http://www.massobs.org.uk/annual_report.htm

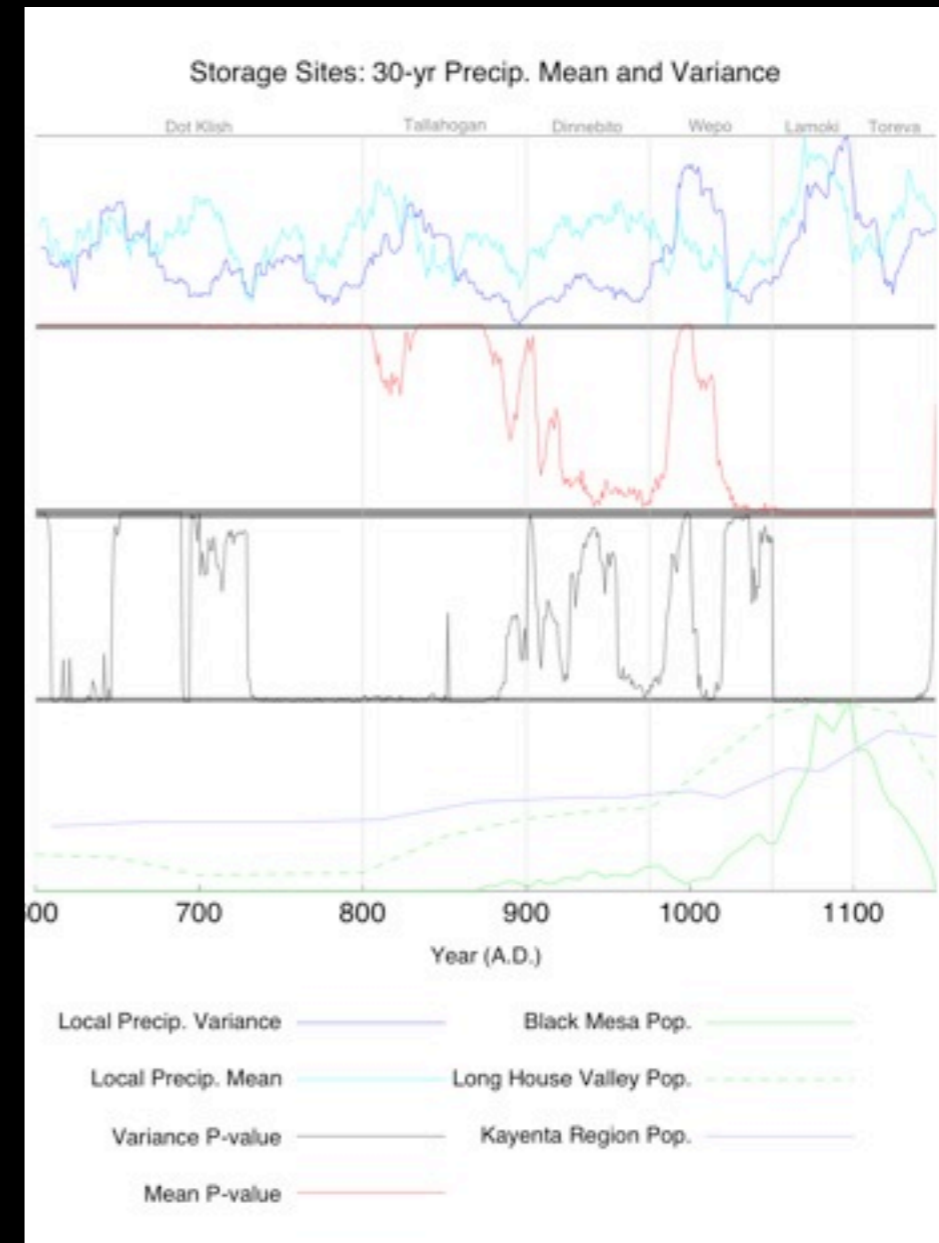http://infoontech.files.wordpress.com/2011/05/archive-shelving.jpg

