# Digital File Format Resources at the Library of Congress

**DPC Briefing Day: Preservation Planning and Technology Watch**

**Kate Murray**
**kmur@loc.gov**

# Sustainability of Digital Formats: Planning for Library of Congress Collections

The Sustainability of Digital Formats Web site provides information about digital content formats. The analyses and resources presented here will increase and be updated over time. The compilers, Caroline R. Arms, Carl 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?

Overview | Formats, Evaluation Factors, and Relationships | Papers and Presentations | Related Resources

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

Disclosure | Adoption | Transparency | Self-documentation | External Dependencies | Impact of Patents | Technical Protection Mechanisms

## 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.

Still Image | Sound | Textual | Moving Image | Web Archive | Datasets | Geospatial | Generic | Browse All Formats
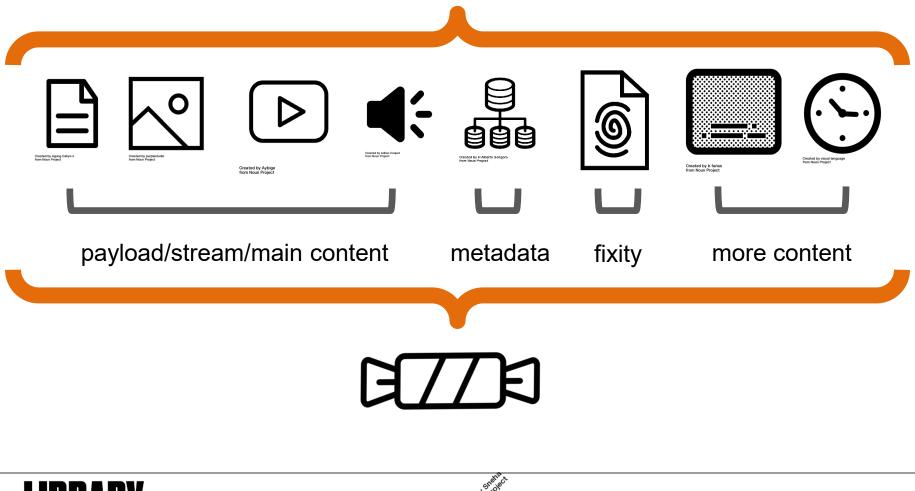
## Format Descriptions
Documents with more information about specific formats.

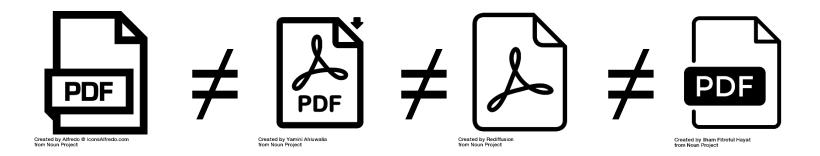Browse categories | Browse alphabetical list | Format Descriptions as XML

