

**CAARA**

Council of Australasian Archives  
and Records Authorities

# **Sustainable Digital File Formats for creating and using records**

---

**Consultation Draft**

Australasian  
Digital  
Recordkeeping  
Initiative

adri

**November 2019**

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	The problem	3
1.2	The proposed solution	3
1.3	This document	4
1.4	What do we want you to do	4
1.5	Who are we?	4
<b>2</b>	<b>Recommended formats (summary)</b>	<b>5</b>
<b>3</b>	<b>Recommended formats (detail)</b>	<b>6</b>
3.1	Plain Text Files	6
3.2	Data	6
3.3	Office Documents (Word Processing)	7
3.4	Office Documents (Spreadsheet)	9
3.5	Office Documents (Presentation)	10
3.6	Project Planning	10
3.7	Email	11
3.8	Website	11
3.9	Website (archival)	13
3.10	Audio	14
3.11	Image (Raster)	15
3.12	Image (Vector)	16
3.13	Image (Scanned Documents)	17
3.14	Image (Graphic Metafiles)	17
3.15	Video (Bare)	17
3.16	Video (Container)	18
3.17	Electronic Publications	19
3.18	Encapsulation	19
3.19	CAD	20
3.20	Geospatial	20
<b>4</b>	<b>Long term preservation format criteria</b>	<b>24</b>
<b>5</b>	<b>Commentary</b>	<b>25</b>
5.1	General	25
5.2	Details of specific criteria	25

# 1 Introduction

1 What this document is, who prepared it, and what we would like you to do

---

## 2 1.1 The problem

3 Will you be able to view a JPEG image in 10 years? In 100 years? In 1000 years?

4 Much of the digital information that is generated by the public service will have a long life. Some of the information  
5 will need to be accessible for decades. A small amount of the information will need to be kept permanently. This is just  
6 the reality of information; it is as true of information expressed in a paper form as information expressed in a digital  
7 form.

8 The particular challenge of digital information is that it is necessary to use software to access the information held in a  
9 digital object. It is necessary to use an image viewer to display a JPEG image, for example. Accessing digital information  
10 means having software to open the digital object for the entire lifespan of the information. For example, it may not be  
11 possible to find software that displays JPEG images in 100 years if the JPEG format has fallen out of use. If software to  
12 access a format is not available, information stored in this format will become inaccessible.

13 Records created by government agencies have a minimum mandated lifespan. This lifespan is set by Disposal  
14 Authorities issued by Government archival authorities. This minimum lifespan is not based on technology, instead it is  
15 based on an analysis on how long the government is likely to need the information. Typical minimum lifespans can  
16 range from 10 to 50 years. Information about people may be required to be kept for up to 100 years. Information  
17 about infrastructure usually must be kept for the life of the asset – major assets can have a life of 100 years or more.  
18 Some information is considered of permanent value and must be kept forever.

19 It is the agency's responsibility to maintain access to records for the minimum authorised retention period,  
20 irrespective of changing technology (unless the information has been transferred to an archive authority). If a format  
21 becomes obsolete or falls out of use, it is the agency's responsibility (and expense) to find a solution to the problem of  
22 ensuring access to the information.

23 This document is a draft of advice on how agencies can retain access to digital information for its expected lifespan.

## 24 1.2 The proposed solution

25 No-one knows what the future will hold; it is consequently impossible to eliminate the risk of a format falling out of  
26 use.

27 What is possible, however, is to reduce the risk by creating information in formats that are likely to remain accessible  
28 for long periods. This still leaves the challenge of selecting these 'low risk' formats.

29 This document presents list of formats that we believe have a low risk of becoming inaccessible over time.

30 We recommend that all information generated by an agency be created in a 'low risk' format, irrespective of how long  
31 it is expected that the information needs to be retained. This is for two reasons:

- 32 • It is a far simpler to implement in an agency, and to communicate to agency staff; i.e. just use these formats.
- 33 • The minimum lifespan of information can change over time.

34 We also recommend that agencies keep low risk file formats in mind for those records which they receive from other  
35 individuals or organisations. Where it is possible and practical, we encourage agencies to promote the use of low risk  
36 file formats for official business both internally and externally.

37 *Please note that the goal of this document is different to that of many similar lists of formats produced by preservation*  
38 *institutions. Other lists are of formats accepted by the preservation institution for preservation. This list, however, is of*  
39 *recommended creation formats. If the staff of your agency are creating particular types of information (e.g. still*

1 *images), they should be encouraged to create the information in one of the recommended formats. We believe that the*  
2 *use of the recommended formats will minimise the risk of losing access to the information due to format obsolescence.*  
3 *As the formats are in widespread use, restricting staff to these formats is unlikely to be restrictive. Indeed, most staff*  
4 *would be unlikely to notice any restriction.*

### 5 **1.3 This document**

6 This document is a draft of advice to agencies on preferred low risk formats for records being created or received. It  
7 contains two parts:

- 8 • A list of formats that we consider have low risk of access failure over a reasonable time, presented in summary and  
9 detailed versions. These are arranged by the type of information.
- 10 • A list of criteria we have used to select the formats. This list can be used to select a format if the type of  
11 information is not covered in the first list.

### 12 **1.4 What would we like you to do?**

13 We would like feedback on the two lists – the recommended file formats list and the criteria list.

### 14 **1.5 How do you comment?**

15 Please provide feedback to Andrew Waugh (Andrew.Waugh@prov.vic.gov.au) by 18 December 2019.

### 16 **1.6 Who are we?**

17 We are a sub-committee of the Australasian Digital Recordkeeping Initiative (ADRI). ADRI is composed of  
18 representatives of the State, Territory and National Archives of Australia and New Zealand. ADRI is a key work program  
19 of the Council of Australasian Archives and Records Authorities (CAARA).

## 2 Recommended formats (summary)

1 These formats have been identified as having low risk of access failure (e.g.  
2 through obsolescence) in the foreseeable future. They are a good choice for  
3 agency staff when creating and maintaining information.

4 When creating content, you should choose to save it in one of the following low risk formats. We have listed the  
5 formats by the type of information, and ordered the list so that common types of formats appear earlier in the list.

