

ISO/IEC JTC 1/SC 34/JWG 7

Joint JTC 1/SC 34-TC 46/SC 4-IEC/TC 100/TA 10 WG: EPUB

Convenorship: KATS (Korea, Republic of)

Document type:	Text for CD ballot or comment
Title:	Text of ISO/IEC PDTS2 22424-1 Preserving Content in EPUB Format - Part 1: Principles
Status:	This document is prepared by project editor, Juha Hakala, to reflect all comments from the CRM of ISO/IEC PDTS 22424-1 Preserving Content in EPUB Format - Part 1: Principles.
Date of document:	2018-07-24
Source:	Juha Hakala (project editor)
Expected action:	INFO
No. of pages:	39
Email of convenor:	samoh21@gmail.com , zzosang@gmail.com
Committee URL:	https://isotc.iso.org/livelink/livelink/open/jtc1sc34jwg7

**Information technology — Digital publishing — Preserving
Content in EPUB Format - Part 1: Principles**

**Technologies de l'information - Édition numérique - Archivage
pérenne de l'EPUB 3 - Partie I : Principes**

DTS stage

Warning for WDs and CDs

This document is not an ISO International Standard. It is distributed for review and comment. It is subject to change without notice and may not be referred to as an International Standard.

Recipients of this draft are invited to submit, with their comments, notification of any relevant patent rights of which they are aware and to provide supporting documentation.

1
2
3
4
5
6
7
8
9
10
11
12

© ISO 2018, Published in Switzerland

All rights reserved. Unless otherwise specified, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
Ch. de Blandonnet 8 • CP 401
CH-1214 Vernier, Geneva, Switzerland
Tel. + 41 22 749 01 11
Fax + 41 22 749 09 47
copyright@iso.org
www.iso.org

13	Contents	
14	Foreword	iv
15	Introduction	v
16	1 Scope	1
17	2 Normative references	1
18	3 Terms and definitions	1
19	4 Abbreviated terms	11
20	5 Packaging standards	11
21	6 Construction of OAIS information packages	13
22	6.1 General	14
23	6.2 Identification of information packages and their content	19
24	6.3 Structure of information packages	20
25	6.4 Generic Information package metadata	21
26	Annex A (informative) EPUB and digital preservation: issues and recommendations	23
27	Bibliography	26
28		

29 Foreword

30 ISO (the International Organization for Standardization) is a worldwide federation of national
31 standards bodies (ISO member bodies). The work of preparing International Standards is normally
32 carried out through ISO technical committees. Each member body interested in a subject for which a
33 technical committee has been established has the right to be represented on that committee.
34 International organizations, governmental and non-governmental, in liaison with ISO, also take part in
35 the work. ISO collaborates closely with the International Electrotechnical Commission (IEC) on all
36 matters of electrotechnical standardization.

37 The procedures used to develop this document and those intended for its further maintenance are
38 described in the ISO/IEC Directives, Part 1. In particular the different approval criteria needed for the
39 different types of ISO documents should be noted. This document was drafted in accordance with the
40 editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives).

41 Attention is drawn to the possibility that some of the elements of this document may be the subject of
42 patent rights. ISO shall not be held responsible for identifying any or all such patent rights. Details of
43 any patent rights identified during the development of the document will be in the Introduction and/or
44 on the ISO list of patent declarations received (see www.iso.org/patents).

45 Any trade name used in this document is information given for the convenience of users and does not
46 constitute an endorsement.

47 For an explanation on the voluntary nature of standards, the meaning of ISO specific terms and
48 expressions related to conformity assessment, as well as information about ISO's adherence to the
49 World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT) see the following
50 URL: www.iso.org/iso/foreword.html.

51 This document was prepared by Technical Committee ISO/IEC JTC 1, Joint Technical Committee,
52 Subcommittee SC 34, Document Description and Processing Language.

53 ISO/IEC DTS 22424 consists of the following parts, under the general title *Information technology —*
54 *Digital publishing — EPUB 3 Preservation*:

55 — *Part 1: Principles*

56 — *Part 2: Metadata requirements*

57

58 Introduction

59 This document facilitates the long-term preservation of EPUB publications by specifying EPUB features
60 which are mandatory for long-term preservation (such as font embedding) and features which should
61 be avoided (including fixed layout properties).

62 This specification is related to EPUB in the same way as PDF/A, specified in ISO 19005-1 – 19005-3, is
63 related to PDF. If the EPUB community develops detailed guidelines for the production of archivable
64 EPUB publications, this document can be used as one of the starting points.

65 Long term preservation in general requires two things:

- 66 • making the object such as EPUB publication fit for preservation – including features to be used
67 and features to avoid;
- 68 • the packaging of the object (and any metadata related to it) together with any additional data
69 such as other versions of the object and other documentation into an OAIS Submission
70 Information Package (SIP).