http://www.loc.gov/preservation/digital/formats/index.html

**LIBRARY** LIBRARY OF CONGRESS

# what's in a file?



payload/stream/main content     metadata    fixity     more content

**even files that look the same, with the same file extension, can be very different**

**pdf**

> mother of all container formats
> its flexibility is its curse

**pdf/a** ("archival" friendly. maybe? kinda?)

> no encryption
> no javascript
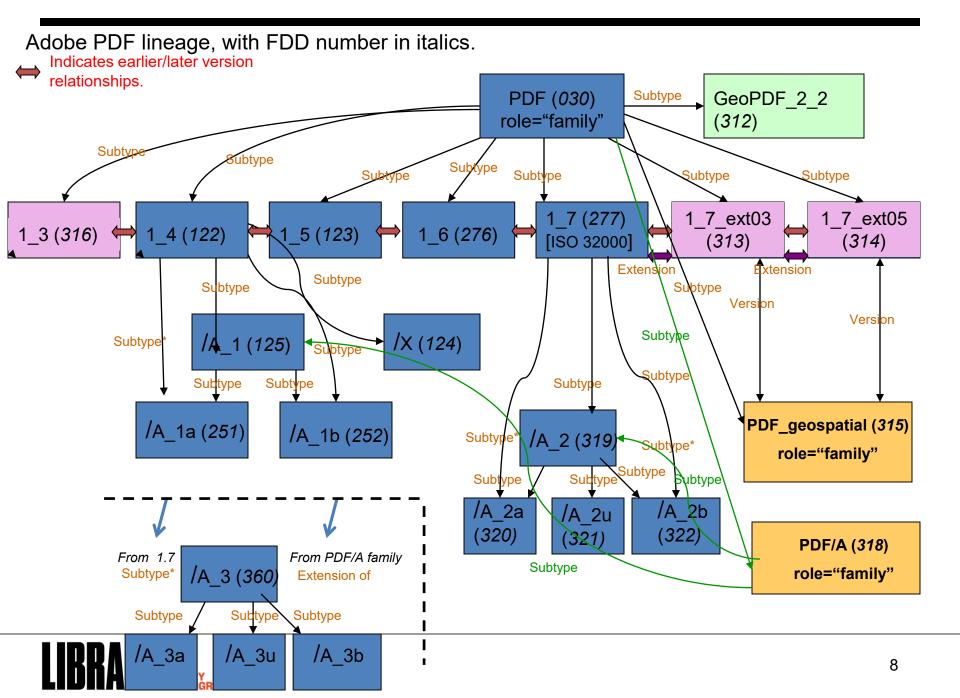> audio/video only permitted in pdf/a-3

**pdf/x**

**pdf/ua**

# relationships are hard.

has subtype | subtype of | contains | may contain | used by | must have component | may have component | component of | defined via | requires | modification of | has modified version | extension of | has extension | has earlier version | has later version | version of | equivalent to | affinity to | additional | other | yada yada yada ...

# PDF is much more than text -- a file format, a wrapper, a bundling format, all in one – note complexity of relationships

| | |
|---|---|
| Has subtype | PDF_1_3, PDF Versions 1.0-1.3 |
| Has subtype | PDF_1_4, PDF Version 1.4 |
| Has subtype | PDF_1_5, PDF, Version 1.5 |
| Has subtype | PDF_1_6, PDF, Version 1.6 |
| Has subtype | PDF_1_7, PDF, Version 1.7 (ISO 32000-1:2008) |
| Has subtype | PDF_1_7_ext03, PDF, Version 1.7, ExtensionLevel 3 |
| Has subtype | PDF_1_7_ext05, PDF, Version 1.7, ExtensionLevel 5 |
| Has subtype | PDF_2_0, PDF, Version 2.0, ISO 32000-2 (2017, 2020) |
| Has subtype | PDF/A_family, PDF for Long-term Preservation. As of November 2012, there are three chronological versions of PDF/A. |
| Has subtype | PDF/A-1, PDF for Long-term Preservation, Use of PDF 1.4 |
| Has subtype | PDF/A-2, PDF/A-2 for Long-term Preservation, Use of ISO 32000-1 (PDF 1.7) |
| Has subtype | PDF/A-3, PDF/A-3 for Long-term Preservation, Use of ISO 32000-1 (PDF 1.7), with Embedded Files |
| Has subtype | PDF/A-4, PDF for Long-term Preservation, Use of ISO 32000-2 (PDF 2.0) |
| Has subtype | PDF/E-1, PDF Engineering Document Format, Use of PDF 1.6 |
| Has subtype | PDF/UA-1, PDF/UA-1, PDF Enhancement for Accessibility, Use of ISO 32000-1 |
| Has subtype | PDF/X, PDF for Prepress Graphics File Exchange |
| Has subtype | PDF/R-1, For raster image transport and storage. Based on PDF 1.4-1.7 (ISO 32000-1) |
| Has subtype | PDF/R-1_enc, For raster image transport and storage. Encrypted, based on PDF 2.0 (ISO 32000-2) |
| Has subtype | GeoPDF_file, GeoPDF File Format (TerraGo) |
| May contain | PDF_geospatial, PDF, Geospatial encoding (Adobe). Supported by version 1.7 ExtensionLevel 3. |
| May contain | GeoPDF_OGC, GeoPDF Encoding (TerraGo 2.2), OGC Best Practice |

# Adobe PDF lineage, with FDD number in italics.

Indicates earlier/later version relationships.