- |    |   |    |   |
|----|---|----|---|
| 6  | • <b>Word processing documents:</b> Microsoft Word          | 26 | • <b>Electronic publications:</b> PDF (.pdf) or EPUB            |
| 7  | (.docx, .doc), Open Office Document (.odt) or PDF           | 27 | (.epub)   |
| 8  | (.pdf)  | 28 | • <b>Encapsulation:</b> ZIP (.zip), GZIP (.gzip), or TAR (.tar) |
| 9  | • <b>Spreadsheets:</b> Microsoft Excel (.xlsx, .xls), Open  | 29 | • <b>CAD:</b> DXF (.dxf), .DWG (.dwg), STEP (.stp, .step, or    |
| 10 | Office Spreadsheet (.ods), PDF (.pdf), CSV (.csv), or       | 30 | .p21), PDF/E (.pdf)   |
| 11 | TSV (.tsv)  | 31 | • <b>Geospatial:</b> ESRI Shapefiles, Geopackage, GEOjson       |
| 12 | • <b>Presentations:</b> Microsoft PowerPoint (.pptx, .ppt), | 32 | (.json), DEM, GML (.gml), KML (.kml or .kmz),                   |
| 13 | Open Document Presentation (.odp), or PDF (.pdf)            | 33 | GEOtiff (.tiff), BigTiff (.tiff), ECW (.ecw), JPEG2000          |
| 14 | • <b>Raster images:</b> TIFF (.tiff or .tif), JPEG (.jpg or | 34 | (.jp2), SVG (.svg), LAS (.las) or IMG (.img)                    |
| 15 | .jpeg), JPEG2000 (.jp2), PNG (.png) or DNG (.dng)           | 35 | • <b>Email:</b> Mime encoding (.eml), UNIX mailbox (.mbx        |
| 16 | • <b>Vector images:</b> SVG (.svg) or Open Document         | 36 | or .mbox), or Microsoft Outlook (.msp or .pst)                  |
| 17 | drawing (.odg)  | 37 | • <b>Websites:</b> HTML (.htm or .html) or XML (.xml)           |
| 18 | • <b>Multipage scanned documents:</b> PDF (.pdf)            | 38 | together with supporting files (.css, .xsd, and .dtd)           |
| 19 | • <b>Audio:</b> WAV (.wav, .bwav, or .bwa), MP3 (.mp3),     | 39 | • <b>Archived Websites:</b> Web archive (.warc) or              |
| 20 | MP4 (.mp4 or .m4a), or FLAC (.flac)                         | 40 | Internet archive (.arc)   |
| 21 | • <b>Video (in a container):</b> Ogg (.ogg or .ogv), DCP    | 41 | • <b>Image (graphic metafile):</b> CGM (.cgm)                   |
| 22 | (.dcp), MPEG2 (.mpg, .mpeg), or MPEG4 (.m4v,                | 42 | • <b>Data files:</b> CSV (.csv), TSV (.tsv), XML (.xml), JSON   |
| 23 | .m4a, f4v, .or f4a)   | 43 | (.json or .jsn), or SIARD                                       |
| 24 | • <b>Video (bare):</b> Motion JPEG2000 (.mjp, .mj2 or       | 44 | • <b>Text files:</b> Plain text (.txt) files                    |
| 25 | .m2v), or DPX (.dpx)  |    |   |

45

# 3 Recommended formats (detail)

1 This section provides brief detail and our rationale for recommending each  
2 format.

## 3 3.1 Plain Text Files

4 Plain text files contain simple text without any markup or formatting.

5 Where the text files contain data (e.g. CSV files), see section 2.2 (Data). Where the text files contain website  
6 information, see Section 2.8 (Websites).

Recommended	Notes	Rationale
TXT	<p>A ".txt" file generally contains only plain text (without formatting), intended for humans to read.</p> <p>The Unicode character set, and its encoding as UTF-8 or UTF-16 is the standard and widespread method of representing textual characters. The UTF-8 encoding was designed to be backwards compatible so that the older ASCII encoding is also valid UTF-8. UTF-16 caters for the additional characters needed for non-English text.</p>	<p>Plain text is public, standardized, and universally readable. The use of plain text provides independence from programs that require their own special encoding or formatting or file format. Plain text files can be opened, read, and edited with countless text editors and utilities.</p>

7

## 8 3.2 Data

9 Data files contain data such as database tables, scientific observations, and metadata.

Recommended	Notes	Rationale
SIARD <sup>1</sup>	<p>Software Independent Archiving of Relational Databases</p> <p>Use for capturing information from relational databases</p>	<p>SIARD is an XML representation of the data in a database and its schema. As a textual representation, it is expected that the information will be accessible, even if the SIARD software becomes obsolete. SIARD is widely used by Archives in Europe to archive databases.</p>

<sup>1</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000426.shtml>

Recommended	Notes	Rationale
CSV/TSV	Comma separated variables, tab separated variables  It is not possible to represent data coding in a CSV/TSV file. This information needs to be documented separately.	CSV/TSV are very widely used textual representation of simple data structures. The mechanism has been used for a very long time – probably over 50 years.
XML <sup>2</sup>	The XML standard describes how data is encoded into text. It does not describe the meaning of the data. This information needs to be documented separately.	XML is an open international standard mechanism for marking up data in textual format. It has been very widely adopted as the basis for interchange of documents and data over the Web. Many generic tools exist, including free and open source software. Major software vendors have all incorporated support for XML in some form.
JSON <sup>3</sup>	JavaScript Object Notation (JSON) is a lightweight, text-based, language-independent data interchange format. The openly documented JSON standard describes how data is encoded into text. It does not describe the meaning of the data. This information needs to be documented separately.	As a simple textual encoding of data, JSON is expected to be easily processible indefinitely. JSON is considered to be a lightweight alternative to XML. JSON is so simple that support for reading or writing it is integrated into almost every system or programming language used for applications on the Web.

1

## 2 3.3 Office Documents (Word Processing)

3 Use these formats for documents produced by word processing programs (e.g. Word).

4 **WARNING:** Office documents may have a range of other format types (e.g. images) embedded within the document.

5 When producing documents, this embedded content should be in an appropriate long term preservation format (e.g. JPEG for photographs).

6

Recommended	Notes	Rationale
DOCX <sup>4</sup>	Microsoft Word Document (XML).  Although this is a format used by a specific proprietary software product, it is not necessary to use this product to create DOCX documents.	This format is an international standard, and is based on XML. Further, Microsoft Word is the clear market leader in office document preparation and very large numbers of documents in this format are created.
DOC	Microsoft Word Document  Although this is a format used by a specific proprietary software product, it is not necessary to use this product to create DOC documents.	This format is an Ecma standard (an industry-based standards organisation) <sup>5</sup> , and is consequently documented. Its main basis for inclusion, however, is simply the extraordinarily large number of documents created in this format. It is consequently unlikely that it would become inaccessible in the foreseeable future.