71 Packaging is covered in Part 2 of this technical specification.

72 EPUB

73 The EPUB standard

74 *defines a distribution and interchange format for digital publications and documents. The*
75 *EPUB® format provides a means of representing, packaging and encoding structured and*
76 *semantically enhanced Web content — including HTML, CSS, SVG and other resources — for*
77 *distribution in a single-file container [EPUB 3.1].*

78 EPUB format was developed by the International Digital Publishing Forum, IDPF, which merged with
79 the World Wide Web Consortium, W3C, in January 2017. Ongoing technical development of the
80 standard, related extension specifications and ancillary deliverables are the responsibility of the W3C
81 EPUB 3 Community Group¹, which published its charter in February 2017. According to the charter,

82 *work on any future major revision of EPUB, e.g. an EPUB 4, is initially out of scope on the*
83 *presumption that this will be taken up by a new W3C WG as a W3C [Recommendation Track](#)*
84 *activity. The EPUB 3 CG will coordinate its work with such new WG, and meanwhile with the*
85 *existing W3C [Digital Publishing Interest Group](#) (DPUB IG). [W3C]*

86 The International Digital Publishing Forum, IDPF, has ceased operations as a membership organization
87 in January 2017, and its website² is now an archive. The latest version of the standard and information
88 about future EPUB developments is available at the Publishing@W3C webpage,
89 <https://www.w3.org/publishing/>.

90 The specification at hand covers EPUB 3 versions up to EPUB 3.1³, which is the first major revision of
91 EPUB 3.0. EPUB 3.0.1⁴, as of this writing the most widely used version of the standard, was a minor
92 update of EPUB 3.

93 Differences between EPUB specification are well documented:

¹ <https://www.w3.org/publishing/groups/epub3-cg/>

² <http://idpf.org/>

³ <https://www.w3.org/Submission/epub31/>

⁴ <http://idpf.org/epub/301>

94

- 95 • EPUB 3 Changes from EPUB 2.0.1⁵
- 96 • EPUB 3.0.1 Changes from EPUB 3.0⁶
- 97 • EPUB 3.1 Changes from EPUB 3.0.1⁷

98

99 All EPUB specifications are available in the Web; 2.0.1 at <http://idpf.org/epub/201>, EPUB 3.0.1 at
100 <http://idpf.org/epub/301>.

101 As a rule the differences between 3.x versions are not critical from long term preservation point of view.
102 There are two exceptions which concern foreign resources and fixed layout documents.

103 In EPUB 3.1 foreign resources do not require fallbacks if they are not in the spine and not embedded in
104 EPUB Content Documents. In EPUB 3.0.1, fallback guarantees that there is a version of the document
105 that can be rendered; in 3.1 such guarantee no longer exists.

106 Fixed layout documents were introduced in EPUB 3.0.1. Since it is difficult to preserve the original look
107 and feel in the long term, reflowable EPUB is easier to preserve than fixed layout documents since the
108 presentation and meaning of the document are not interconnected.

109 This specification recommends that even foreign resources should have fallbacks and that fixed layout
110 specification should be avoided, except the SIP contains both a fixed layout version and a reflowable
111 version of the document. These limitations apply only to specific EPUB versions (3.1 and 3.0.1,
112 respectively), other requirements to all three 3.x versions.

113 EPUB 3.1 is also the last revision of the EPUB 3 standard which was prepared by the IDPF. There are six
114 EPUB 3.1 base specifications, each defining a component of a general EPUB publication:

- 115 • EPUB 3.1 [EPUB specification]⁸, a blanket document providing a good point of entry to the EPUB
116 standard; includes e.g. common terms and definitions.
- 117 • EPUB Packages 3.1⁹ defines the package semantics and conformance requirements.
- 118 • EPUB Content Documents 3.1¹⁰ defines the usage of HTML, SVG, and CSS optimized for the
119 representation of structured, composable, and accessible documents.
- 120 • EPUB Open Container Format (OCF) 3.1¹¹ defines the file format and processing model for
121 encapsulating a set of related resources into a single-file container.
- 122 • EPUB Media Overlays 3.1¹² defines the usage of SMIL, the Package Document; CSS Style Sheets;
123 and EPUB Content Documents for the representation of audio, synchronized with an EPUB
124 Content Document.

125 There are several extension specifications to the EPUB base standards. The list below is incomplete, as
126 it only contains the specifications that are relevant from the long-term preservation point of view. Some
127 of them are still drafts:

- 128 • EPUB 3 Fixed Layout Documents¹³ defines a set of metadata properties to allow declarative
129 expression of intended rendering behaviors of fixed-layout documents in the context of EPUB 3.

⁵ <http://www.idpf.org/epub/30/spec/epub30-changes-20111011.html>

⁶ <http://www.idpf.org/epub/301/spec/epub-changes-20140626.html>

⁷ <http://www.idpf.org/epub/31/spec/epub-changes-20170105.html>

⁸ <http://www.idpf.org/epub/31/spec/epub-spec-20170105.html>

⁹ <https://www.w3.org/Submission/2017/SUBM-epub-packages-20170125/>

¹⁰ <https://www.w3.org/Submission/2017/SUBM-epub-contentdocs-20170125/>

¹¹ <https://www.w3.org/Submission/2017/SUBM-epub-ocf-20170125/>

¹² <https://www.w3.org/Submission/2017/SUBM-epub-mediaoverlays-20170125/>

¹³ <http://www.idpf.org/epub/fxl/>. This specification has been superseded; fixed-layout metadata is now defined in EPUB Publications 3.1, chapter 4.4.2 Fixed layout properties

- 130 • EPUB Canonical Fragment Identifiers 1.1¹⁴ defines a standardized method of referencing
- 131 content within an EPUB publication through the use of URI fragments.
- 132 • EPUB Previews 1.0¹⁵ describes how content previews can be included in EPUB publications.
- 133 • EPUB Distributable Objects 1.0¹⁶ is a draft specification that defines a method for the
- 134 encapsulation, transportation, and integration of distributable objects in EPUB publications.
- 135 • EPUB Scriptable Components 1.0¹⁷ provides an interoperable publish and subscribe (pubsub)
- 136 pattern by which interactive content can be created and incorporated into EPUB publications.
- 137 Same as EPUB Distributable Objects, it is as of this writing¹⁸ a draft.
- 138 • EPUB Scriptable Components Packaging and Integration 1.0¹⁹ is a draft that defines a method
- 139 for the creation and inclusion of dynamic and interactive components in EPUB publications.
- 140 • EPUB Multiple-Rendition Publications 1.0²⁰ defines the creation and rendering of EPUB
- 141 publications consisting of more than one rendition.
- 142 • EPUB Dictionaries and Glossaries 1.0²¹ provides a means for expressing dictionary and glossary
- 143 semantics in EPUB publications.

