# Format Standards:

What Do I Need To Know?

# Overview for Today:

1. **What are Formats** and Why Should You Care?

2. Formats and Preservation Strategies for **Textual Records, Structured Records, and Email**.

3. Brief Overview of **Strategies and Available Resources**.

**BORN-DIGITAL**

**DIGITIZED**

born dig·it·al (adjective)
digital information
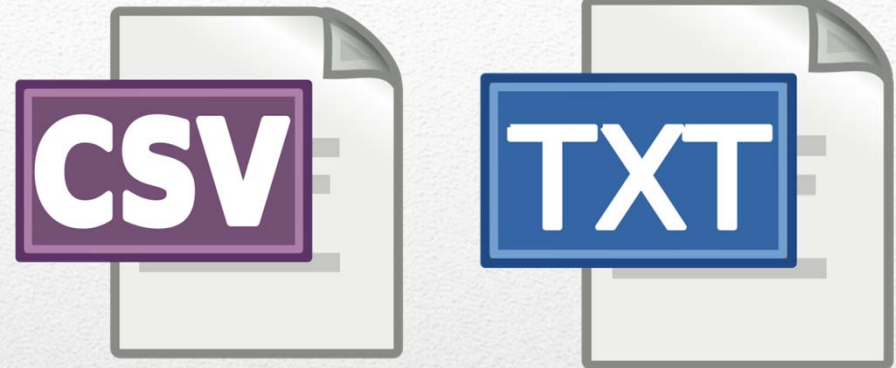originally created
in electronic form

**FORMAT DECISIONS AND ISSUES APPLY TO BOTH!**

# Understanding Formats

**MEDIA**

**CONTENT**



# Understanding Formats

File formats are the "rules that specify how the bytes that make up [a given] file are organized, interpreted, and rendered."

- Ciran B. Trace

# Understanding Formats

# The Biggest Risk with Formats?



## Understanding Formats

**Ask questions about format choices now…**



**…to help insure future access.**

# Understanding Formats

# QUESTION YOUR FORMATS

**WHO OWNS?**

**WHO SUPPORTS?**

**WHO (OR WHAT) GOVERNS?**

# Understanding Formats

## NATIONAL ARCHIVES

- Refers to born digital textual data.

- Digitized text is a concern for image formats.

- Generally two types:
  - Formatted
  - Unformatted

### 7. Textual Data

The textual data category refers to two general content types: unformatted (plain text) or formatted. Unformatted plain text (defined in MIME as text/plain) contains basic character information and control or non-printing characters but lacks styling information. Formatted text files include all of the attributes of plain text files but have extended formatting capabilities, for "stylized" or "rich" text features including italics, bold, colors, hyper-linking, etc.

Agencies must identify the character encoding method used with each text file.

**Preferred Formats**

| Preferred Formats | Format Versions | Format Specifications |
|---|---|---|
| ASCII Text | 7 bit | ISO/IEC 646:1991 Information technology -- ISO 7-bit coded character set for information interchange: (http://www.iso.org/iso/catalogue_detail.htm?csnumber=4777) |
| Unicode Text | UTF-8 <br><br> UTF-16 | RTF 3629: UTF-8, A Transformation Format of ISO 10646: (http://tools.ietf.org/html/rfc3629) <br><br> RFC 2781 UTF-16: An Encoding of ISO 10646: (http://www.ietf.org/rfc/rfc2781.txt) |
| OpenDocument Text Format (ODF) | OpenDocument 1.0 | ISO/IEC 26300:2006 Information technology -- OpenDocument Format for Office Applications (OpenDocument) v1.0: (http://www.iso.org/iso/iso_catalogue/catalogue_tc/catalogue_detail.htm?csnumber=43485) |
| PDF/A-1 | PDF/A-1 | ISO 19005-1:2005 Document management -- Electronic document file format for long-term preservation -- Part 1: Use of PDF 1.4 (PDF/A-1): (http://www.iso.org/iso/catalogue_detail?csnumber=38920) |
| PDF/A-2 | PDF/A-2 | ISO 19005-2:2011 Document management -- Electronic document file format for long-term preservation -- Part 2: Use of ISO 32000-1 (PDF/A-2): (http://www.iso.org/iso/home/store/catalogue_tc/catalogue_detail.htm?csnumber=50655) |

# Textual Records

NATIONAL ARCHIVES

- Concerns data that is organized and stored in defined fields. This can include:
  - Databases
  - Spreadsheets
  - Statistical Data
  - Scientific data

- Typically requires transfer of other associated files and metadata (such as schemas or data dictionaries) that help make the structured data valid and comprehensible.

## 8. Structured Data Formats

Structured data comprises the broad category of data that is stored in defined fields. Categories for structured data are as follows:

- Database formats are organized collections of associated data that conform to a logical structure. Database formats are determined by "data models" that describe specific data structures used to model an application and generally include navigational, relational, and hybrid models;

- Spreadsheets are tables made up of columns and rows and which contain cells of data. Relationships between cells can be pre-defined as mathematical formulas;

- Statistical data is the result of quantitative research and analysis. Statistical data formats contain collections of data presented in both tabular and non-tabular form; and

- Scientific data refers to research data collected by instrumentation tools during the scientific process. Scientific data formats are either domain specific within a single field of study, or are multi-domain formats used for transfer of scientific data between domains.