# Scientific Data Preservation

Karl Benedict

EDAC, University Libraries, Dept. of Geography

University of New Mexico

# Why Me?

- Earth Data Analysis Center

- University Libraries

- Geography Department

- Federation of Earth Science Information Partners

- Foundation for Earth Science

# A Story

- Dissertation work with two large data collections from the late-60s through early 80s

- Over 2200 archaeological sites, paleoclimate and modern meteorological data, publicly available environmental data

- Focus on data integration and modeling - not feasible without well documented data in understandable and usable formats



Storage Sites: 30-yr Precip. Mean and Variance

# Roadmap

- Core principles

- Strategies

- Resources

# Core Principles

- Data quality/safety
  - During research
  - Following research
- Documentation
  - Discovery
  - Use
  - Understanding
- Sustainability
  - Data and metadata formats/standards

# Data Quality and Safety



*I learn with great satisfaction that you are about committing to the press the valuable historical and State papers you have been so long collecting.* Time and accident are committing daily havoc on the originals deposited in our public offices. *The late war has done the work of centuries in this business. The last cannot be recovered, but* let us save what remains; not by vaults and locks which fence them from the public eye and use in consigning them to the waste of time, but by such a multiplication of copies, as shall place them beyond the reach of accident.

Letter to Ebenezer Hazard
Philadelphia, 2/18/1791.
Thomas Jefferson

# Documentation



*Data creators will provide sufficient metadata (defined as all the information necessary for data to be independently understood by users and to ensure proper stewardship of the data) to the data repositories responsible for long-term archival.*

- Interagency Data Stewardship Guidelines. http:// commons.esipfed.org/node/419

# Sustainability

# Sustainability

Seven sustainability factors



- Library of Congress - http://www.digitalpreservation.gov/formats/intro/format_eval_rel.shtml

# Sustainability

Seven sustainability factors

1. <u>Disclosure</u>.  Degree to which complete specifications and tools for validating technical integrity exist and are accessible to those creating and sustaining digital content. A spectrum of disclosure levels can be observed for digital formats.  What is most significant is not approval by a recognized standards body, but the existence of complete documentation.

# Sustainability

Seven sustainability factors

1. Disclosure.  Degree to which complete specifications and tools for validating technical integrity exist and are accessible to those creating and sustaining digital content. A spectrum of disclosure levels can be observed for digital formats.  What is most significant is not approval by a recognized standards body, but the existence of complete documentation.

2. Adoption.  Degree to which the format is already used by the primary creators, disseminators, or users of information resources.  This includes use as a master format, for delivery to end users, and as a means of interchange between systems.

# Sustainability

Seven sustainability factors

1. <u>Disclosure</u>.  Degree to which complete specifications and tools for validating technical integrity exist and are accessible to those creating and sustaining digital content. A spectrum of disclosure levels can be observed for digital formats.  What is most significant is not approval by a recognized standards body, but the existence of complete documentation.

2. <u>Adoption</u>.  Degree to which the format is already used by the primary creators, disseminators, or users of information resources.  This includes use as a master format, for delivery to end users, and as a means of interchange between systems.

3. <u>Transparency</u>.  Degree to which the digital representation is open to direct analysis with basic tools, such as human readability using a text-only editor.

- Library of Congress - http://www.digitalpreservation.gov/formats/intro/format_eval_rel.shtml

# Sustainability

Seven sustainability factors

1. Disclosure.  Degree to which complete specifications and tools for validating technical integrity exist and are accessible to those creating and sustaining digital content. A spectrum of disclosure levels can be observed for digital formats.  What is most significant is not approval by a recognized standards body, but the existence of complete documentation.

2. Adoption.  Degree to which the format is already used by the primary creators, disseminators, or users of information resources.  This includes use as a master format, for delivery to end users, and as a means of interchange between systems.

3. Transparency.  Degree to which the digital representation is open to direct analysis with basic tools, such as human readability using a text-only editor.

4. Self-documentation.  Self-documenting digital objects contain basic descriptive, technical, and other administrative metadata.