PDF (*030*) role="family"

Subtype → GeoPDF_2_2 (*312*)

Subtype → 1_3 (*316*)

Subtype → 1_4 (*122*)

Subtype → 1_5 (*123*)

Subtype → 1_6 (*276*)

Subtype → 1_7 (*277*) [ISO 32000]

Subtype → 1_7_ext03 (*313*)

Subtype → 1_7_ext05 (*314*)

Extension (between 1_7 and 1_7_ext03)

Extension (between 1_7_ext03 and 1_7_ext05)

Subtype → /A_1 (*125*)

Subtype → /X (*124*)

Subtype* → /A_1a (*251*)

Subtype → /A_1b (*252*)

Subtype → /A_2 (*319*)

Subtype* → /A_2a (*320*)

Subtype → /A_2u (*321*)

Subtype → /A_2b (*322*)

Subtype*

Version → PDF_geospatial (*315*) role="family"

Version

Subtype → PDF/A (*318*) role="family"

From 1.7 Subtype* → /A_3 (*360*)

From PDF/A family Extension of

Subtype → /A_3a

Subtype → /A_3u

Subtype → /A_3b

# FDD

## format description document



**Microsoft Office Word 97-2003 Binary File Format (.doc)**

>> Back

**Table of Contents**
- Identification and description
- Local use
- Sustainability factors
- Quality and functionality factors
- File type signifiers
- Notes
- Format specifications
- Useful references

http://www.loc.gov/preservation/digital/formats/fdd/fdd000509.shtml

LIBRARY **LIBRARY OF CONGRESS**

- CFB header (usually 512 bytes):
  - Header Signature for the CFB format with 8-byte Hex value D0CF11E0A1B11AE1. Gary Kessler notes that the beginning of this string looks like "DOCFILE"
  - 16 bytes of zeroes
  - 2-byte Hex value 3E00 indicating CFB minor version 3E
  - 2-byte Hex value 0300 indicating CFB major version 3 or value 0400 indicating CFB major version 4. [Note: All DOC files created by compilers of this resource (in various versions of Word since 2003) and examined with a Hex dump utility have been based on CFB major version 3. Comments welcome.]
  - 2-byte Hex value FEFF indicating little-endian byte order for all integer values. This byte order applies to all CFB files.
  - 2-byte Hex value 0900 (indicating the sector size of 512 bytes used for major version 3) or 0C00 (indicating the sector size of 4096 bytes used for major version 4)
  - 480 bytes for remainder of the 512-byte header, which fills the first sector for a CFB of major version 3
  - Note: For a CFB of major version 4, the rest of the first sector would be 3,584 bytes of zeroes.
- Internal identifier for Word binary file (usually at byte offset 512 from beginning of DOC file):
  - 2-byte wIdent: Hex value ECA5
  - 2-byte version identifier: Hex value C100 [Note: The specification indicates that this is the value (equivalent to the integer 193) that *should* be used in this location, as *FibBase.nFib*, but indicates that some versions of Word had used other values. Hex value C000 has been used for new empty documents. Hex value C200 was used by the BiDi (bi-directional) build of Word 97.]
- Usually observed near end of file in documents created by recent versions of Microsoft Word:
  - More detailed version info, e.g., "Microsoft Word 97-2003 Document" or "Microsoft Word 97-2004 Document". See Note below on Identification of Microsoft Word version in CompObj stream.