144
 145 EPUB 3 Core Media Types are listed at <https://idpf.github.io/epub-cmt/v3/>. The EPUB working group
 146 may approve new media types or deprecate old ones at any time, and the possible changes will apply to
 147 all the EPUB versions. As of this writing [2018-01-12], the latest update has been made on October 6,
 148 2016.

149
 150 In 2014, EPUB 3.0 specifications were republished as a standard, ISO/IEC TS 30135 parts 1-7, by the
 151 International Standards Organization. Each of these seven ISO specifications is identical to its IDPF
 152 equivalent, for example TS-30135-1 has exactly the same content as the EPUB 3.0 Overview.

153
 154 ISO/IEC JTC 1/SC 34 is currently updating the ISO standard to match the version 3.0.1. No EPUB
 155 extension specifications such as EPUB 3 Fixed Layout Documents (see the list above) have not been
 156 included in the ISO standard yet.

157
 158 EPUB is a rich document format with a lot of features. From the digital preservation point of view this is
 159 a challenge. Preserving all aspects and features of EPUB publications may be difficult, since there are
 160 features which are difficult to preserve and therefore should be avoided, such as fixed layout properties.
 161 Moreover, EPUB reading systems usually do not support all features of the specification and finding
 162 tools supporting rare features can be difficult.

163
 164 In spite of these challenges EPUB is generally regarded as a suitable format for digital archiving. For
 165 instance, the Finnish National Digital Library initiative has selected just eight archivable file formats for
 166 text, EPUB being one of them. The selection criteria were openness/transparency, adoption as a
 167 preservation standard, degree of forward/backward compatibility, degree of protection against file
 168 corruption, frequency of version releases, dependencies/interoperability, and standardization. EPUB
 169 got an A, the best grade, from everything else except the second and third criterion. For those, the grade
 170 was the second best, a B [File formats, p. 42]. Based on these generic criteria, EPUB seems to provide a
 171 good basis for long-term preservation, although additional guidelines on how to use the standard are
 172 needed to guarantee EPUB files can be preserved efficiently.

173
 174 The British Library's Digital Preservation Team has published an assessment of EPUB as a preservation
 175 format [Day]. It covers EPUB versions 3.0.1 and 2 and the overall view of EPUB is positive [Day, p. 2]:
 176

¹⁴ <http://www.idpf.org/epub/linking/cfi/epub-cfi.html>

¹⁵ <http://www.idpf.org/epub/previews/epub-previews-20150826.html>

¹⁶ <https://w3c.github.io/publ-epub-revision/do/epub-do.html>

¹⁷ <https://w3c.github.io/publ-epub-revision/sc/sc-api.html>

¹⁸ 2017-10-10

¹⁹ <https://w3c.github.io/publ-epub-revision/sc/sc-packaging.html>

²⁰ <http://www.idpf.org/epub/renditions/multiple/>

²¹ <http://www.idpf.org/epub/dict/>

177 *EPUB 3 is currently the closest thing available to an open standard for e-books. In 2013,*
 178 *Bläsi and Rothlauf concluded that EPUB 3 had the “highest expressive power” of all formats*
 179 *in the e-book ecosystem, and that it included the superset of all features used in proprietary*
 180 *formats like KF8, Fixed Layout EPUB, and iBooks.*

182 EPUB is enjoying reasonable support in the e-book market. Many suppliers, publishers, and application
 183 developers who have supported EPUB 2 have implemented version 3.0.1. According to the EPUBTest
 184 web site²², EPUB 3 support in reading systems is far from exhaustive, but market coverage is good – in
 185 January 2018, there were 59 reading systems which support at least some of the features specified in
 186 EPUB 3.0.

187
 188 E-book suppliers have produced EPUB 3 based formats that incorporate Digital Rights Management
 189 (DRM), and EPUB modifications that may restrict using the format on other than the suppliers’ own
 190 platforms. For example, the Kindle Fire eReader, released in 2015, uses a new format called Kindle
 191 Format 8 (KF8), which is partly based on EPUB 3, with Amazon’s DRM. [Day, 3]. Publisher/supplier
 192 specific DRM often restricts the use of e-books to that publisher’s/supplier’s rendering devices and/or
 193 applications, and is therefore a major obstacle to digital preservation [Day, p. 7].

194
 195 The EPUB specification does not enforce a particular Digital Rights Management scheme, but DRM may
 196 be layered on top of the EPUB specifications. A producer can, for instance, use one of the three major
 197 rights management systems in the market (Amazon DRM, Apple FairPlay DRM for books bought from
 198 iBooks, and Adobe DRM), or some other DRM system along with some additional platform-targeting.