- Library of Congress - http://www.digitalpreservation.gov/formats/intro/format_eval_rel.shtml

# Sustainability

Seven sustainability factors

1. <u>Disclosure</u>.  Degree to which complete specifications and tools for validating technical integrity exist and are accessible to those creating and sustaining digital content. A spectrum of disclosure levels can be observed for digital formats.  What is most significant is not approval by a recognized standards body, but the existence of complete documentation.

2. <u>Adoption</u>.  Degree to which the format is already used by the primary creators, disseminators, or users of information resources.  This includes use as a master format, for delivery to end users, and as a means of interchange between systems.

3. <u>Transparency</u>.  Degree to which the digital representation is open to direct analysis with basic tools, such as human readability using a text-only editor.

4. <u>Self-documentation</u>.  Self-documenting digital objects contain basic descriptive, technical, and other administrative metadata.

5. <u>External Dependencies</u>.  Degree to which a particular format depends on particular hardware, operating system, or software for rendering or use and the predicted complexity of dealing with those dependencies in future technical environments.

- Library of Congress - http://www.digitalpreservation.gov/formats/intro/format_eval_rel.shtml

# Sustainability

Seven sustainability factors

1. <u>Disclosure</u>.  Degree to which complete specifications and tools for validating technical integrity exist and are accessible to those creating and sustaining digital content. A spectrum of disclosure levels can be observed for digital formats.  What is most significant is not approval by a recognized standards body, but the existence of complete documentation.

2. <u>Adoption</u>.  Degree to which the format is already used by the primary creators, disseminators, or users of information resources.  This includes use as a master format, for delivery to end users, and as a means of interchange between systems.

3. <u>Transparency</u>.  Degree to which the digital representation is open to direct analysis with basic tools, such as human readability using a text-only editor.

4. <u>Self-documentation</u>.  Self-documenting digital objects contain basic descriptive, technical, and other administrative metadata.

5. <u>External Dependencies</u>.  Degree to which a particular format depends on particular hardware, operating system, or software for rendering or use and the predicted complexity of dealing with those dependencies in future technical environments.

6. <u>Impact of Patents</u>.  Degree to which the ability of archival institutions to sustain content in a format will be inhibited by patents.

# Sustainability

Seven sustainability factors

1. <u>Disclosure</u>.  Degree to which complete specifications and tools for validating technical integrity exist and are accessible to those creating and sustaining digital content. A spectrum of disclosure levels can be observed for digital formats.  What is most significant is not approval by a recognized standards body, but the existence of complete documentation.

2. <u>Adoption</u>.  Degree to which the format is already used by the primary creators, disseminators, or users of information resources.  This includes use as a master format, for delivery to end users, and as a means of interchange between systems.

3. <u>Transparency</u>.  Degree to which the digital representation is open to direct analysis with basic tools, such as human readability using a text-only editor.

4. <u>Self-documentation</u>.  Self-documenting digital objects contain basic descriptive, technical, and other administrative metadata.

5. <u>External Dependencies</u>.  Degree to which a particular format depends on particular hardware, operating system, or software for rendering or use and the predicted complexity of dealing with those dependencies in future technical environments.

6. <u>Impact of Patents</u>.  Degree to which the ability of archival institutions to sustain content in a format will be inhibited by patents.

7. <u>Technical Protection Mechanisms</u>.  Implementation of mechanisms such as encryption that prevent the preservation of content by a trusted repository.