<sup>2</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000264.shtml>

<sup>3</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000381.shtml>

<sup>4</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000397.shtml>

<sup>5</sup> <http://www.ecma-international.org/>

Recommended	Notes	Rationale
ODT <sup>6</sup>	<p>OpenDocument Text Document Format. The OpenDocument format is a generic format for office documents, such as text, spreadsheets, presentations, drawings, and databases.</p>	<p>The format is an open standard, based on XML, and is supported by a variety of office software packages.</p>
PDF	<p>The Portable Document Format (PDF) is a file format developed by Adobe to deliver final ‘published’ documents that are primarily intended to be read by end users. It is difficult to extract data from PDF files and to modify them.</p> <p>These documents may be structured or simple. They can include text, images, graphics, and other multimedia content, such as video and audio, and are independent of application software, hardware, and operating systems.</p> <p>For recordkeeping purposes, we recommend using variants of PDF-A<sup>7</sup>, but PDF 1.7<sup>8</sup> is acceptable.</p> <p>This format offers several forms of technical protection, including encryption, which could prevent custodians of digital content ensuring accessibility in future technological environments. These technical mechanisms should not be used.</p>	<p>Fully documented. Most members of the PDF family were developed by Adobe Systems Incorporated, which makes the specifications available openly and at no charge. Several members of the family have been adopted as ISO international standards, e.g. PDF/X, PDF/A, and PDF version 1.7. Extremely widely adopted for disseminating page-oriented documents.</p>

1

<sup>6</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000428.shtml>

<sup>7</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000318.shtml>

<sup>8</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000277.shtml>



## 1 3.4 Office Documents (Spreadsheet)

2 Use these formats for spreadsheets (e.g. Excel).

Recommended	Notes	Rationale
XLSX <sup>9</sup>	Microsoft Excel Spreadsheet (XML).  Although this is a format used by a specific proprietary software product, it is not necessary to use this product to create XLSX spreadsheets.	This format is an international standard, and is based on XML. Further, Microsoft Excel is the clear market leader in spreadsheet preparation and very large numbers of spreadsheets in this format are created.
XLS	Microsoft Excel Spreadsheet.  Although this is a format used by a specific proprietary software product, it is not necessary to use this product to create XLS spreadsheets.	This format is an Ecma standard, and is consequently documented. Its main basis for inclusion, however, is simply the extraordinarily large number of spreadsheets created in this format. It is consequently unlikely that it would become inaccessible in the foreseeable future.
ODS <sup>10</sup>	OpenDocument Spreadsheet Document Format. The OpenDocument format is a generic format for office documents, such as text, spreadsheets, presentations, drawings, and databases.	The format is an open standard, based on XML, and is supported by a variety of office software packages.
PDF	See PDF in Section 2.3.  Spreadsheets saved to this format have restricted functionality – it is difficult to extract data from them. The underlying formulae that are used to calculate the information are not captured. This format is most appropriate for final ‘published’ spreadsheets that are exclusively intended to be read by end users.	
CSV <sup>11</sup>	See CSV/TSV in Section 2.3.  Note that spreadsheets in this format will have restricted functionality – the format only represents the data. It cannot store the underlying formulae used to calculate the data, nor can it represent the formatting of particular cells. It should only be used when only the data in the spreadsheet is of value.	As it is simple text, this format is extremely easy to access and process.

3

<sup>9</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000398.shtml>

<sup>10</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000439.shtml>

<sup>11</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000323.shtml>

## 3.5 Office Documents (Presentation)

Use these formats for presentations (e.g. PowerPoints).

WARNING: Office documents may have a range of other format types (e.g. images) embedded within the document. When producing documents, this embedded content should be in an appropriate long term preservation format (e.g. JPEG for images).

Recommended	Notes	Rationale
PPTX <sup>12</sup>	Microsoft PowerPoint (XML). Although this is a format used by a specific proprietary software product, it is not necessary to use this product to create PPTX presentations.	This format is an international standard, and is based on XML. Further, Microsoft PowerPoint is the clear market leader in simple presentation preparation and very large numbers of presentations in this format are created.
PPT	Microsoft PowerPoint Although this is a format used by a specific proprietary software product, it is not necessary to use this product to create PPT presentations.	This format is an Ecma standard, and is consequently documented. Its main basis for inclusion, however, is simply the extraordinarily large number of presentations created in this format. It is consequently unlikely that it would become inaccessible in the foreseeable future.
ODP <sup>13</sup>	OpenDocument Presentation Document Format. The OpenDocument format is a generic format for office documents, such as text, spreadsheets, presentations, drawings, and databases.	The format is an open standard, based on XML, and is supported by a variety of office software packages.
PDF	See PDF in Section 2.3.  Presentations in this format will have restricted functionality – it is difficult to extract information from them and to modify them. It is most appropriate for final ‘published’ presentations that are primarily intended to be read by end users.	

## 3.6 Project Planning

There are no recommended formats to represent project planning data.

This is because we know of no project planning formats sufficiently widely used, or sufficiently well documented, to make the format likely to survive for the required periods.

<sup>12</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000399.shtml>

<sup>13</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000435.shtml>

## 1 3.7 Email

2 These formats are used to capture sent or received emails. Some recommended formats can either represent a single  
3 email, while others capture an entire mailbox including the contained emails.

4 NOTE. Emails usually have attachments. The formats contained in this section will only ensure that the basic text in the  
5 email body and the email headers are accessible, and that the attachments can be extracted as they were originally  
6 sent. They will not ensure that the information in an attachment is accessible. Attachments can be of any format. In  
7 order that they remain accessible, it is important that each attachment be an appropriate long term preservation  
8 format.

Recommended	Notes	Rationale
EML <sup>14</sup>	MIME encoded emails  Each instance of this format represents a single email.	This is a textual representation that is identical to the standard representation of an email while it is being transferred between computers. It is consequently extraordinarily widely used and well documented.
MBX, MBOX <sup>15</sup>	UNIX style mailbox format  Each instance of this format represents a collection of emails, the mailbox (or part of a mailbox) of a user.	This is a textual representation that is an extension to the standard representation of an email while it is being transferred between computers. It is widely used.
MSG <sup>16</sup>	The Microsoft Outlook Item MSG file format is a syntax for storing a single Message object, such as an email, an appointment, a contact, a task, and so on, in a file. A MSG file may also include email attachments.	A proprietary file format developed by Microsoft but no patents are claimed. Full documentation is available. This format is widely used.
PST <sup>17</sup>	Microsoft mailbox format (generic)  Each instance of this format represents a collection of emails, the mailbox (or part of a mailbox) of a user.	This format is widely used.

## 9 3.8 Website

10 The formats in this section represent a captured website. These represent the files that are transmitted to a web  
11 browser to display a web page.

12 Note that while several formats are presented here, many of these formats capture only an aspect of the website (e.g.  
13 CSS files capture style information that is applied to other HTML or XML files). Therefore, a website capture typically  
14 contains a range of these file types. It is not sufficient to preserve only the content (HTML or XML) files if the intent is  
15 to capture the core aspects of a website.