199
 200 DRM protection should be removed from EPUB publications during pre-ingest by producer or as a part
 201 of the ingest process by the OAIS archive. In practice, only national libraries may be able to do this,
 202 provided that legal deposit act and / or copyright act guarantee them such privilege. If migration is the
 203 chosen preservation strategy, existing EPUB publications will be converted into more modern EPUB
 204 versions when rendering tools for old versions are no longer available, and (eventually) migrated into
 205 other formats.

206
 207 If archival copies of EPUB publications are not directly accessible by the public, removing DRM, digital
 208 watermarking, and other protection mechanisms from the archived documents is not a risk. When
 209 publications are delivered to the customers as Dissemination Information Packages (DIPs), the archive
 210 shall use a combination of administrative and technical means to protect the documents as required in
 211 the submission agreement. These means may include DRM protection mechanism into the package
 212 submitted to the user according to the requirements of the submission agreement. The agreement may
 213 also specify the customers the archive is entitled to serve; for instance, it is possible to require that the
 214 preserved documents can only be disseminated to the producer, and the producer will serve the end-
 215 users who do not have access the repository system.

216 **Digital preservation**

217
 218 The information society is dependent on successful long-term digital preservation. When an increasing
 219 percentage of information is produced and published only in a digital format, it is important to make
 220 sure that this information remains available in the distant future.

221 Digital preservation is not about preserving just bits, but about preserving access. The “business logic”
 222 is as follows:

- 223 • we need software and hardware to render content for human users
- 224 • software changes over time; there are new versions from old applications, and entirely new
- 225 applications

²² <http://epubtest.org/testsuite/epub3/>

- 226 • new or updated applications may not be able to render outdated file formats or format versions
227 correctly
- 228 • digital preservation makes an effort to have all archived content in stable formats. Publications
229 should also contain the smallest possible amount of features which are not commonly
230 supported in software packages used to render the content in these formats, and also avoid
231 adding links to external resources since then the long-term access to the publication requires
232 also persistence of these external resources.
- 233 • when necessary, data in old formats may be migrated into more modern formats or updated
234 versions of the same format. For instance, an e-book in EPUB 3.0.1 format may be migrated to
235 EPUB 5.2. when version 3.0.1 is no longer widely supported by reading systems.
- 236 • since the aim is to preserve the content, not the bits, the bits may change as a result of version
237 updates and format migrations.
- 238 • Many OAIS archives preserve successive versions of archives publications, because migration
239 may change the look and feel of the original document, or even its intellectual content.

240 In many countries, national libraries are responsible for preserving the published cultural heritage for
241 the future generations, while national archives take care of governmental publications, irrespective of
242 which format they are available in. All of these resources have to be preserved for decades, centuries
243 even. Then again, publishers may guarantee continuous access to the subscribers of electronic serials
244 and other licensed content. If this is so, either the publisher or a third-party should look after the
245 publications and make sure they remain accessible or at least available.

246 Ordinary digital asset management systems are not suitable for long-term preservation; therefore it is a
247 normal practice to separate short-term and long-term information management into different systems.
248 However, this does not mean that digital archiving is independent of the routine life cycle of documents.
249 Digital preservation is a long process that begins when publications are created.

250 The presence of preservation metadata such a, which allows the publication to be found, rendered and
251 authenticated correctly, is a prerequisite for digital preservation. Some preservation metadata elements
252 can only be provided by the original creator of the publication. It is also important to keep preservation
253 requirements in mind when preparing a publication. Any feature in a file format can be either essential,
254 useful, neutral, questionable, or even downright counterproductive from a long-term preservation point
255 of view. However, publishers are likely to use the features that let them achieve their own goals, and
256 preservation may not be among them.

257 There already are archivable versions of some file formats. PDF/A (ISO 19005-1:2005)²³ is probably the
258 best known example. It specifies how to use the Adobe Portable Document Format (PDF) for long-term
259 preservation. An example of a counterproductive feature for preservation in PDF is font referencing;
260 therefore in PDF/A all fonts shall be embedded in order to guarantee that the document can be
261 rendered correctly.

262 PDF/A forbids also the use of encryption, because encryption is generally regarded as a risk for long-
263 term preservation. But storing unencrypted documents is a risk as well, because if they are stolen, non-
264 authorized usage is easy. Therefore, according to the Digital preservation handbook [Digital]:

265 *Information security methods such as encryption add to the complexity of the preservation*
266 *process and should be avoided if possible for archival copies. Other security approaches may*
267 *therefore need to be more rigorously applied for sensitive unencrypted files; these might*

²³ <https://www.iso.org/standard/38920.html>

268 *include restricting access to locked-down terminals in controlled locations (secure rooms),*
269 *or strong user authentication requirements for remote access.*
270

271 In order to guarantee the correct processing of PDF/A files, there are specific requirements for PDF/A
272 reading systems, such as support for embedded fonts. There are three versions of the specification:
273 PDF/A-1 is based on PDF 1.4, PDF/A-2 adds features from PDF 1.5, 1.6 and 1.7, and PDF/A-3 contains
274 all the features of PDF/A-2 as well as allows the embedding of other file formats into PDF/A conforming
275 documents [PDF/A].