- Library of Congress - http://www.digitalpreservation.gov/formats/intro/format_eval_rel.shtml

# Dataset Considerations

*Significant for all datasets is that they be represented in a structure that reveals the characteristics of individual data items and the relationships among them. A dataset format suitable for preservation must retain the syntactical integrity of both the **structure** and individual values, so that automated analysis is possible. Also essential for future usability is an understanding of the semantics of the data elements and their relationships within the dataset. The semantics may be described explicitly within the dataset, described explicitly in an ancillary document (preferably itself machine-processable), or implicit through compliance with a community best practice or external specification.*

http://www.digitalpreservation.gov/formats/content/dataset_quality.shtml

# Dataset Considerations

*Significant for all datasets is that they be* **represented in a structure** *that reveals the characteristics of individual data items and the* **relationships** *among them. A dataset format suitable for preservation must retain the* **syntactical integrity** *of both the* **structure** *and* **individual values**, *so that automated analysis is possible. Also essential for future usability is an* **understanding of the semantics** *of the data elements and their relationships within the dataset. The semantics may be described explicitly within the dataset, described explicitly in an ancillary document (preferably itself machine-processable), or implicit through compliance with a community best practice or external specification.*

http://www.digitalpreservation.gov/formats/content/dataset_quality.shtml

# Dataset Considerations

Significant for all datasets is that they be represented in a structure that reveals the characteristics of individual data items and the relationships among them. A dataset format suitable for preservation must retain the syntactical integrity of both the structure and individual values, so that automated analysis is possible. Also essential for future usability is an understanding of the semantics of the data elements and their relationships within the dataset. The semantics may be described explicitly within the dataset, described explicitly in an ancillary document (preferably itself machine-processable), or implicit through compliance with a community best practice or external specification.

represented in a structure

relationships

syntactical integrity

structure

individual values

understanding of the semantics

http://www.digitalpreservation.gov/formats/content/dataset_quality.shtml

# Strategies

# Preservation Resources



- Use existing individual resources
- Leverage existing shared resources
- Acquire new individual or shared resources
- Use commodity/hosted (e.g. cloud) resources

# Use Existing Resources

- Pros
  - Little if any initial cost
  - Resources already well understood
- Cons
  - Potential lack of alignment with preservation needs
  - May be insufficient for long-term



http://www.flickr.com/photos/27784972@N07/

# Shared Resources



http://www.flickr.com/photos/ben_grey/

- Pros
  - May be able to better use underutilized resource
  - Resource is a known commodity
- Cons
  - Resource decisions are shared
  - Potential contention for resource as it is consumed

# Acquire New Resources

- Pros
  - Resources can be specified for need
  - Growth may be built into resource planning
- Cons
  - Requires expenditure of funds
  - Adequate funds may not be available for acquisition
  - New resource may have related costs (i.e. labor, administration)

# Hosted Resources

- Pros

  - Cost can scale with demand/need

  - Local administrative and infrastructure costs can be reduced

- Cons

  - Reduced control over data

  - Risk of provider going out of business

  - May not exactly match needs

# Documentation

- Start Early
- Elements to Record
  - Who
  - Where
  - When
  - What
  - Why
  - How
- Know target documentation content standard at start to ensure coverage

# Format Selection

# Resources

- UNM LoboVault / UNM Libraries data curators

- UNM Research Storage Consortium

- NM Resource Geographic Information System / NM EPSCoR

- DataONE Member Node(s), OneShare

- Community Repositories

- Further Reading



http://www.flickr.com/photos/kalexanderson/

# UNM LoboVault & Data Curators



http://libguides.unm.edu/data

http://repository.unm.edu

# UNM Research Storage Consortium

# NM RGIS & EPSCoR Data Portals



http://rgis.unm.edu

http://nmepscor.org/dataportal

# DataONE



http://www.dataone.org/

# Community Repositories



Check out Databib for a list of over 300 repositories:
http://databib.org/index.php

# Additional Background Material

- ESIP Federation *Interagency Data Stewardship Guidelines*
  http://commons.esipfed.org/node/419

- ESIP Federation *Data Citation Guidelines for Data Providers and Archives*
  http://commons.esipfed.org/node/308

- Library of Congress *Sustainability of Digital Formats Planning for Library of Congress Collections*
  http://www.digitalpreservation.gov/formats/index.shtml

- UNM Libraries *Digital Data Management, Curation and Archiving* Research Guide
  http://libguides.unm.edu/data

- DataONE *Best Practices database*

  http://www.dataone.org/best-practices

Questions?