16 Further, a website typically delivers various types of content (e.g. images, audio, etc.). This section is only concerned  
17 with the formats which structure and display web information. Other delivered content may be covered by formats  
18 mentioned elsewhere in this chapter.

<sup>14</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000388.shtml>

<sup>15</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000383.shtml>

<sup>16</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000379.shtml>

<sup>17</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000378.shtml>

1 Finally, note that these web formats only relate to ‘static’ web pages. Many web pages also include dynamic content (a  
 2 trend which is increasing). This could include code that is executed in the web browser (e.g. client side JavaScript) to  
 3 achieve certain effects. Or it could include API calls to servers to retrieve content as the user performs actions (e.g.  
 4 Google map embedded in the page). A preserved website which depends partly upon client-side executed scripts will  
 5 likely vary in completeness and sometimes pull content from the live web (“leakage”) which is likely to be  
 6 chronologically inconsistent<sup>18</sup>. It is essentially impossible to fully capture these types of web pages in a form that is  
 7 suitable for long term preservation but the results might still be acceptable, depending on the extent and nature of the  
 8 client-side scripting used.

Recommended	Notes	Rationale
HTM, HTML <sup>19</sup>	<p>Any variant acceptable, with a preference to the later versions.</p> <p>Different web browsers will display HTML pages differently. Later versions of HTML are more consistent in their display.</p> <p>Many HTML pages have an associated CSS resource that controls how the page is displayed. It is consequently necessary to capture both the HTML page and the CSS resource to properly capture the page.</p>	<p>HTML is a text based standard. The various versions of HTML are standards which are widely available. HTML is ubiquitous on the internet.</p>
XML <sup>20</sup>	<p>Any variant acceptable, with a preference to the later versions.</p> <p>XML pages have an associated CSS resource that controls how the page is displayed. It is consequently necessary to capture both the XML page and the CSS resource to properly capture the page.</p>	<p>XML is an international standard that is extremely widely used in a range of applications.</p>
CSS <sup>21</sup>	<p>Used to represent the styles of elements in XML and HTML. There are a large range of versions of the standard, with more features being added continually.</p>	<p>A text based international standard, it is very widely used to style web pages.</p>
XSD <sup>22</sup>	<p>XML Schema Definition. These files are used to control the structure of XML documents. They are not necessary to display or use XML documents, but some tools will require them.</p>	<p>This is a text based international standard.</p>
DTD <sup>23</sup>	<p>XML Document Type Definition. These files are used to control the structure of XML documents. They are not necessary to display or use XML documents, but some tools will require them.</p>	<p>DTD documents are part of the XML international standard.</p>

9

<sup>18</sup> Brunelle, J. F., Kelly, M., Weigle, M. C., & Nelson, M. L. (2016). The impact of JavaScript on archivability. *International Journal on Digital Libraries*, 17(2), 95-117. doi:<http://dx.doi.org/10.1007/s00799-015-0140-8>

<sup>19</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000475.shtml>

<sup>20</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000075.shtml>

<sup>21</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000482.shtml>

<sup>22</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000077.shtml>

<sup>23</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000076.shtml>

## 1 3.9 Website (archival)

2 The previous section described the types of files that must be captured in order to preserve a web page (or collection  
 3 of webpages. This section describes formats that are used to encapsulate the files representing a webpage or website  
 4 into a single archival object. They are equivalent to the Encapsulation formats described below, but are targeted  
 5 specifically at encapsulating web sites.

Recommended	Notes	Rationale
WARC <sup>24</sup>	Web ARChive file format. This format was developed from the ARC format.	This format is commonly used internationally by archives and libraries to perform large scale harvest and capture of web sites. It is consequently well supported by tools. It is a simple encapsulation format and would be easy to migrate if necessary.
ARC <sup>25</sup>	Internet Archive ARC file format  This format is acceptable but the WARC format is preferred.	Documentation and tools to use files in the format freely available.

6

<sup>24</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000236.shtml>

<sup>25</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000235.shtml>

## 1 3.10 Audio

2 These formats represent sound (only), such as recorded music, radio broadcasts, telephone conversations, and  
3 podcasts.

4 High quality audio files are extremely large. For this reason, most audio formats discard information, either by using  
5 compression, reducing the sampling rate, or the sample quality. A great deal of skill has been used to produce formats  
6 (and compression techniques) that minimise the size of audio files while retaining sufficient audio quality for the  
7 intended purpose. In most cases, the creator of a recording can select the quality of the recording. In selecting an  
8 audio format, it is consequently important to consider the type of audio material being stored to ensure that the  
9 quality is sufficient. Telephone conversations, for example, can be stored with far lower quality than the master  
10 recording of a performance of a piece of music.

Recommended	Notes	Rationale
WAV <sup>26</sup> , BWAV, BWF	WAV encoding is normally used when it is desired to encode a high quality representation of audio information. WAV and its Broadcast WAVE variants are container formats, which contain an audio bitstream that is encoded using LPCM (as used on audio CDs) or an MPEG audio format. For format planning purposes, it is useful to know which encoding is to be used.	Proprietary format, with full documentation freely available. The WAV format and its variants, and the audio encodings mentioned are extremely widely used in the sound industry for containing high quality audio information.
MP3 <sup>27</sup>	MPEG-1 or -2 audio layer III. MP3 is the standard audio format for consumer consumption of audio. It has been carefully constructed to give an acceptable quality of audio with a small file size (high compression).  It is expected that MP3 audio files would give acceptable quality for most business purposes. Where high quality recordings are necessary, a WAV file is recommended.	MP3 audio files are an international standard. Beyond that, however, they extraordinarily widely used to distribute audio over the internet. It is extremely unlikely that they will become inaccessible for the foreseeable future.
MP4, M4A <sup>28</sup>	MPEG-4 Audio (AAC encoding). This audio format is intended to replace MP3, giving better quality for the same level of compression.  Compared with MP3, this format is not as widely used, but AAC has seen considerable adoption by industry <sup>29</sup> .	These formats are part of the international standard for representing video information. As such they are unlikely to become inaccessible, even though they are not as widely used for audio.
FLAC (Free Lossless Audio Codec) <sup>30</sup>	Open source bit stream encoding format designed for lossless compression of LPCM audio data, with many of its default parameters tuned to CD-quality music data. Generally used for final-state, end-user delivery.	FLAC files are much smaller than WAV (about a third the size of the uncompressed audio). FLAC is an open format well supported by free software.