276 The TI/A (Tagged Image for Archival) standard initiative intends to create an ISO recommendation to
277 optimize the format specification for archival purposes. The motivation behind the initiative applies
278 perfectly to other image formats, but there are valid points to the EPUB community as well [TI/A]:

279 *The versatility of the TIFF format has made it very attractive for memory institutions for*
280 *long term archival of their digital images. However, since the TIFF format offers such a great*
281 *flexibility, it is not guaranteed that in the future a standard TIFF reader will be able to read*
282 *some TIFF images.*

283 *The limitations of the baseline TIFF are too severe for many applications in digital archiving.*
284 *It is important that, besides crucial technical metadata such as ICC color profiles (in case of*
285 *color images) also important descriptive metadata is stored within the image file. Having*
286 *descriptive metadata available (such as content description, iconography, copyright and*
287 *ownership information etc.) is crucial for every archive. Having this information in the same*
288 *file as the image data guarantees that this information will always be associated with the*
289 *image.*

290 TIFF is as of this writing not an EPUB core media type, but four other image types have been listed; GIF,
291 JPEG, PNG, and SVG. It is significant from a digital preservation point of view how these formats and
292 other core media types are used in the EPUB context. Image and audio files embedded in an EPUB
293 publication may require migration before the EPUB publication itself has to be migrated into a more
294 modern file format, if commonly available EPUB reading systems no longer support these file formats.
295 This specification does not provide guidelines for creating archivable files in EPUB 3 core media types,
296 due to the magnitude of such a task. But EPUB community SHOULD follow the archival file format lists
297 of national archives or libraries (for example the Library of Congress file format list²⁴ and the U.S.
298 National Archives list²⁵) when the core media file format list is updated. Publishers SHOULD also
299 consider the persistence of file formats used when creating EPUBs for which the need for long-term
300 preservation is foreseen. .

301 This specification does not require any changes to be made to the EPUB standard or to any future
302 versions of it. However, with each new EPUB standard version it is necessary to check if the ISO 22424
303 needs to be revised, since any new EPUB features may be either useful, counterproductive, or irrelevant
304 from a long-term digital preservation point of view. A similar approach is already in place for PDF/A:
305 ISO 19005-1 applies to PDF 1.4, and ISO 19005-2 covers the subsequent PDF versions up to 1.7.

306 **OAIS and related standards**

307 ISO 22424 provides guidance on how to utilize the Open Archival Information System (OAIS) and
308 current practices of OAIS archives in preservation of EPUB publications. The OAIS [ISO 14721] is
309 equally relevant to both parts of the ISO 22424.

310 OAIS is a reference model for long-term data storage systems. It is used by memory institutions
311 (libraries, archives, and museums) and many other organizations that need to preserve digital

²⁴ <http://www.loc.gov/preservation/digital/formats/>

²⁵ <https://www.archives.gov/records-mgmt/policy/transfer-guidance-tables.html>

312 resources in the long-term. Although an ISO standard, the OAIS was originally developed by the CCSDS,
 313 The Consultative Committee for Space Data Systems²⁶, which still maintains the specification.

314 The model has five functional units:

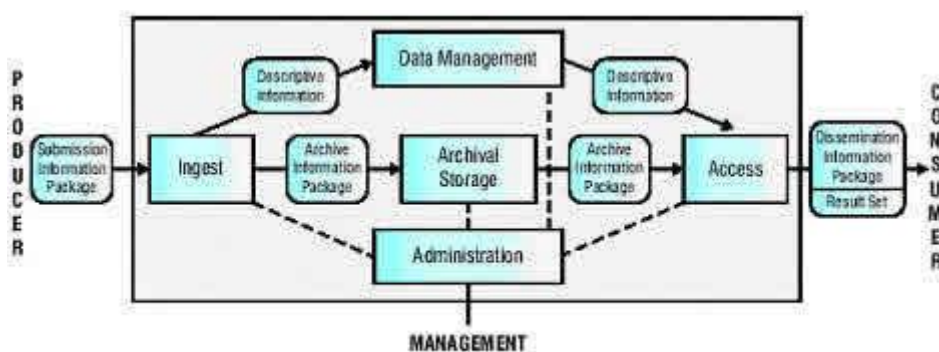


Figure 3

315

316

Figure 1. OAIS Model [Lavoie]

317 In the model, the *Ingest function* is responsible for receiving information from producers and preparing
 318 it for storage and management within the OAIS archive. The Ingest accepts information – in this case,
 319 EPUB publications – from producers in the form of Submission Information Packages (SIPs), performs
 320 quality assurance checks on the SIP, and generates an Archival Information Package (AIP) from one or
 321 more SIPs (or multiple AIPs from a single SIP). Finally, the Ingest function transfers the new AIPs to
 322 Archival Storage and the associated Descriptive Information (metadata) to Data Management.

323 Modifying an EPUB publication so that it is suitable for digital archiving is from the OAIS point of view a
 324 part of pre-ingest and as such not a part of the OAIS model. The importance of the OAIS to the ISO
 325 22424 is that the model provides a terminology, information package data model and an overall
 326 framework within which digital preservation can be performed.

327 Neither OAIS nor this specification describe the interface between a repository system used by the
 328 archive and systems used by producers. The Producer-Archive Interface Methodology Abstract
 329 Standard, also known as PAIMAS [ISO 20652], covers the first stages of the ingest process defined by
 330 the OAIS. It provides a basis for detailed specifications on how production systems communicate with
 331 OAIS archives. One such specification is DEPIP, the Data Exchange Protocol for Interoperability and
 332 Preservation [ISO/FDIS 20614]. The DEPIP is intended for systems used by libraries, archives, and
 333 museums. Other domains are likely to create their own API specifications.

