert Review
  - The outside review was crucial in identifying and resolving mapping inconsistencies

# Linking the Data

- Goals
  - Link the entities in the museum data to other resources
  - Capability for museums to curate the automatically generated links
  - Demonstration: linking artists to Getty ULAN

- Approach
  - Attempted to use existing linking tools, but they either didn't scale or students found them difficult to configure
  - Wrote a specialized script to generate high recall & precision candidates
  - Built a link review tool that the museums used to curate their links
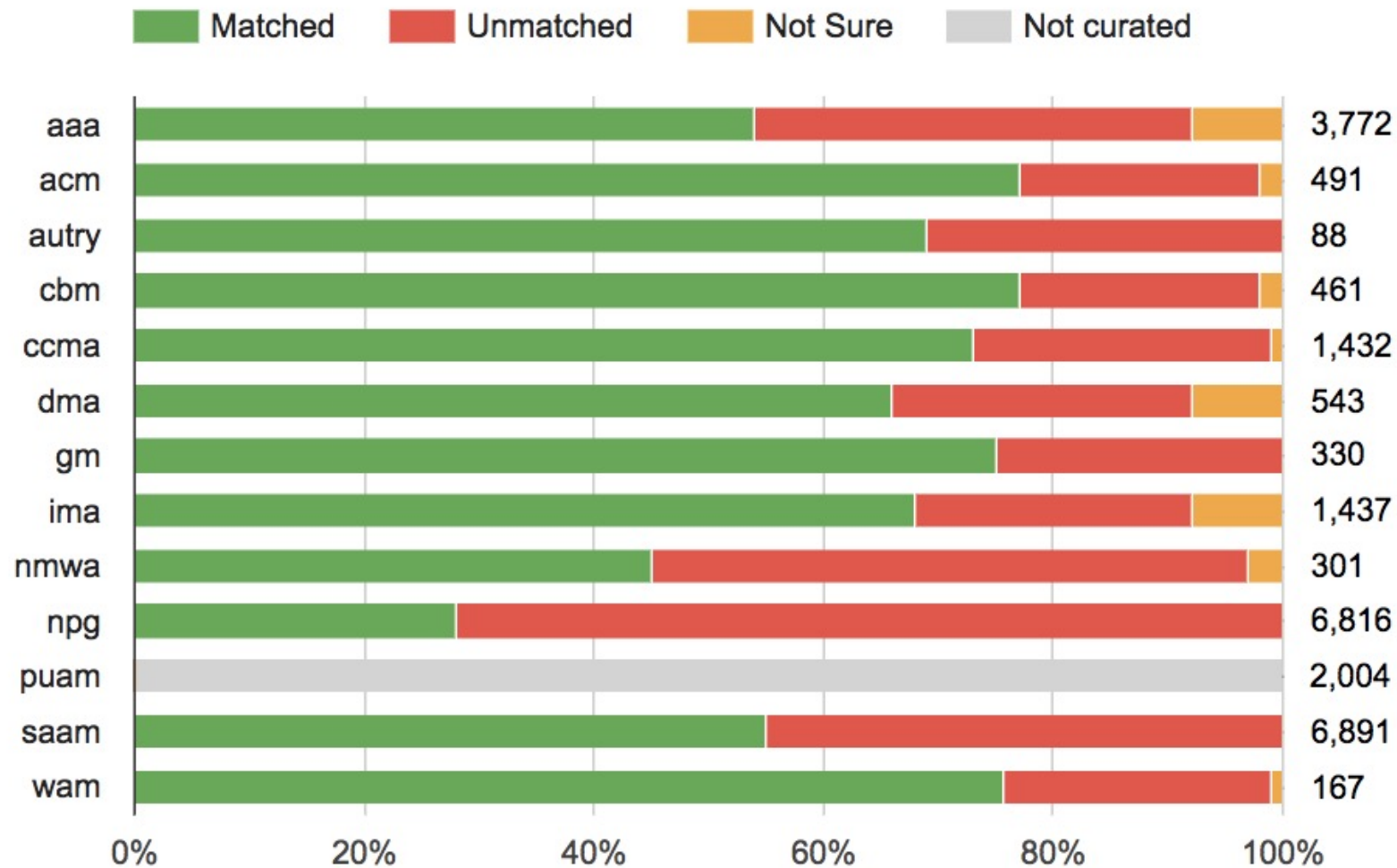
# Link Review Screenshot

| acm | ulan |
|---|---|
| Laton A. Huffman | Huffman, L. A. |
| **Similarity Score: 0.920** ||