11

<sup>26</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000001.shtml>

<sup>27</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000105.shtml>

<sup>28</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000155.shtml>

<sup>29</sup> <https://mpeg.chiariglione.org/standards/mpeg-4/audio>

<sup>30</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000198.shtml>

## 1 3.11 Image (Raster)

2 Raster images are still images typically produced by cameras or document scanners. Raster images are composed of a  
3 grid of dots.

4 High quality images are extremely large. For this reason, most image production processes discard information to  
5 produce an 'acceptable' image at a reasonable size. Information may be discarded by using compression (i.e. lossy  
6 compression), reducing the resolution (down sampling) or reducing the colour width (going from 48 bit colour to 24  
7 bit, or from greyscale to bitonal).

8 The formats in this section should be used where the file contains a single image. For multi-page scanned images, refer  
9 to section 2.13 (Image - Scanned Documents).

Recommended	Notes	Rationale
TIFF <sup>31</sup>	TIFF 6.0 is a preferred format for long term preservation use. The image in a TIFF file can be either uncompressed, or compressed using a variety of compression formats. LZW compression is a widely used lossless technique. Note that uncompressed images can be extremely large.	A widely supported format by many applications. Fully documented specification. Proprietary but patents not enforced for wrapper format. Patents for the common LZW compression have long expired.
JPEG, JPG <sup>32</sup>	JPEG images are usually expressed as a JFIF file, but they can be a raw JPEG image. Either approach is acceptable.	JPEG is an international standard, and it is extremely widely used – it is the default capture format for digital cameras (including phones). There would be billions of JPEG images. It is unlikely that JPEG would be unreadable at any time in the foreseeable future.
JPEG2000, JP2 <sup>33</sup>	JPEG2000 was intended to replace JPEG, but JPEG is simply too well entrenched in the market. JPEG2000 achieves higher rates of compression for similar quality, and has a lossless compression option.	JPEG2000 is an international standard, though is not very widely adopted.
PNG <sup>34</sup>	Portable Network Graphics. This format was also intended to replace JPEG, with similar benefits to JPEG2000. Like JPEG2000 it has been unable to make significant inroads into the market, though it is becoming more widely used in the web.	PNG is an international standard, though is not very widely adopted.

<sup>31</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000022.shtml>

<sup>32</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000017.shtml>

<sup>33</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000138.shtml>

<sup>34</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000153.shtml>

Recommended	Notes	Rationale
DNG <sup>35</sup>	The Digital Negative (DNG) Camera Raw Format can be used for storing camera raw files, which are data captured from the camera sensor and requiring further processing to generate usable images. The advantage of raw files is the increased artistic control over the resulting image, as opposed to JPEG and TIFF files. Storing the DNG file with an embedded JPEG preview image may help with digital asset management. DNG is an extension of TIFF 6.0 and is compatible with TIFF/EP (Digital Photography).	A fully documented non-proprietary format, created as a standard by Adobe Systems. The DNG format can be used by a wide variety of hardware and software applications which generate, process, manage or archive camera raw files.

1

## 2 3.12 Image (Vector)

3 Vector images are still images composed of lines and uniformly coloured areas. They are easily identifiable as the  
4 quality of the image does not degrade as the image is zoomed in. Typically vector images are produced by drawing  
5 programs such as Adobe Illustrator.

6 It should be noted that the most widely used Image (Vector) formats are the proprietary formats AI (Adobe Illustrator)  
7 and CDR (Corel Draw).

8 CAD files and Geospatial files are also typically vector images, but they have additional capabilities. For more  
9 information on these formats see the relevant sections (2.19 CAD; 2.20 Geospatial).

Recommended	Notes	Rationale
SVG <sup>36</sup>	Scalable Vector Graphics (SVG) is a language for describing two-dimensional graphics in XML. SVG allows for three types of graphic objects: vector graphic shapes (e.g., paths consisting of straight lines and curves), images and text.	This is an open standard from the W3C using XML. (SVG) is supported by all modern browsers for desktops and mobiles. Some features, such as SMIL animation and SVG Fonts are not as widely supported. There are many SVG authoring tools, and export to SVG is supported by all major vector graphics authoring tools <sup>37</sup> .
ODG <sup>38</sup>	Open Document graphics format. SVG is preferred due to its greater support and more powerful features.	This format is an international standard, and is based on XML.

10

<sup>35</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000188.shtml>

<sup>36</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000020.shtml>

<sup>37</sup> <https://www.w3.org/Graphics/SVG/>

<sup>38</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000436.shtml>

### 1 3.13 Image (Scanned Documents)

2 The formats in this section should be used when a multipage document (e.g. a book or report) has been scanned.

3 For simple single images, please use the formats in section 2.11 (Image - Raster). For geo-referenced images, please  
4 use section 2.20 (Geospatial).

5 Note that while JPG and TIFF are not typically recommended scanning formats for documents (neither being widely  
6 used for holding multipage documents), some large document and image repositories do use multi-page TIFF.

Recommended	Notes	Rationale
PDF	See PDF in Section 2.3.	

7

### 8 3.14 Image (Graphic Metafiles)

9 Graphic metafiles are file formats that can contain a mixture of different types of images (e.g. both raster and vector  
10 images)

Recommended	Notes	Rationale
CGM <sup>39</sup>	Computer Graphics Metafile	An open, platform-independent format for the exchange of raster and vector data for technical applications.

### 11 3.15 Video (Bare)

12 These formats are used for moving images (which may or may not include sound). Video formats are divided into two  
13 classes: bare formats which contain only the video and sound information; and container formats that allow multiple  
14 different types of object to be included in one file. This section is about bare video formats.

15 High quality video streams are extraordinarily large. For this reason, actual video files almost invariably have the  
16 quality reduced to allow the video stream to be a manageable size. Normally, the video is compressed using a lossy  
17 compression algorithm, the image size (resolution) is strictly limited, and the frame rate is adjusted to the minimum  
18 acceptable. It is for this reason that decisions about the quality of the video should be based on business needs and  
19 industry standards, rather than archival decisions. In general however, agencies should avoid complete dependence on  
20 proprietary formats/compression techniques that are natively produced by capture devices; standard formats and  
21 compression techniques are preferred.

22

Recommended	Notes	Rationale
MJP, MJ2, M2V <sup>40</sup>	Motion JPEG 2000. The extension MJP is preferred	Open international standard.
DPX <sup>41</sup>	Digital Moving Picture Exchange (SMTPE). Used for high quality representation intended for	Reasonably widely adopted and fully documented

<sup>39</sup> <https://www.nationalarchives.gov.uk/documents/graphic-file-formats.pdf>

<sup>40</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000127.shtml>

<sup>41</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000178.shtml>

Recommended	Notes	Rationale
	theatrical distribution. The standard does not control how individual images are represented.	open standard.