334 Of all the functional units of the OAIS model, this specification covers only the Ingest unit. In addition
 335 there are tasks that are part of non-OAIS unit Pre-ingest, or things a producer shall take care of when
 336 preparing a SIP. Other OAIS units are beyond the scope, and therefore archival or dissemination related
 337 functions such as migration or creation of dissemination information packages are discussed only in
 338 passing. It is assumed that Ingest does not require any major changes, although if EPUB for some reason
 339 were no longer approved as preservation format, the archive would be obliged to migrate the EPUB
 340 publications into eligible file format. Even then the submission agreement might require the archive to
 341 disseminate the publication back to consumers in the original EPUB format.

342 OAIS submission agreements specify the principles of how documents should be prepared and
 343 submitted to the repository system. If the archive uses migration as the preservation method²⁷,

²⁶ <https://public.ccsds.org/default.aspx>

344 submission agreements should specify file formats (and metadata formats) suitable for submission
345 and/or archival, or refer to external documents listing these formats. File formats suitable for
346 submission but not for archival are migrated during the ingest process, although the original files may
347 be included in the AIP.

348 The submission agreements may also refer to SIP schema specifications, which provide more guidelines
349 for document producers. Schemas may utilize long-term preservation standards such as METS
350 (Metadata Encoding and Transmission Standard). Together the submission agreement and related
351 documents should give a producer a clear idea on when and which publications should be sent to the
352 repository system, which file formats and metadata specifications should be used, means of data
353 transfer available etc. These requirements should cover both ingest and dissemination; that is,
354 submission of documents to the repository system by the producer, and retrieval of the archived
355 documents by customers.

356 This specification (ISO 22424 Part 1: Principles) outlines the general principles for the submission of
357 EPUB publications from digital asset management systems to repository systems. The principles of
358 archival storage or dissemination of archived documents are not covered here, because OAIS archives
359 may apply various methods and processes to meet the requirements of submission agreements. Bit
360 level preservation is also out of scope; the purpose of this specification is to make it easier for
361 producers and OAIS archives to preserve access to EPUB documents.

362 The second part of this specification (ISO 22424 Part 2: Metadata requirements) provides a technical
363 basis to meet the principles listed in this document by specifying metadata required for long term
364 preservation, and a method for packaging this metadata with the original EPUB container.

365 This specification is applicable to EPUB versions 3.x and as such it should be used cautiously with other
366 (previous or later) versions of the standard. If there is a need to preserve documents that are in earlier
367 EPUB versions, they do not need to be migrated, provided that a) submission agreement specifies those
368 EPUB versions as archivable formats, and b) there are reading systems for these EPUB versions.
369 Additional features in future EPUB versions should be analyzed from long-term preservation point of
370 view. If such analysis reveals that they may constitute a risk, they should be avoided in submitted EPUB
371 publications, or removed during ingest.

372 Annex A in this specification provides a summary of issues and recommendations related to the EPUB
373 standard and its usage from long term preservation point of view.

374

²⁷ In this document, preservation method is assumed to be migration. In practice, emulation may also be applied if it is important to preserve the original look and feel of the publication. In an ideal world such migrations between the file formats would be lossless; in practice that may not be the case. Migrated document may look different even if the content is the same, and in the worst case semantics changes as well. Therefore archives often preserve also the original version of the archived resource, alongside more modern versions.

375 Information technology — Digital publishing — Preserving Content in 376 EPUB 3 Format - Part 1: Principles

377 **1 Scope**

378 This document supports long term preservation of EPUB publications via a dual strategy. First, it
379 considers EPUB features from long term preservation point. Some EPUB features are forbidden and
380 some others required, depending on they relate to long term preservation. An EPUB document
381 constructed according to these guidelines are suitable for preservation. In this respect, this specification
382 is related to EPUB in the same way than PDF/A is related to PDF.

383 Second, this specification makes EPUB compliant with current practices of OAIS archives and technical
384 requirements of repository systems. The former tend to rely on Open Archival Information Systems
385 (OAIS) in their operations; the latter prefer to ingest electronic documents only in containers
386 conforming to standards such as METS (Metadata Encoding and Transmission Standard).

387 **2 Normative references**

388 The following documents are referred to in the text in such a way that some or all of their content
389 constitutes requirements of this document. For dated references, only the edition cited applies. For
390 undated references, the latest edition of the referenced document (including any amendments) applies.

391 ISO/IEC TS 30135, *Information technology — Digital publishing — EPUB3*

392 ISO 14721. *Space data and information transfer systems – Open archival information system (OAIS) –*
393 *Reference model*

394 **3 Terms and definitions**

395 For the purposes of this document, the following terms and definitions apply. Unless stated otherwise,
396 the terms have been adopted from ISO 14721:2012.

397 ISO and IEC maintain terminological databases for use in standardization at the following addresses:

398 — IEC Electropedia: available at <http://www.electropedia.org/>

399 — ISO Online browsing platform: available at <https://www.iso.org/obp>

400 **3.1**

401 **access functional entity**

402 OAIS functional entity that contains the services and functions, which make the archival information
403 holdings and related services visible to Consumers

404 **3.2**

405 **administrative metadata**

406 metadata that provides information to help manage a resource, such as when and how it was created,
407 file type and other technical information, and access rights

408 [SOURCE: Understanding metadata]

409 **3.3**
410 **archival information package**
411 **AIP**
412 Information Package consisting of Content Information and associated Preservation Description
413 Information (PDI), which is preserved within an OAIS

414 **3.4**
415 **archive**
416 **OAIS archive**
417 organization that intends to preserve information for access and use by a Designated Community

418 **3.5**
419 **authenticity**
420 property than an entity is what it claims to be

421 [SOURCE: ISO/IEC 27000]

422 Note 1 to entry: Authenticity is judged on the basis of evidence.

423 **3.6**
424 **bit preservation**
425 term used to denote a very basic level of preservation of digital resource as it has been submitted
426 (literally the preservation of the **bits** forming a digital resource)

427 Note 1 to entry: This may include maintaining onsite and offsite backup copies, virus checking, fixity-checking, and
428 periodic refreshing to a new storage medium.

