

# Quality Control and Peer Review of Data Sets:

## Mapping Data Archiving Processes to Data Publication Requirements

Matthew S. Mayernik<sup>1</sup>, Mike Daniels<sup>1</sup>, Christopher Eaker<sup>2</sup>, Gary Strand<sup>1</sup>, Steven F. Williams<sup>1</sup>, Steven J. Worley<sup>1</sup>

1. National Center for Atmospheric Research (NCAR), 2. University of Tennessee-Knoxville, School of Information Sciences



### Data Peer Review

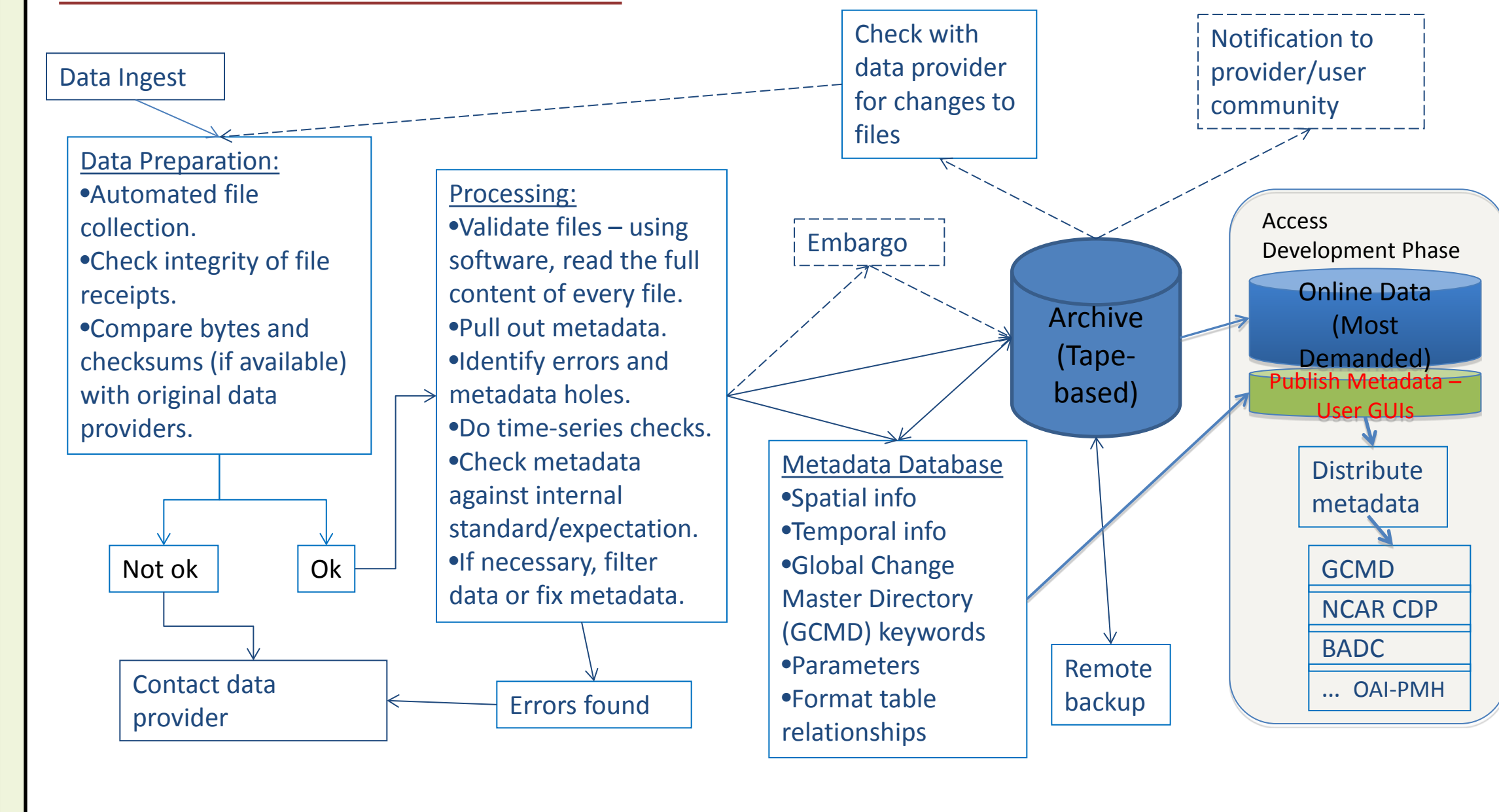
- Peer review is a central way to assess research quality
- How does peer review apply to data publication and citation initiatives?
- Peer review of growing volumes of digital data will increase the stress on the scholarly publication system.