1

## 2 3.16 Video (Container)

3 These formats are used for moving images (which may or may not include sound). Video formats are divided into two  
 4 classes: bare formats which contain only the video and sound information; and container formats that allow multiple  
 5 different types of object to be included in one file. These could include multiple audio files (for different languages),  
 6 closed captions, etc. This section is about container video formats.

7 It is important to note that with container formats, the actual video stream may be encoded in many different ways.  
 8 This needs to be considered.

9 High quality video streams are extraordinarily large. For this reason, actual video files almost invariably have the  
 10 quality reduced to allow the video stream to be a manageable size. Normally, the video is compressed using a lossy  
 11 compression algorithm, the image size (resolution) is strictly limited, and the frame rate is adjusted to the minimum  
 12 acceptable. It is for this reason that decisions about the quality of the video should be based on business needs and  
 13 industry standards, rather than archival decisions. In general however, agencies should avoid complete dependence on  
 14 proprietary formats/compression techniques that are natively produced by capture devices; standard formats and  
 15 compression techniques are preferred.

Recommended	Notes	Rationale
AVI <sup>42</sup>	Audio Video Interleave. AVI containers can use a variety of video formats (WAVE, MP3 Audio, DivX)	Widely adopted and fully documented proprietary format.
OGG, OGV <sup>43</sup>	Ogg container with video	Fully documented. Developed as an open source and patent-free project.
DCP <sup>44</sup>	Digital Cinema Package. Format is widely used in final-state for use in a cinema distribution chain; may also serve as a middle-state format for archiving. Uses DCDM to actually encode the video stream. This, in turn, uses MXF_UNC to actually encode the video stream	Fully disclosed proprietary format. Documentation available at Digital Cinema Initiatives Web site. No licensing or patents identified.
MPG, MPEG	MPEG2 video (i.e. as used on DVDs). Lossy compression. Many software tools exist for encoding and decoding.	Open standard. Widely adopted for filmmaking, DVD disks, and other applications. Most significant is the format's required use in digital terrestrial TV broadcasting.
M4V, M4A, F4V, F4A	MPEG4 video. Huge range of options. Generally a final-state (end-user delivery) format. Widely used. MPEG-4_AVC based on H.264 is a subtype that appears to be more widely adopted.	Open standard. Fully documented as ISO standard.

16

<sup>42</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000059.shtml>

<sup>43</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000026.shtml>

<sup>44</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000200.shtml>

## 1 3.17 Electronic Publications

2 eBooks are used to publish books in a form that they can easily be read on a digital device such as a tablet.

Recommended	Notes	Rationale
PDF	See PDF in Section 2.3.  Presentations in this format will have restricted functionality – it is difficult to extract information from them and to modify them.	Open standard. Fully documented. Widely used.
EPUB <sup>45</sup>	EPUB is a format for electronic publications with reflowable text in marked up document structure with associated images for illustrations, all in a container format. Reflowable text allows the text display to dynamically resize to fit screen size.	Popular recent standard. Fully documented. Widely used.

3

## 4 3.18 Encapsulation

5 Encapsulation formats group a number of files together into one ‘archive’ file. Probably the best known example of an  
6 encapsulation format is the ‘ZIP’ file format. Another name for this type of file is an ‘archive file’, but we prefer not to  
7 use this term as these formats are not designed for archival use.

8 Encapsulation formats are particularly easy to migrate, as they provide a simple representation of a file system. It is  
9 consequently easy to unpack the encapsulation and re-encapsulate in another format.

10

Recommended	Notes	Rationale
ZIP <sup>46</sup>	GZIP and ZIP files are completely different. A program that reads one format will not necessarily be able to read the other.	ZIP is very widely used. Although not a standard, the format is openly published, and ZIP is used by a number of International standards.
GZIP <sup>47</sup>	GZIP is primarily intended as a compression program. Although multiple files can be included in a GZIP encapsulation, this is not its most common use.  GZIP and ZIP files are completely different. A program that reads one format will not necessarily be able to read the other.	GZIP is widely used, particularly in UNIX/LINUX environments
TAR	The age of TAR can be judged by its meaning: ‘Tape Archive’; it was originally intended to encapsulate files for writing onto magnetic tape.	TAR is widely used in UNIX/LINUX environments. It is defined in an open standard (POSIX).

11

<sup>45</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000278.shtml>

<sup>46</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000354.shtml>

<sup>47</sup> <https://en.wikipedia.org/wiki/Gzip>

### 1 3.19 CAD

2 At their basic, Computer Aided Design (CAD) files represent plans of objects (technical drawings). However, modern  
3 CAD files can model complex 3D shapes and can associate properties and metadata with objects in the plans. CAD files  
4 are consequently far more complex than simple representations of paper plans.

5 Because of their complexity of the data, and the dominance of a few commercial products, users generally have a  
6 choice between widely used proprietary formats, and standard formats that are not widely used, or which lack  
7 features of the proprietary formats.

Recommended	Notes	Rationale
DXF <sup>48</sup>	(Autodesk) Drawing Exchange Format. This is a proprietary format.  DXF files come in two versions: ASCII and binary. The ASCII version of this format should be used, as the binary is not widely used.  Should only be used for 2D models, not 3D models, as 3D models are not fully documented in the public releases of the specification.	The format is proprietary, but is openly published by Autodesk and made available under a CC attribution- non commercial – sharealike 3.0 license.
DWG <sup>49</sup>	(Autodesk) Drawing. Extremely widely used as an exchange format, particularly for 3D models.  This format is proprietary, and has not been published by Autodesk. The format has, however, been reverse engineered by a consortium and this specification is available.	While proprietary and not formally published, the format is extremely widely supported by CAD products. It is also extremely widely used. A published, reverse engineered, specification does exist.
STP, STEP, P21 <sup>50</sup>	Standard for the Exchange of Product Model Data (ISO 10303-28:2007)	An international standard. Does not appear to have the product or usage support that DWG/DXF have.
PDF/E	PDF/Engineering (ISO 24517-1:2008).  This format is relatively new and little experience is available as to the ability to represent all the information required in a CAD file – particularly where the CAD information is active and continues to be used.	An international standard.

8

### 9 3.20 Geospatial

10 Geospatial formats capture information related to the Earth. The classic piece of geospatial data is a map, however,  
11 geospatial data is best thought of as a set of geographical labelled features (points, lines and polygons) with associated  
12 properties.

13 Like CAD data, there are a number of widely used proprietary geospatial formats and a set of standardised, but less  
14 powerful, geospatial formats.