**Matching Values**

| gender | male |
|---|---|

**Different Values**

| object_links | http://www.cartermuseum.org/imu/acm/#details=ecatalogue.28344<br>http://www.cartermuseum.org/imu/acm/#details=ecatalogue.92682<br>http://www.cartermuseum.org/imu/acm/#details=ecatalogue.31996<br>http://www.cartermuseum.org/imu/acm/#details=ecatalogue.51417<br>http://www.cartermuseum.org/imu/acm/#details=ecatalogue.187882 | None |
|---|---|---|
| death_year | 1931-12-28 | 1931 |
| uri | http://data.americanartcollaborative.org/acm/artist/6026 | http://vocab.getty.edu/ulan/500016161 |
| nationality | American | American (North American) |
| birth_year | 1854-10-31 | 1854 |

YES ➤          NO ➤          NOT SURE ➤

# Screenshot with Linking Review Results



Legend: Matched (green), Unmatched (red), Not Sure (orange), Not curated (gray)

| Institution | Count |
|---|---|
| aaa | 3,772 |
| acm | 491 |
| autry | 88 |
| cbm | 461 |
| ccma | 1,432 |
| dma | 543 |
| gm | 330 |
| ima | 1,437 |
| nmwa | 301 |
| npg | 6,816 |
| puam | 2,004 |
| saam | 6,891 |
| wam | 167 |

# Summary Statistics on the Linking Process

- Number of constituents in data: 42,685
- Previously existing links to ULAN: 3,349
- Linking based on previously existing links:
  - Precision: .96
  - Recall: .88
  - F1-measure: .92
- Candidate matches: 24,733
- Incorrect links in museums datasets: 19
- New links to ULAN: 9,357
- Incorrect links after review: 2
- Previously existing links not discovered: 136

# Lessons Learned on Linking

- Lesson 6: Linking Tools
  - Difficult to configure and use the existing linking tools and get them to scale to large datasets (e.g., ULAN, DBPedia, & VIAF)
  - We need easy to work with and scalable libraries for linking tasks

- Lesson 7: Manual Review
  - Users are willing to invest significant time and energy to ensure the final data is accurate
  - The museums reviewed almost 25K links!
  - A few weeks of effort almost tripled the number of links to ULAN

# Using the Data: The Browse Application



American Art Collaborative Demonstration Application | About the AAC | Settings

## AAC COLLECTIONS

INSTITUTIONS    EXPLORE ARTISTS    EXPLORE BY CATEGORIES    COLLECTION PROFILE

Gilcrease Museum

## CRUCITA - TAOS INDIAN GIRL IN OLD HOPI WEDDING DRESS AND DRY FLOWERS (WINTER BOUQUET)

### Joseph Henry Sharp

| | |
|---|---|
| ALTERNATE TITLE | Crucita - Taos Indian Girl |
| CREATION DATE | circa 1926 |
| OBJECT # | 0137.2194 |
| PARTNER URL | https://collections.gilcrease.org/object/01372194 |
| TYPES | Oil Painting<br>Painting & Drawing<br>Paintings |
| MATERIAL | Oil on canvas |
| DIMENSION | Overall: 47 1/2 × 55 1/2 × 3 in. (120.7 × 141 × 7.6 cm)<br>Framed: 47 1/2 × 55 3/8 × 3 1/4 in. (120.7 × 140.7 × 8.3 cm) |
| DIMENSIONS | **Framed**<br>Width: 140.65<br>Depth: 8.26<br>Height: 120.65<br>**Overall**<br>Depth: 7.62<br>Width: 140.97<br>Height: 120.65 |
| SUBJECTS | Hopi Indians<br>Hopi |