General requirements for structured data include the following:

- Agencies must transfer structured data that is both well-formed according to the syntactical conventions of the format, and valid according to the structural rules defined in any associated schemas or document type definitions (DTDs);

- Value Separated Files, e.g. CSV or comma separated value files, may use a character other than the comma. The pipe or caret are recommended delimiters because they are not commonly found in free text fields. Alternatively, text files encoded with ASCII characters and where each field is a fixed width, is also an acceptable transfer format for use with structured data, even though ASCII is technically a data encoding type. ASCII text files must be accompanied by complete documentation of the record lengths and field widths;

- Data files and databases shall be transferred as flat files or as rectangular tables, that is, as two-dimensional arrays, lists or tables. All records in a database, or rows (tuples) in a relational database, should have the same logical format. Each data element within a record should contain only one data value. A record should not contain nested repeating groups of data items; and

- Structured data must be transferred together with any associated files necessary to verify the validity of the data, e.g., DTDs, schemas, and data dictionaries.

# Structured Records

COMMON FORMATS

# Structured Records

**NATIONAL ARCHIVES**

- Consider the traditional functions email has come to replace.

- Critical that email is preserved and transferred according to set retention schedules and/or collection policies.

- Email is unique in how it can seem to transition between textual record and structured record.

# Email

## 9. Email

Email is defined as discrete electronic communications transmitted over the Simple Mail Transfer Protocol (SMTP), between two or more people or entities, in compliance with applicable IETF's Request for Comments (RFC) specifications. Email does not include other functions commonly available via email programs such as calendars, tasks, appointments, newsgroups, or instant messaging. In order for information in a calendar, contact list, address book etc. to be transferred to NARA, it must be scheduled as a separate item.

Please note that NARA considers email attachments to be a component of the email record and does not require that unseparated email attachments meet the transfer standards specified by the format category under which the attachment alone would fall.

General requirements for email:

- Transfers of email records must consist of an identifiable, organized body of records (not necessarily a traditional series);

- Email messages should include delimiters that indicate the beginning and end of each message and the beginning and end of each attachment, if any. Each attachment must be differentiated from the body of the message, and uniquely identified;

- Email messages transferred as XML files must be accompanied by any associated document type definitions (dtds), schemas, and/or data dictionaries;

- Labels to identify each part of the message (Date, To [all recipients, including cc: and bc: copies], From, Subject, Body, and Attachment) including transmission and receipt information (Time Sent, Time Opened, Message Size, File Name, and similar information, if available). To ensure identification of the sender and addressee(s), agencies that use an email system that identifies users by codes or nicknames, or identifies addressees only by the name of a distribution list should include information with the transfer-level documentation; and

- Email converted to formats not natively used by the email program, and which do not maintain header information (such as RTF or Word documents) are not accepted. Printouts of emails are also not accepted under this Bulletin.
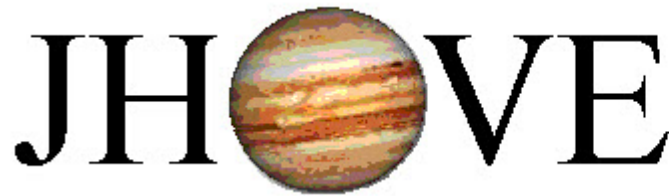
# COMMON FORMATS
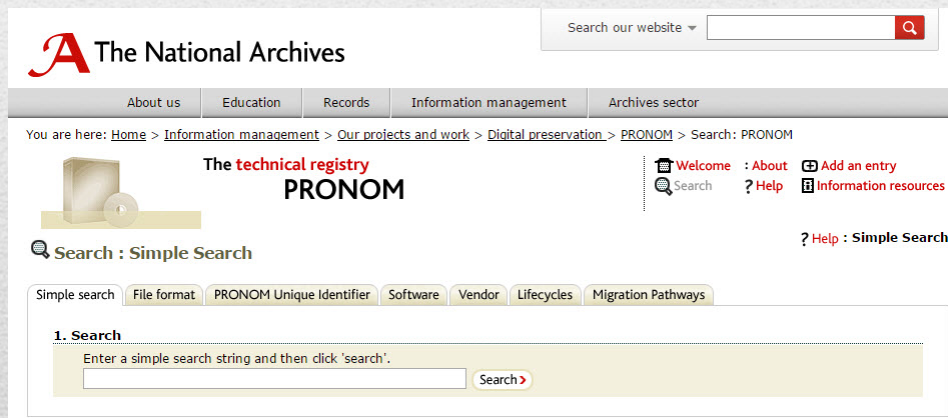