<sup>48</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000446.shtml>

<sup>49</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000445.shtml>

<sup>50</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000448.shtml>

Recommended	Notes	Rationale
ESRI Shapefile. <sup>51</sup>	<p>An open representation of geospatial information. It is important to note that, unlike other file formats discussed in this document, ESRI geospatial information is distributed amongst a number of different files. Related files have the same name, but different file extensions (e.g. SHP, SHX, DBF, CPG, PRJ, SBN, and SBX).</p> <p>The ESRI format dates from the 1990s and is a defacto standard for the exchange of geospatial information. The age of the specification means that some features of modern geospatial systems are not supported.</p>	<p>Fully documented, and the specification is openly available. Widely used for the exchange of geospatial information.</p> <p>Open libraries are freely available for reading and writing to these geospatial formats.</p>
GeoPackage <sup>52</sup>	<p>GeoPackage is an open, standards-based, platform-independent, portable, self-describing, compact format for transferring geospatial information.</p>	<p>Fully documented, with support for many geospatial data types. GeoPackage has strong open standards support, through OGC.</p> <p>Has particular application to devices where internet connectivity and bandwidth are limited.</p>
GeoJSON <sup>53</sup>	<p>GeoJSON is an open standard data interchange format designed for representing simple geographical features, along with their non-spatial attributes. It is based on JSON, the JavaScript Object Notation. TopoJSON is an extension of GeoJSON that encodes topology.</p>	<p>This fully documented format is an open standard, maintained by the Open Geospatial Consortium. It is used by many open source spatial systems, and is being rapidly adopted for Web applications involving mapping, such as those processing high volumes of Internet of Things 'IoT' data.</p>
DEM <sup>54 55</sup>	<p>A digital elevation model (DEM) represents terrain elevations for ground positions at regularly spaced horizontal intervals (i.e. 'bare earth' with no vegetation or man-made features, referenced to a common vertical datum). Single file ASCII text format. DEM is often used as a generic term for digital terrain models (DTM) and digital surface models (DSM), which include features on the earth's surface.</p>	<p>It is an open standard, maintained by the US Geographical Survey (USGS) and is used throughout the world. It was superseded by the SDTS format (no longer preferred) but the format remains popular due to large numbers of legacy files, self-containment, relatively simple field structure and broad, mature software support.</p>
GML <sup>56</sup>	<p>Geography Markup Language. This is a full featured representation of geospatial information.</p>	<p>Openly documented format by the Open Geospatial Consortium, based on XML. The standard is relatively new, and adoption is still occurring.</p>
KML <sup>57</sup>	<p>Keyhole Markup Language. XML based representation of geographic data. KMZ is the Zip</p>	<p>Openly documented format by the Open Geospatial Consortium. Textual format based on</p>

<sup>51</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000280.shtml>

<sup>52</sup> <https://www.goepackage.org/>

<sup>53</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000382.shtml>

<sup>54</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000285.shtml>

<sup>55</sup> <https://gisgeography.com/dem-dsm-dtm-differences/>

<sup>56</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000296.shtml>

<sup>57</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000340.shtml>

Recommended	Notes	Rationale
	<p>version.</p> <p>Primarily used for the publication of data ready for visualisation, particularly for non-specialists.</p> <p>GML and KML are complementary. Whereas GML is a way of modelling or encoding geographic content, KML is a way of presenting that content visually.</p>	<p>XML. Widely used (particularly in Google Earth and Google Maps).</p>
GeoTIFF, BigTIFF <sup>58</sup>	<p>Geospatial Tagged Image File Format. This is a profile of TIFF (raster image). The profile adds the ability to store georeferenced and geocoding information to a raster image (e.g. map, aerial photograph, satellite image). BigTIFF is a variant that supports very large raster images.</p> <p>Both are openly available specifications.</p>	<p>Openly available, widely supported in geospatial systems.</p>
ECW <sup>59</sup>	<p>ECW (Enhanced Compression Wavelet) is a proprietary wavelet compression image format optimized for aerial and satellite imagery. The lossy compression format efficiently compresses very large images with fine alternating contrast while retaining their visual quality.</p>	<p>ECW is a ubiquitous and stable format. While proprietary, the specification is available and the format is vendor-supported. Backwards-compatible.</p>
JPEG2000	<p>Open-source raster format (refer also section 2.11 Image - Raster). Similar functionality to ECW, apart from compression types.</p>	<p>Openly available, widely supported in geospatial systems.</p>
SVG <sup>60</sup>	<p>Scalable Vector Graphics</p>	<p>This is an open standard from the W3C using XML. Reasonably widely supported, although an older format now and less relevant in the GIS industry these days.</p> <p>Vector caches are the current way to render vector out as image tiles. An example is the Open GIS Consortium's (OGC) GeoPackage:  <a href="https://www.geopackage.org/">https://www.geopackage.org/</a></p>

<sup>58</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000279.shtml>

<sup>59</sup> [https://www.gdal.org/frmt\\_ecw.html](https://www.gdal.org/frmt_ecw.html)

<sup>60</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000020.shtml>

Recommended	Notes	Rationale
LAS <sup>61</sup>	LAS is a binary file format for the exchange of LIDAR (Light Detection and Ranging) 3D point cloud data. Although point cloud data can be described in ASCII text files, such files are very clumsy to use, because of the time required to import data and convert the numbers to binary for analysis. In practice, LIDAR point cloud datasets are usually shared in LAS files or losslessly compressed derivatives of LAS (e.g. LAZ). Used to generate other geospatial products, such as digital elevation models, canopy models, building models, and contours.	Most widely supported LIDAR format. Maintained as a public specification by the American Society for Photogrammetry and Remote Sensing (ASPRS).
IMG <sup>62</sup>	This format is used for processing remote sensing data, since it provides a framework for integrating sensor data and imagery from many sources.	A widely used proprietary, partially documented format for multi-layer geo-referenced raster images.

1

2

---

<sup>61</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000418.shtml>

<sup>62</sup> <https://www.loc.gov/preservation/digital/formats/fdd/fdd000420.shtml>

# 4 Long term preservation format criteria

## 1 What characteristics distinguish low risk preservation formats?

---

2 These criteria are aspirational and apply to born-digital records and formats for digitisation. Not all formats will meet  
3 all of these criteria all of the time:

- 4 1. **Ubiquitous:** The format is widely used and supported around the world.
- 5 2. **Unrestricted:** The format is free from patents and legal encumbrance (including intellectual property rights)  
6 and preferably embodies open-source principles.
- 7 3. **Well Documented:** The format is identifiable and is well documented (the specification is publicly available  
8 and sufficiently detailed that groups or individuals *with sufficient skills and incentive can build software to*  
9 *render it accurately*.
- 10 4. **Stable:** The file format is stable (rare releases of newer versions) and is backwards and forwards compatible,  
11 or has a clear migration path.
- 12 5. **Platform Independent:** The format should be supported by a wide range of software or is platform-  
13 independent.
- 14 6. **Uncompressed:** Ideally the format should be uncompressed. If compression is used, lossless compression is  
15 preferred.
- 16 7. **Supported:** Technical support is readily available from vendors, community or third parties.
- 17 8. **Metadata Friendly:** File formats with metadata support are preferred.