**RELATED WORKS**

72 Other works by this artist in this institution

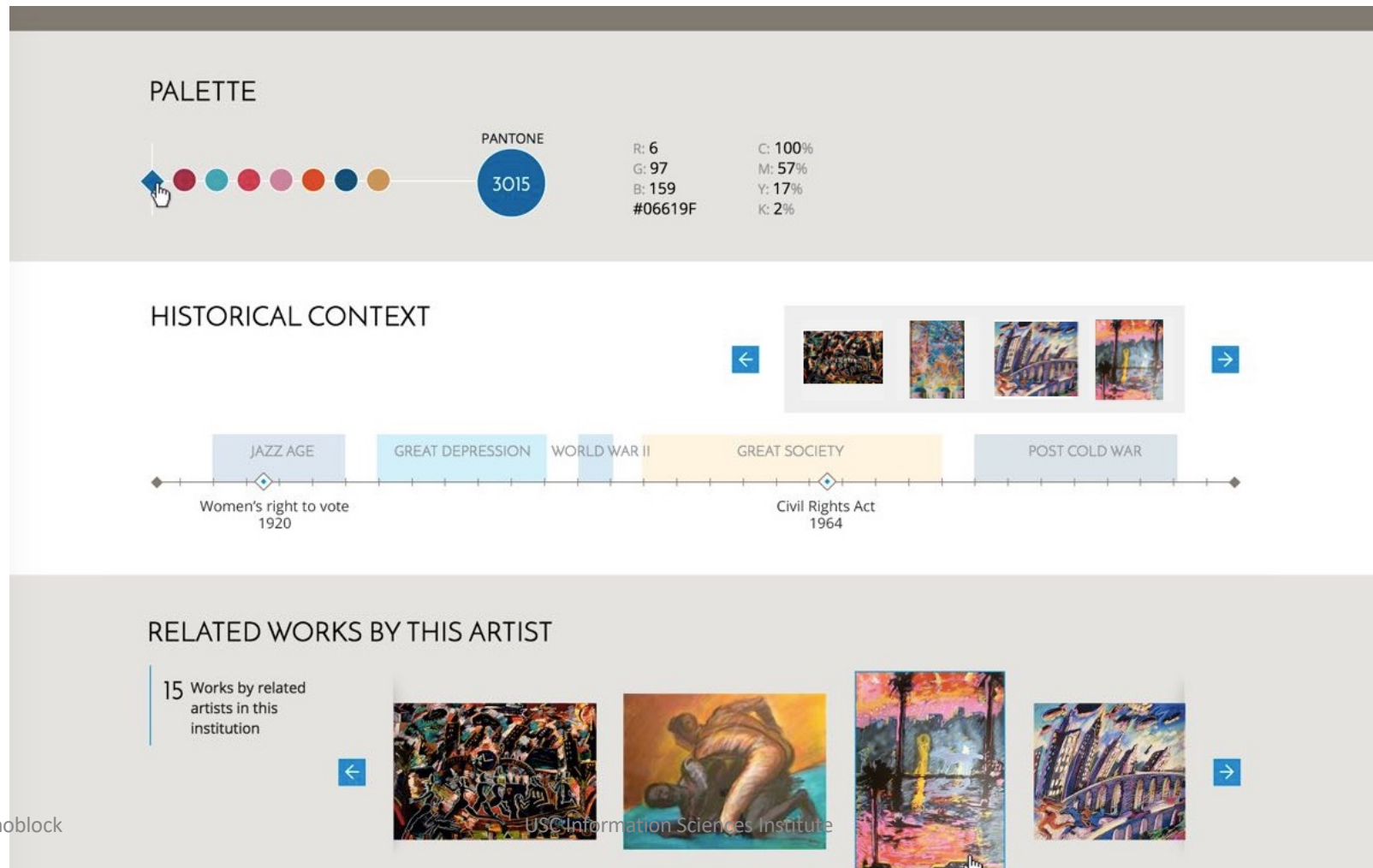34 Works by this artist in other institutions

___ Works on a similar subject in this institution

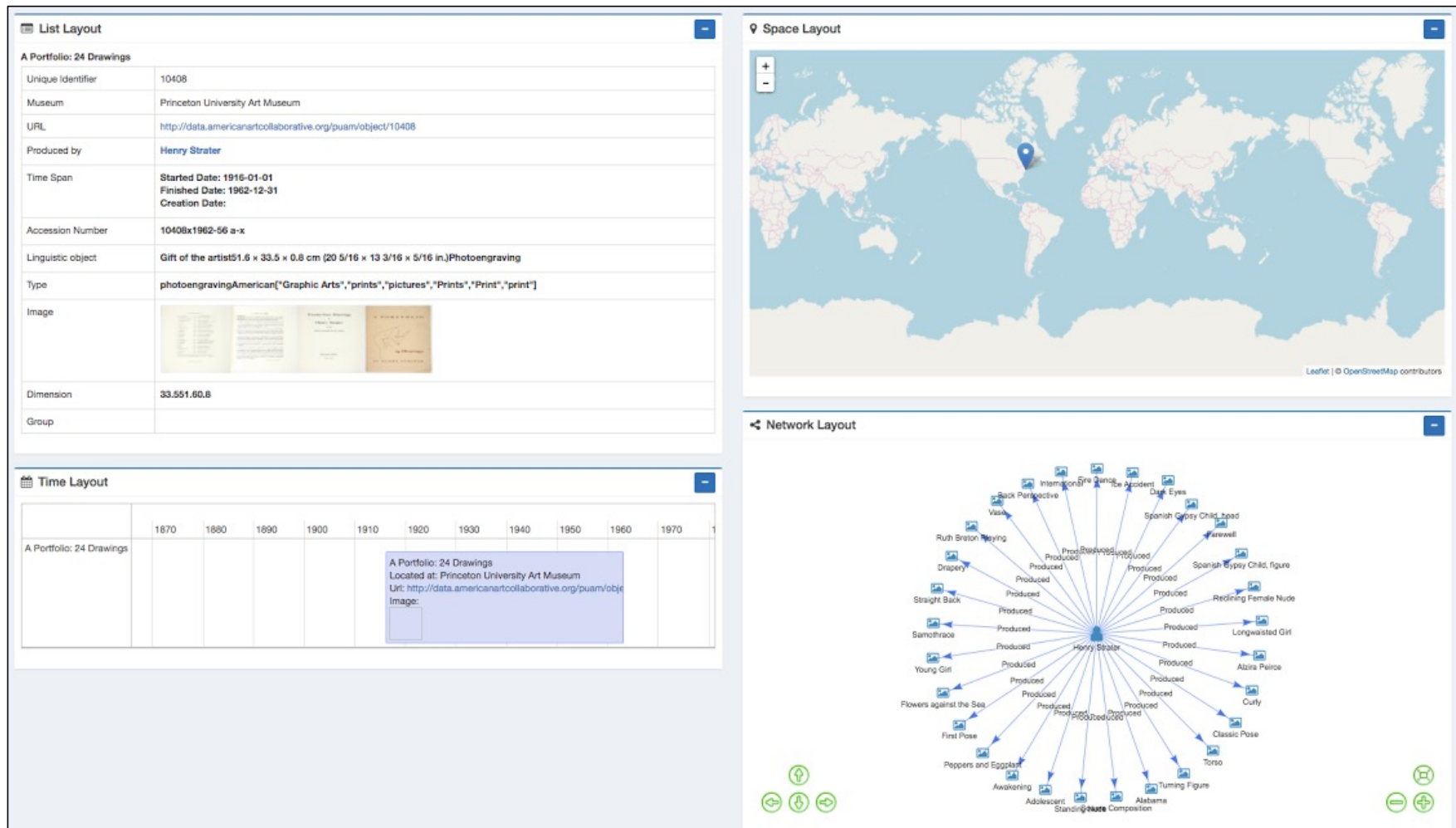___ Works on a similar subject in other institutions