### Issues

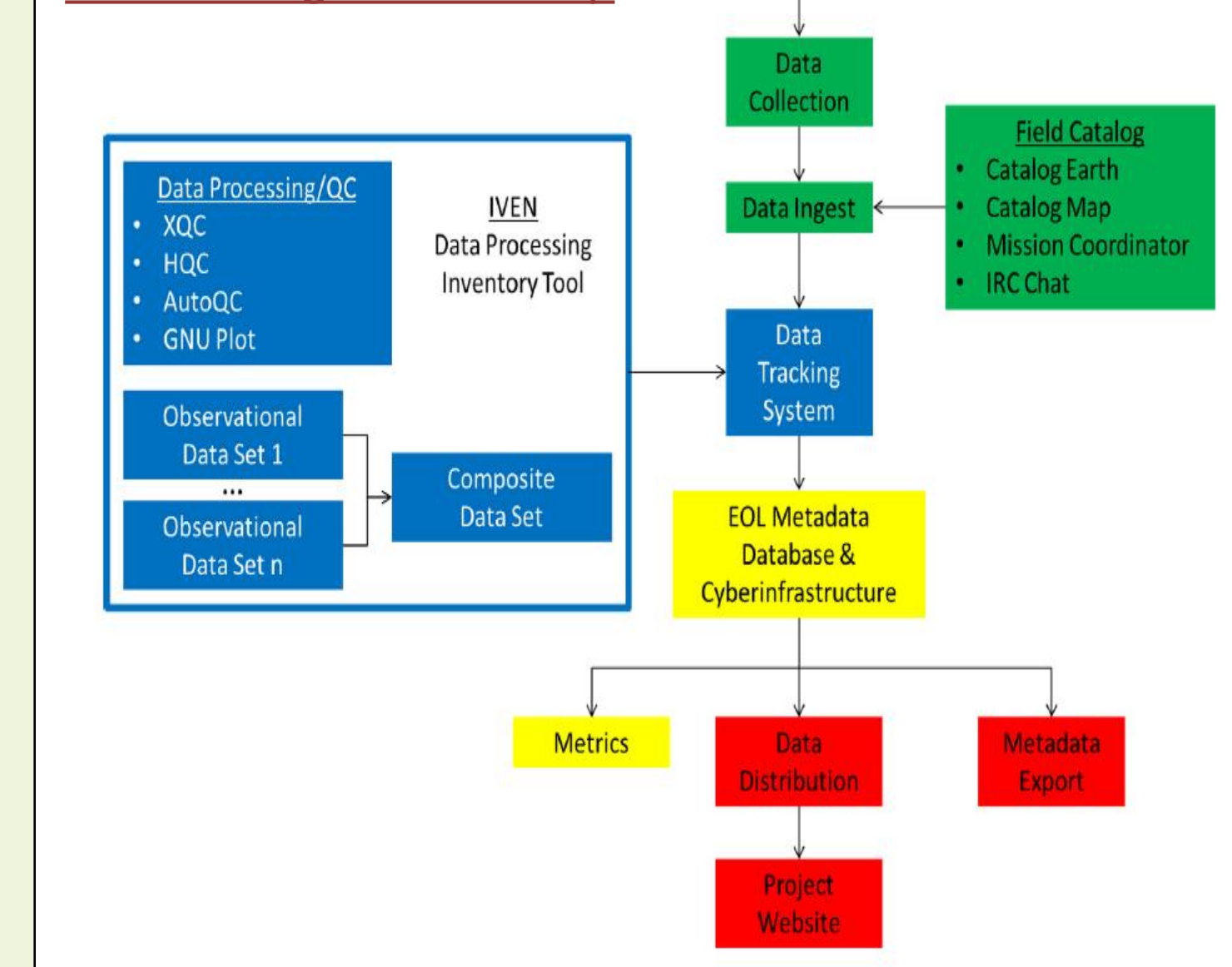
- Data QC processes and software are very specific to data types, experimental designs, and systems
- Human examination vs. Automated review
  - Human examination of data access interfaces, documentation, and metadata is essential to assess suitability for users.
  - Visual exams of data and metadata characteristics are often very important to identify systematic flaws.
  - If the data and metadata are published in standard form, readily available tools can be used to automate some data evaluation.
  - When possible, automation is desired to reduce the time and effort on the part of the human reviewer.
- Research timelines
  - Pre-publication review vs. post-publication review
    - Data users commonly find data errors that can only be found through intensive analysis
    - Repositories must have way to receive, evaluate, and respond to user-discovered errors
    - Reviewers from outside of a project need more time and often assistance from project members
  - There is a growing demand for real-time data. Checks on real-time data quality can be done initially, but quality control timelines have to be responsive to researchers' desire to access and use the data.
- Peer review might be best conceptualized as review of the data collection, assessment, metadata, and archiving processes vs. review of the data themselves.

### NCAR Data Repository Archiving Workflows

#### NCAR CISL Research Data Archive



#### NCAR Earth Observing Lab Data Management Group



Repository workflows have data quality assessment processes integrated throughout the data ingest and archiving workflows

### Repository Data Quality Control Processes

- Flag questionable or faulty data by creating new metadata. Always maintain original data.
- Provide mechanisms for feedback loops between users, the archive, and data providers.
  - Sometimes data quality problems are found by external users. External users are excellent data reviewers.
  - People who are knowledgeable about the project are more likely to find actual problems with the model and data, whereas users are likely to find smaller scale anomalies that may or may not be errors. Shared evaluation is sometimes required.
- Compare data with other data, or model runs with other data/models
- Develop standard sets of diagnostics tools and methods over time
- Technical review vs. scientific peer review
  - For model data, use control runs to evaluate the model functionality
  - For observational field campaign data, keep "housekeeping parameters", like battery life, ambient or equipment temperature ranges, etc., to evaluate the equipment functionality
  - Is the data set well constructed, e.g. following conventions and standards?