429 Note 2 to entry: Bit preservation is not digital preservation but it does provide a building block for the more
430 complete set of digital preservation practices and processes that ensure the survival of digital content and also its
431 usability, display, context and interpretation over time.

432 [SOURCE: Digital preservation handbook, Glossary]

433 **3.7**
434 **consumer**
435 role played by those persons or client systems, who interact with OAIS services to find preserved
436 information of interest and to access that information in detail

437 Note 1 to entry: This can include other OAISs, as well as internal OAIS persons or systems.

438 **3.8**
439 **content information**
440 set of information that is the original target of preservation or that includes part or all of that
441 information

442 Note 1 to entry: It is an Information Object composed of its Content Data Object and its Representation
443 Information.

444 **3.9**
445 **context information**
446 information that documents the relationships of the Content Information to its environment

447 Note 1 to entry: This includes reasons why the Content Information was created and how it relates to other
448 Content Information objects.

449 **3.10**
450 **core media type**
451 a set of publication resource for which no fallback is required.

452 [SOURCE: EPUB Publications 3.0 Recommended Specification 11 October 2011]

453 Note 1 to entry: Core media types have been specified in chapter 5.1. of the EPUB publications specification,
454 version 3.0.1.

455 EXAMPLE core media types for still images are image/gif, image/jpg, image/png and image/svg+xml. Any
456 other still image file format is foreign and requires a fallback, meaning the same resource expressed in another
457 foreign format or core media type.

458 3.11

459 data, pl

460 reinterpretable representation of information in a formalized manner suitable for communication,
461 interpretation, or processing

462 [SOURCE: ISO 5127:2017]

463 Note 1 to entry: Data are often understood as taking the form of a set of values of qualitative or quantitative
464 variables.

465 3.12

466 data dictionary

467 organized and constructed (electronic data base) compilation of descriptions of data concepts that
468 provides a consistent means for documenting, storing and retrieving the syntactical form (i.e.
469 representational form) and the meaning and connotation of each data concept

470 [SOURCE: ISO 24531:2013]

471 Note 1 to entry: PREMIS²⁸ is a data dictionary.

472 3.13

473 descriptive metadata

474 descriptive information

475 metadata about a resource for example for discovery and identification

476 Note 1 to entry: These can include elements such as title, abstract, author, and keywords.

477 [SOURCE: Understanding metadata]

478 3.14

479 designated community

480 identified group of potential Consumers who should be able to understand a particular set of
481 information

482 Note 1 to entry: A Designated Community may be composed of multiple user communities. The community is
483 defined by an Archive, though this definition may change later on.

484 3.15

485 digital preservation

486 series of managed activities necessary to ensure continued access to digital materials for as long as
487 necessary

488 Note 1 to entry: Digital preservation refers to all of the actions required to maintain access to digital materials
489 beyond the limits of media failure or technological and organizational change

²⁸ PREMIS Data Dictionary for Preservation Metadata (<https://www.loc.gov/standards/premis/>) is a leading metadata specification for metadata needed for long-term preservation.

490 Note 2 to entry: Those materials may be records created during the day-to-day business of an organization; "born-
491 digital" materials created for a specific purpose (e.g. teaching resources); or the products of digitisation projects.

492 EXAMPLE 1 **Short-term preservation** - Access to digital materials either for a defined period of time while
493 use is predicted but which does not extend beyond the foreseeable future and/or until it becomes inaccessible
494 because of changes in technology.

495 EXAMPLE 2 **Medium-term preservation** - Access to digital materials beyond changes in technology for a
496 defined period of time but not indefinitely.

497 EXAMPLE 3 **Long-term preservation** - Access to digital materials, or at least to the information contained in
498 them, indefinitely.

499 [SOURCE: Digital preservation handbook, Glossary]

500 **3.16**

501 **digital rights management**

502 **DRM**

503 packaging, distributing, controlling, and tracking content based on rights and licensing information

504 [SOURCE: ISO 19153:2014]

505 **3.17**

506 **digital signature**

507 **signature**

508 data appended to, or a cryptographic transformation of, a data unit that allows the recipient of the data
509 unit to prove the source and integrity of the data unit and protect against forgery, e.g. by the recipient

510 [SOURCE: ISO/IEC 19784-1:2006]

511 **3.18**

512 **dissemination information package**

513 **DIP**

514 information package, derived from one or more AIPs, sent by an Archive to a Consumer in response to a
515 request in the OAIS

516 **3.19**

517 **distributable object**

518 component of an EPUB publication that can be reused in other contexts

519 Note 1 to entry: A Distributable Object can be a complete EPUB Content Document (e.g., a chapter of a book), a
520 section of such a document (e.g., an exercise or a promotional excerpt), a media resource (e.g., a video or
521 interactive feature), or a combination of such resources that are not necessarily contiguous within the parent
522 EPUB publication but are intended to be able to be distributed as a unit.