___ Works by related artists in this institution

___ Works by related artists in other

19

# Using the Data: Under Development

# SemUI Visualization [Giunchiglia, Ojha, & Das, ICSC 2017]

# SemSpect (http://aac.semspect.de/)
## Thorsten Liebig (derivo)

# Lessons Learned in Using the Data

- Lesson 8: Visualization
  - Easy to understand visualization is needed for non-technical users

- Lesson 9: Simple Schema
  - CRM ontology may be useful for research, but challenging for applications
  - Created a set of SPARQL queries to create JSON objects
    - Loaded the objects into Elasticsearch for complex analysis

# Related Work

- Consortiums of museums
  - Europeana – 1500 cultural heritage institutions, 17 million items
  - CHIN – 8 Canadian museums, 85,000 items
  - LODAC – 114 museums in Japan
  - Published a fixed schema and mapped all museums/institutions to the schema

- Organizations using the CRM ontology
  - Research space (British museum & Yale Center for British Art)
  - Pharos – 14 historical photo archives
  - Each organization is responsible for publishing their own data to CRM

- Mapping data to CRM
  - X3ML – maps XML data to CRM using manually written rules

- Linking
  - Silk, Dedup, etc. – focus is on automatic linking, but no curation of the links
  - Mix'n'match, OpenRefine – support link review, but targets highly technical users

# Discussion

- Collaborated with 14 American art museums to publish 5* Linked Data
- Created a set of tools to create the linked data
  - Karma – clean and map the data, publish directly to Gitub
  - Mapping Validation tool – review the mappings to ensure consistency
  - Karma execution tool – applies Karma mappings and published both RDF and JSON-LD directly to Github
  - Link Review tool – allows non-technical users to quickly and easily review links to other sources
  - Browse application – allows museum staff, art historians, and the general public to verify and explore the data

# Success?

- 14 additional museums have now released their data as linked data
- Three museums have already learned how to use the tools to create their own mappings
  - Indianapolis Museum of Art
  - Smithsonian Archives of American Art
  - Colby College Museum of Art
- Researchers outside the project have applied their visualization tools to the data
  - Sajan Raj Ojha (Univ. of Trento): SemUI
  - Thorsten Liebig (derivo): SemSpect http://aac.semspect.de/

# Future Work

- Automate the addition of new museums to the AAC
  - Gather, map, and link the data directly form their online web pages
- Extend the types of information supported
  - E.g., exhibition data & bibliographies
  - Improve the ability in Karma to automate complex mappings
- Link the existing data to other sources
  - E.g., VIAF, Geonames, & DBpedia
  - Build a library of linking functions to support easy and scalable linking

# More Info

Karma: karma.isi.edu

AAC: americanartcollaborative.org

Github: github.com/american-art

- Thanks to the Mellon Foundation & Institute of Museum and Library Services for their financial support of this project

# Thanks!