# Email

# Policies & Procedures



# Strategies and Resources

# Validation Tools



**http://jhove.sourceforge.net/**





**https://www.nationalarchives.gov.uk/PRONOM/Default.aspx**

# Strategies and Resources

# Open Source Software



https://www.openoffice.org/



https://sourceforge.net/projects/pdfcreator/

# Strategies and Resources

# Professional Guidance
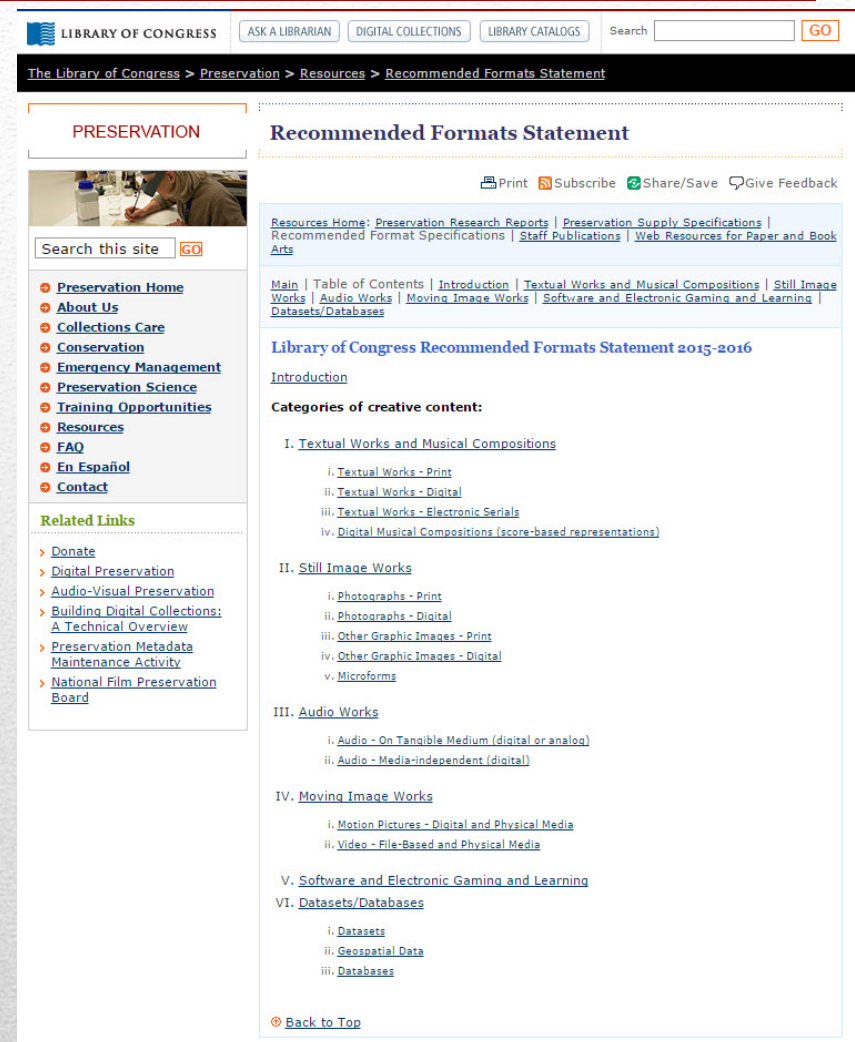


## NARA 2014-04: Appendix A

- **Provides format standards for a variety of digital objects encountered in an archival digital preservation environment.**

- **Tiered categories for formats that provides greater flexibility.**

- **https://www.archives.gov/records-mgmt/policy/transfer-guidance-tables.html#presentationformats**

# Strategies and Resources

# Professional Guidance

**LoC Recommended Formats Statement:**

- "...purpose is to inform the creative and library communities on best practices for ensuring the preservation of, and long-term access to, the creative output of the nation and the world."

- Provides format recommendations and standards for both analog and digital materials.

- https://www.loc.gov/preservation/resources/rfs/



# Strategies and Resources

## Professional Guidance



**Sustainability of Digital Formats**
**Planning for Library of Congress Collections**

Introduction | Sustainability Factors | Content Categories | Format Description

The Digital Formats Web site provides information about digital content for Fleischhauer, and Kate Murray invite feedback on the content.

### Introduction
Background information and overview: What is a format? How shall we evaluate formats? What projects in other organizations are addressing these questions? >>

### Sustainability Factors
What affects the ability of the Library to preserve content in a given format? These sustainability factors apply to all formats. >>

### Content Categories
The evaluation of formats must take into account quality and functionality. These factors vary according to the type of content under consideration and the categories will be expanded as time passes. >>

### Format Descriptions
Documents with more information about specific formats. >>

**LoC Sustainability of Digital Formats**

- **Useful tool for thinking about long-term format sustainability, and building corresponding institutional format policies.**

- **Discusses sustainability factors, content categories, and provides descriptions for individual content formats.**

- **http://www.digitalpreservation.gov/formats/**

# Strategies and Resources

# James Kichas
[jkichas@utah.gov](mailto:jkichas@utah.gov)
### 801-531-3844