18  
19 A discussion about the meaning of these criteria can be found in the next section.

# 5 Commentary

## 1 Explanatory information about the characteristics

---

### 2 **5.1 General**

#### 3 **5.1.1 Must a format satisfy all criteria for it to be considered low risk?**

4 No, a format could be considered low risk even though it does not satisfy all the criteria. In general, the more criteria a  
5 format satisfies, the lower the risk.

### 6 **5.2 Details of specific criteria**

#### 7 **5.2.1 Ubiquitous**

8 A format is considered to be ubiquitous if it dominates its information category. For example, the vast majority of word  
9 processing documents are in the Microsoft Word (.doc, .docx) format. Similarly, almost all digital photographs captured  
10 are in the JPEG format.

11 Ubiquitous formats are considered to be low risk for longevity because the sheer amount of information in that format  
12 means that there is a significant economic incentive for applications to support that format. For example, image  
13 processing software that does not support the JPEG format is unlikely to sell many copies.

#### 14 **5.2.2 Unrestricted**

15 Legal restrictions on implementing software that reads or writes a format is likely to reduce the range of software that  
16 handles the format. This, in turn, increases the risk that the format will become unsupported.

17 Some formats may have different restrictions on the production of format readers vs writers (e.g. PDF). Restrictions on  
18 software that writes the format are less serious than those for readers.

19 Open source software is preferred but not required. By open source, we mean that the source code for software that  
20 handles the format can be downloaded and is sufficiently free from restrictions that it can be used to maintain in  
21 house (or could be used to reverse engineer a specification that can be used to write a new piece of software). Open  
22 source is not required because economic demand for very widely used software often results in reverse engineered  
23 competing products that will read proprietary formats (a good example is the Autodesk DWG CAD format). The fact  
24 that a piece of software is open source does not mean that software is being actively maintained, or is of sufficient  
25 quality to provide accurate access to the information in a format.

#### 26 **5.2.3 Well Documented**

27 A well-documented format allows future users to accurately implement software that reads the format. A format that  
28 has no public documentation can be extremely difficult to re-implement.

29 The best documentation has been formally published by independent bodies (e.g. ISO, CCITT, W3C, or the IETF).  
30 However, format descriptions produced by vendors themselves can still be valuable (e.g. the ZIP specification).

31 It is important that the documentation accurately documents the format. Some vendors extend the documented  
32 format and use un-documented extensions.

1 The best indication of the quality of a specification is that multiple independent implementations have been produced  
2 from the specification (this links to interoperability, and ubiquity). Some formats, e.g., QuickTime and MPEG-4, allow  
3 for a very wide range of implementations compared to say, MPEG-2<sup>63</sup>.

4 Note that open source software code itself is an excellent form of documentation.

#### 5 5.2.4 Stable

6 A stable or robust format means that the economic investment in writing software to deal with the format is long-  
7 lived. This encourages the production of software. It also means that the software itself has more chance of being  
8 robust and stable; that is bugs have been found and eliminated.

9 Where formats do change, but are backwards and forwards compatible, this means that software can read documents  
10 from newer versions of the software (possibly with some loss of functionality) and older versions. This minimises the  
11 risk that documents can become unreadable.

12 Formats which have a clear migration path also have lower risk, for example migrations to Office 365 from Microsoft  
13 Outlook PST files and Exchange Server are relatively straightforward whereas Office 365 migrations from Lotus Notes  
14 are possible but somewhat more complex.

#### 15 5.2.5 Platform Independent

16 Formats that have been independently implemented by a range of vendors are lower risk than those that have only  
17 been implemented by one or two. Implementation by a range of vendors shows that 1) the format is documented, 2)  
18 the documentation can be used to accurately implement the format, and 3) that the documentation actually describes  
19 the format.

20 One concern is that a format may appear to be implemented by multiple vendors, when, in fact, the different products  
21 actually use a single common implementation.

#### 22 5.2.6 Uncompressed

23 Formats that are compressed have an additional layer of complexity. Applications that need to use the data must first  
24 decompress it before it can be used. This decompression adds a risk to future use. The compression algorithm can be  
25 considered as an additional data format layered over the top of the actual data format. It can consequently be  
26 evaluated just as any other format: is it ubiquitous (e.g. JPEG); has it any restrictions (e.g. patents); and so on.

27 It should be noted that some types of information are routinely compressed to reduce their large size (e.g. audio and  
28 video files).

29 Some forms of compression are lossy; that is, they discard information to achieve a small file size and the  
30 decompressed object will have lesser quality compared with the original. Other compression algorithms are lossless,  
31 and the decompressed object will be identical to the original.

32 In general, lossless compression is preferred as this allows the original object to be retrieved. However, lossy  
33 compression is acceptable in the following circumstances:

- 34 • The original object is already compressed using a lossy compression. Examples of these would be most still images,  
35 sound files, and almost all video. All of these would normally be received from the capture device (e.g. camera)  
36 already compressed using lossy compression and usually in a proprietary format. Care must be taken to ensure  
37 that files taken directly from a capture device are in an interoperable format, or that they are transcoded to a  
38 suitable format as part of the capture workflow. In other words, agencies should avoid proprietary  
39 formats/compression techniques that are natively produced by capture devices; standard formats and  
40 compression techniques such as LZW, JPEG and JPEG2000 are preferred.

---

<sup>63</sup> [https://www.loc.gov/preservation/digital/formats/content/video\\_quality.shtml](https://www.loc.gov/preservation/digital/formats/content/video_quality.shtml)

- 1 • The original object is too large to be used for business purposes and the original will not be required. Again, the  
2 examples would be still images, sound files, and video which may need to be reduced in size in order to be usable.  
3 Consideration should be given to retaining both the high-quality original and the lower-quality reduced size copy.

#### 4 5.2.7 Supported

5 Technical support means that it is possible to get someone to assist with the format, particularly with regard to solving  
6 technical problems. Some formats are actively supported, particularly those that are commercial and public domain  
7 formats that have an active user community.

#### 8 5.2.8 Metadata Friendly

9 A digital object which incorporates all the metadata needed to render its data as usable information or understand its  
10 context are likely to be easier to sustain over the long term and less vulnerable to catastrophe than digital objects that  
11 are stored separately from related metadata. An example of a format with metadata embedded within the file is JPEG  
12 (where the image includes metadata that describes how the image was captured). This metadata can be a fixed set of  
13 attributes, or it can be arbitrary.