

ENABLING DATA DISCOVERY & ANALYTICS THROUGH CURATION IN A FEDERATED ENVIRONMENT

Presenter: Kathy Rondon



WHY A FEDERATED SYSTEM?

- AUTHORITY IS “DIVIDED BETWEEN A CENTRAL...AUTHORITY AND SMALLER, LOCALLY AUTONOMOUS UNITS...USUALLY FORMED THROUGH THE POLITICAL UNION OF TWO OR MORE FORMERLY INDEPENDENT STATES UNDER ONE SOVEREIGN GOVERNMENT THAT DOES NOT, IN ANY CASE, ARROGATE THE INDIVIDUAL POWERS OF THOSE STATES.”
- A FEDERATED SYSTEM ALLOWS BOTH AUTONOMY IN SOME RESPECTS AND CENTRALIZATION IN OTHERS.
- LARGE ORGANIZATIONS GENERALLY BENEFIT FROM SUCH A FEDERATED GOVERNANCE SYSTEM.

WHAT TO CENTRALIZE?

- THE DATA ASSETS THEMSELVES ALMOST NEVER NEED TO BE CENTRALIZED, STORED OR ACCESSED IN A SINGLE REPOSITORY
- CENTRALIZED TASKS OR CONTENT: REFERENCE DATA, MASTER DATA, AND/OR METADATA
- THE OPEN ARCHIVAL INFORMATION SYSTEM IS ONE EXAMPLE OF A FEDERATED ARCHITECTURE



- BORROWING FROM OAIS TENETS TO CREATE A FEDERATED GOVERNANCE ENVIRONMENT:
 - NEGOTIATE WITH AND ACCEPT DATA FROM DATA PROVIDERS
 - OBTAIN SUFFICIENT CONTROL OF DATA ASSETS TO ENSURE LONG-TERM PRESERVATION
 - ENSURE THAT DATA IS UNDERSTANDABLE TO THE USER COMMUNITY WITHOUT HAVING TO CONSULT WITH DATA CREATORS OR PROVIDERS
 - ***FOLLOW DOCUMENTED POLICIES AND PROCEDURES WHICH ENABLE INFORMATION TO BE DISSEMINATED AS AUTHENTICATED COPIES AND/OR TRACEABLE TO THE ORIGINAL***

- THE SUBMISSION INFORMATION PACKAGE (SIP): INFORMATION THAT PROVIDERS SEND TO THE CENTRAL REPOSITORY
- THE DISSEMINATION INFORMATION PACKAGE (DIP): INFORMATION THE CENTRAL REPOSITORY SENDS TO DATA REQUESTERS
- THE ARCHIVAL INFORMATION PACKAGE (AIP): INFORMATION STORED BY THE CENTRAL REPOSITORY FOR PRESERVATION AND ARCHIVAL PURPOSES

- A DISCIPLINE THAT ENABLES DATA DISCOVERY AND RETRIEVAL, MAINTAINING ITS QUALITY, ADDING VALUE, AND PROVIDING FOR REUSE OVER TIME.*
- THE ONGOING DOCUMENTATION OF **CONTEXTUAL METADATA** ABOUT KEY DATA ASSETS ACCORDING TO STANDARD OPERATING PROCEDURES AND CONSISTENT TERMS OF REFERENCE

*Source: The University of Illinois
School of Library and Information
Sciences

CURATION provides the appropriate procedures to help develop conflict resolution strategies, data provenance, data consistency, data reliability of a dataset and identify gaps in data for possible additional collection and/or articulation of caveats and margins of error

- DEVELOP AND MAINTAIN TERMS OF REFERENCE FOR DESCRIBING YOUR ORGANIZATION'S DATA.
- DOCUMENT AND COMMUNICATE TO USERS DATA HANDLING POLICIES (SUCH AS RESTRICTIONS ON USE, COPYRIGHTS AND LICENSES).
- DETERMINE THE METADATA (STRUCTURED AND UNSTRUCTURED) THAT YOUR ORGANIZATION'S DATA USERS NEED—AND THINK AHEAD TO FUTURE REUSABILITY OF THE DATA.
- TRAIN YOUR WORKFORCE ON CONSISTENT DOCUMENTATION STANDARDS.
- DATA CURATION IS ESSENTIALLY A DATA MANAGEMENT MATURITY ISSUE.