# description

| Relationship to other formats | |
|---|---|
| Subtype of | CFB_3, Microsoft Compound File Binary File Format, Version 3. The compilers of this resource have experimented with saving Word documents as DOC files in several recent versions of Word. In all cases, the resulting file was in version 3 of CFB. Comments welcome. |
| Has later version | DOCX/OOXML_2012, DOCX Transitional (Office Open XML), ISO 29500:2008-2016, ECMA-376, Editions 1-5 |

## MS-DOC (509) is based on CFB 3 (380) but CFB 3 is the basis for MANY other formats – all the old Microsoft formats – not just MS-DOC

| Relationship to other formats | |
|---|---|
| Has subtype | MSG, Microsoft Outlook Item |
| Has subtype | MS-DOC, Microsoft Office Word 97-2003 Binary File Format (.doc) |
| Has subtype | MS-PPT, Microsoft Office Powerpoint 97-2003 Binary File Format (.ppt) |
| Has subtype | MS-XLS, Microsoft Office Excel 97-2003 Binary File Format (.xls, BIFF8) |
| Has later version | CFB_4, Microsoft Compound File Binary File Format, Version 4 |
| Affinity to | AAF_1_1, Advanced Authoring Format (AAF) Object, Version 1.1.<br><br>Early versions of the AAF format detailed use of the structured storage systems outlined in CFB to store the objects on disk. |

# 2 types of evaluation factors

**sustainability factors** for all formats

influence feasibility and cost of preserving content in the face of future change

**quality & functionality factors** that vary by content category

reflect considerations that will be expected by future users

# kate's favorite section: identifiers

**File type signifiers and format identifiers** ⓘ

| Tag | Value | Note |
|-----|-------|------|
| Filename extension | doc | Documented in the specification and elsewhere by Microsoft in many locations, for example at Office 2007 File Format MIME Types for HTTP Content Streaming. |
| Internet Media Type | application/msword | Documented by Microsoft at Office 2007 File Format MIME Types for HTTP Content Streaming. Also listed at IANA. See 1993 registration at https://www.iana.org/assignments/media-types/application/msword. Note that, unlike file formats for other proprietary Microsoft applications, the media type for the file with .doc as extension was assigned prior to establishment of the vendor (vnd) convention for media types. |
| Magic numbers | Hex: D0 CF 11 E0 A1 B1 1A E1 | Documented in the CFB specification, in 2.2 Compound File Header. Applies to all files in CFB format; see GCK'S File Signatures Table entry for Compound Binary File format (aka OLECF). |
| File signature | Hex: 3E 00 03 00 FE FF 09 00 | At byte offset 24 from beginning of file. Indicates CFB (Compound File Binary format) major version 3, minor version 3e. Assumes that all DOC files use this version of CFB. Comments welcome. |
| File signature | Hex: ECA5 | From specification. Indicates that this CFB file is a Word document. Usually at byte offset 512 from beginning of file. |
| Pronom PUID | fmt/40 | PRONOM has a number of entries for Microsoft Word format variants with the .doc extension. The PRONOM entry that corresponds to the scope of this format description is http://www.nationalarchives.gov.uk/PRONOM/fmt/40. |
| Wikidata Title ID | Q686498 | See https://www.wikidata.org/wiki/Q686498 for Word Binary File Format, all versions |
| Wikidata Title ID | Q28858035 | See https://www.wikidata.org/wiki/Q28858035 for Word Binary File Format, version nFib=0x00C1. Entry refers to [MS-DOC] as source reference and thus corresponds to the DOC format described here. |

**LIBRARY** LIBRARY OF CONGRESS

# hex view

# lots and lots of references



*   format specifications    **  useful references

Main | Table of Contents | Introduction | Textual Works | Still Image Works | Moving Image Works | Audio Works | Musical Scores | Datasets | GIS, Geospatial and Non-GIS Cartographic | Design and 3D | Software and Video Games | Web Archives

## Library of Congress Recommended Formats Statement - 2020-2021

Recommended Formats Statement identifies hierarchies of the physical and technical characteristics of creative formats, both analog and digital, which will best meet the needs of all concerned, maximizing the chances for survival and continued accessibility of creative content well into the future.