- BEST PRACTICES
 - **TRY NOT TO BOIL THE OCEAN:** CHOOSE THE MOST IMPORTANT AND POSSIBLY MOST MISUNDERSTOOD TERMS AND FUNCTIONS AND DEFINE THEM IN AS STRAIGHTFORWARD AND SIMPLE LANGUAGE POSSIBLE
 - **COORDINATE:** ENSURE THAT STAKEHOLDERS AGREE ON THE TERMS OF REFERENCE, BUT DON'T LET A SINGLE HOLD OUT TAKE YOUR DATA DICTIONARY HOSTAGE
 - **COMMUNICATE:** ONCE YOU HAVE A DATA DICTIONARY, DON'T STICK IT IN A DRAWER OR KEEP IT A SECRET; PUBLISH IT, SHARE IT AND TRAIN IT IN AS MANY FORUMS AS POSSIBLE

- TO CONSISTENTLY CURATE DATA, THE ORGANIZATION MUST HAVE A METADATA STANDARD, THE CONTEXTUAL METADATA FIELDS THAT APPLY TO EVERY CURATED DATASET
- CONTEXTUAL METADATA CAN BE STRUCTURED, UNSTRUCTURED, OR BOTH, DEPENDING ON THE NEEDS OF THE ORGANIZATION
- IF A RELEVANT METADATA STANDARD ALREADY EXISTS, DON'T REINVENT THE WHEEL, BUT DON'T BE AFRAID TO CREATE SOMETHING NEW IF YOU NEED TO

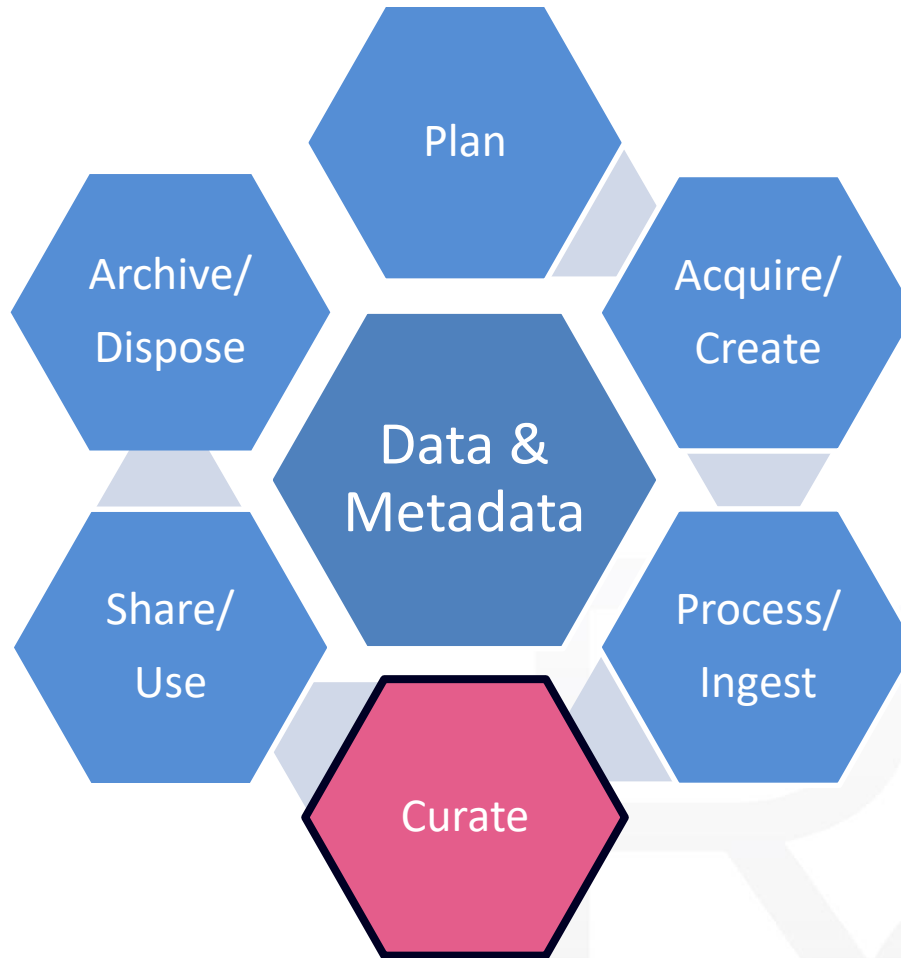
“STANDARDS ARE LIKE TOOTHBRUSHES. EVERYONE AGREES THAT THEY’RE A GOOD IDEA, BUT NOBODY WANTS TO USE ANYONE ELSE’S.” *

* FROM *METADATA* BY JEFFREY POMERANTZ; ATTRIBUTED TO MURTHA BACA, GETTY RESEARCH INSTITUTE

ALL WILL BE FOR NAUGHT IF...

- DATA CURATION: EVERYONE THINKS IT'S A GREAT IDEA...FOR SOMEONE ELSE TO DO.
- TRAINING AND PROFESSIONALIZATION OF A CURATION WORKFORCE IS KEY TO THE SUCCESS OF A CURATION EFFORT
- SKILL SET OF A DATA CURATOR IS ONE THAT BRIDGES THE GAP BETWEEN IT AND BUSINESS OR MISSION.
 - A CURATOR SHOULD UNDERSTAND DATA FORMATS AND DATA EXPLOITATION SYSTEMS AT A BASIC LEVEL.
 - THE CURATOR MUST HAVE EXCELLENT WRITING AND COMMUNICATION SKILLS
 - THE CURATOR MUST IN DEPTH UNDERSTANDING OF THE INTENDED BUSINESS OR MISSION USES OF THE DATA.

CURATION IN THE DATA LIFECYCLE



QUESTIONS?