The 2020-2021 version includes significant changes from the 2019-2020 version. Specific changes are detailed in the Change Log and the Introduction.

# Recommended Formats Statement
https://www.loc.gov/preservation/resources/rfs/index.html

# content categories

textual works | still image works | moving image works | audio works | musical scores | datasets | GIS, geospatial and non-GIS cartographic | design and 3D | software and video games | web archives

orange = new or changed content categories for 2020

# evaluation criteria

## Global/Community Format Sustainability Factors

Disclosure
Adoption
Transparency
Self-documentation
External dependencies
Impact of patents
Technical protection
mechanisms

## LC Local/Institutional Factors

Staff experience and expertise

Software/Hardware/OS available

Representation/extent in LC collections/storage

Established workflow/functionality

# PREFERRED

**Global/community**: Meets or exceeds benchmarks for all relevant sustainability factors

**Local/institutional:** The Library of Congress has the skills, experience, workflows, tools and systems to manage and preserve these formats in current systems with confidence.

LIBRARY
LIBRARY
OF CONGRESS

# ACCEPTABLE

**Global/community:** Meets minimum acceptability across benchmarks or does not meet all relevant sustainability factors.

**Local/institutional:** The Library of Congress can manage this format at a basic level of acquisition, management and preservation; and a greater ability for management and preservation is within the Library's capacity with further investment.

## Global/Community Format Sustainability Factors

Each of these factors may have different emphasis or importance depending on the community of practice and content type. Some may not be applicable or essential for every format.

| Disclosure | Adoption | Transparency | Self-documentation | External dependencies | Impact of patents | Technical protection mechanisms | Notes | Community / Sustainability Summary |
|---|---|---|---|---|---|---|---|---|
| Is technical information about the format available through complete and open documentation and specifications? | Is the format widely used, especially in peer institutions? Is it integrated into multiple toolsets and not locked into specific vendor implementations? Are community user groups available for advice and support? | Can the format be analyzed with basic tools? Is standard character encoding supported? Is lossy compression or encryption enforced? | Is a file in this format able to describe its own content and structure with embedded metadata? If applicable, does this format have accessibility options for ADA compliance such as closed captioning? | Is this format free of dependence on particular hardware, operating system, or software for rendering or use? | Is the format free from patents with terms which might impede long term use? For example, when license terms include royalties based on use, costs could be high and unpredictable. | If this format *requires* the use of DRM, encryption or other protection mechanisms, is it possible for custodians to maintain future access to content reliably? | Any other mitigating factors to consider | **Preferred:** Meets or exceeds benchmarks for all relevant sustainability factors<br><br>**Acceptable:** Meets minimum acceptability across benchmarks or does not meet all relevant sustainability factors. |
| Yes | Yes | Maybe | Maybe. Few examples found in practice. | Yes | Yes | Yes | More common as distribution format than as master. | Preferred |
| Yes (new ISO spec available soon) | Yes | Maybe - varies with encoding/manufacturer | Yes | Maybe - depending on implementation - like Canon CR2. Consider in RFS adding qualifiers for tech specs | Yes | Yes | Release of upcoming ISO spec (2019 June). Need to update fdd188 | Acceptable |

**Green = PNG**
**Yellow = DNG**

| LC Local/Institutional Factors | | | | | | | Final Designation for RFS: Preferred or Acceptable |
| To the best of your knowledge, estimate the level of resources at LC available to preserve and manage the format. | | | | | | | |
| Staff experience and expertise | Software/Hardware/OS available | Representation/extent in LC collections/storage | Established workflow/functionality | Notes | Local Factors Summary | |
|---|---|---|---|---|---|---|
| Does LC have expertise with this format? For example, does LC staff: participate in standards efforts; format-related research and testing; have proficiency in tools and applications. | What are LC's functionality capabilities for this format? For example: Do staff have the software to analyze, describe, manage, and render this format? Are licenses available to all staff that need it? Is it approved by WCC process? Can it be installed/run on the networks/domains needed to process/view materials? | Does LC already have a meaningful number of files in this format in collections/managed storage? Basic check via Kibana: http://reportingvlp01.loc.gov/goto/171d80871e9c1adc1202b33a81770913 | Does LC's managed storage systems (such as CTS) have the resources to perform technical actions for this format such as format characterization, identification and validation; allow a "viewing copy"/QA by producing a thumbnail or the like? | Any other mitigating factors to consider? | Preferred: LC has the skills, experience, workflows, tools and systems to manage and preserve these formats in current systems with confidence. | Acceptable: LC can manage this format at a basic level of acquisition, management and preservation. Better management and preservation | |
| Maybe | Yes / Ample LC resources | Yes / Ample LC resources | Yes / Ample LC resources | | Preferred | **Preferred** |
| Yes / Ample LC resources | Maybe (site software not current - Photoshop) | Yes | Maybe - but not customized to LC specs so not using validation | Works well as an ingest format | Acceptable | **Acceptable** |

**Green** = PNG
**Yellow** = DNG

LIBRARY
LIBRARY OF CONGRESS

## ii. Photographs - Digital

| ii. Photographs - Digital | Preferred | Acceptable |
|---|---|---|
| **A. Faithful representation of the work** | › Equal in quality to the published version, best edition or master copy<br>› In the same format as the master copy | |
| **B. Technical Characteristics** | › Highest resolution available, not rescaled or interpolated<br>› Highest bit depth available, 16 bits per channel if available<br>› Embedded color profile or specified color space used in published version<br>› Uncompressed<br>› Unlayered | › Lossless compression or lower compression ratios<br>› Discrete wavelet transform (DWT) preferred to discrete cosine transform (DCT)<br>› Layered, if supported by preferred or acceptable format |
| **C. Formats** | › TIFF (*.tif)<br>› JPEG2000 (*.jp2)<br>› PNG (*.png)<br>› JPEG/JFIF (*.jpg) | › Photoshop (*.psd)<br>› JPEG2000 Part 2 (*.jpf, *.jpx)<br>› Digital Negative DNG (*.dng)<br>› Proprietary Camera Raw formats (*.nef, *.crw)<br>› GIF (*.gif) |

LIBRARY
LIBRARY OF CONGRESS

| D. Metadata | 1. As supported by format:<br><br>  a. Title<br>  b. Creator<br>  c. Creation Date<br>  d. Place of publication<br>  e. Publisher/producer/distributor<br>  f. Contact information<br><br>2. Include if available:<br><br>  a. Common embedded schema (e.g., IPTC)<br>  b. Language of work<br>  c. Other relevant identifiers (e.g., DOI, LCCN, etc.)<br>  d. Subject descriptors<br>  e. Abstracts<br>  f. Key or reference to each data field and technical production information (e.g. EXIF metadata from digital camera | Metadata provided separately in external text of XML-based file |
| E. Technological Measures | Files must contain no measures (such as digital rights management technologies or encryption) that control access to or prevent use of the digital work. | |

https://www.loc.gov/preservation/resources/rfs/stillimg.html

# acknowledgements

Sustainability of Digital Formats

https://www.loc.gov/preservation/digital/formats/index.html

- Caroline Arms
- Marcus Nappier
- Laurel Gassie

Recommended Formats Statement

https://www.loc.gov/preservation/resources/rfs/index.html

- Ted Westervelt
- Jesse Johnston
- Marcus Nappier
- All content team leaders/members

# thank you

**Kate Murray**
[kmur@loc.gov](mailto:kmur@loc.gov) | @fileformatology

**LIBRARY** LIBRARY OF CONGRESS