523 [SOURCE: EPUB Distributable Objects 1.0]

524 **3.20**

525 **DRM**

526 **digital rights management**

527 packaging, distributing, controlling, and tracking content based on rights and licensing information

528 [SOURCE: ISO 19153:2014]

529 **3.21**

530 **electronic book**

531 **e-book**

532 non-serial digital document, licensed or not, where searchable text is prevalent, and which can be seen
533 in analogy to a print book

534 Note 1 to entry: The use of e-books is, in many cases, dependent on a dedicated device and/or a special reader or
535 viewing software.

536 [SOURCE: ISO 2789:2013]

537 **3.22**

538 **EPUB container**

539 ZIP based packaging and distribution format for EPUB publications

540 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

541 **3.23**

542 **EPUB content document**

543 publication resource that conforms to one of the EPUB content document definitions

544 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

545 **3.24**

546 **EPUB navigation document**

547 specialization of the XHTML content document, containing human- and machine-readable global
548 navigation information

549 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

550 **3.25**

551 **EPUB publication**

552 collection of one or more renditions conforming to the EPUB specifications, packaged in an EPUB
553 container

554 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

555 **3.26**

556 **EPUB reading system**

557 system that processes EPUB publications for presentation to a user in a manner compliant with EPUB
558 specifications

559 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

560 **3.27**

561 **fallback**

562 mechanism with which versions of the same resource in different file formats can be linked to one
563 another

564 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

565 Note 1 to entry: A reading system that does not support the file format of a foreign resource shall traverse the
566 fallback chain until it finds a version it can render.

567 **3.28**

568 **fixity information**

569 information that documents the authentication mechanisms and provides authentication keys to ensure
570 that the Content Information object has not been altered in an undocumented manner

571 [SOURCE: ISO 13527:2010]

572 **3.29**

573 **foreign resource**

574 publication resource that is not a core media type

575 [SOURCE: EPUB Publications 3.0 Recommended Specification 11 October 2011]

576 Note 1 to entry: According to EPUB 3.1, foreign resources require at least one fallback if they are in the spine or
577 embedded in EPUB Content Documents.

578 **3.30**

579 **identifier**

580 data string or pointer that establishes the identity of an item, institution, or person alone or in
581 combination with other elements.

582 [SOURCE: ISO 8459:2009]

583 Note 1 to entry: EPUB 3 specifies Unique Identifiers and Release Identifiers; the latter is a combination of a Unique
584 Identifier and the last modification data of the rendition of the resource.

585 **3.31**

586 **independently understandable**

587 characteristic of information that is sufficiently complete to allow it to be interpreted, understood, and
588 used by the Designated Community without having to resort to special resources not widely available,
589 including named individuals

590 **3.32**

591 **information**

592 any type of knowledge that can be exchanged

593 Note 1 to entry: In an exchange, this is represented by data

594 **EXAMPLE** a string of bits (the data) accompanied by a description on how to interpret the string of
595 bits as numbers representing temperature observations measured in degrees Celsius (the
596 representation information)

597 **3.33**

598 **information package**

599 logical container composed of optional content information and optional associated preservation
600 description information

601 **3.34**

602 **ingest functional entity**

603 OAIS functional entity that contains the services and functions that accept SIPs from producers,
604 prepares AIPs for storage, and ensures AIPs and their supporting descriptive information become
605 established within the OAIS

606 **3.35**

607 **long-term**

608 period of time long enough to raise concerns about the impact of changing technologies, including
609 support for new media and data formats, and of a changing designated community, on the information
610 being held in an OAIS

611 Note 1 to entry: This period extends into the indefinite future.

612 **3.36**

613 **long-term preservation**

614 act of maintaining information, independently understandable by a designated community, with
615 evidence supporting its authenticity over the long term

616 **3.37**

617 **manifest**

618 EPUB manifest element provides an exhaustive list of the Publication Resources that constitute the
619 given Rendition, each represented by an item element.

620 [SOURCE: EPUB Publications 3.0.1]

621 **3.38 metadata**

622 data about other data, documents, or records that describe their content, context, structure, format,
623 provenance, and/or rights.

624 [SOURCE: ISO 5127:2017]

625 **3.39**

626 **METS**

627 Metadata Encoding and Transmission Standard, a standard for presenting metadata using XML.

628 [SOURCE: Digital preservation handbook, Glossary]

629 **3.40**

630 **migration**

631 means of overcoming technological obsolescence by transferring digital resources from one
632 hardware/software generation to the next

633 Note 1 to entry: The purpose of migration is to preserve the intellectual content of digital objects and to retain the
634 ability for clients to retrieve, display, and otherwise use them in the face of constantly changing technology.

635 Note 2 to entry: Migration differs from the refreshing of storage media in that it is not always possible to make an
636 exact digital copy or replicate original features and appearance and still maintain the compatibility of the resource
637 with the new generation of technology.

638 [SOURCE: Digital preservation handbook, Glossary]

639 **3.41**

640 **Open Archival Information System**

641 **OAIS**

642 archive, consisting of an organization, which may be a part of a larger organization, of people and
643 systems, that has accepted the responsibility to preserve Information and make it available to a
644 Designated Community. It has a set of responsibilities, as defined in section 4, which allow an OAIS
645 Archive to be distinguished from other uses of the term 'Archive'.

646 Note 1 to entry: The term 'Open' in OAIS is used to imply that this Recommendation and future related
647 Recommendations and standards are developed in open forums, but it does not imply access to the Archive is
648 unrestricted.

649 Note 2 to entry: The OAIS abbreviation is also commonly used to refer to the Open Archival Information System
650 Reference Model standard which defined the term. The standard is a conceptual framework describing the
651 environment, functional components, and information objects associated with a system responsible for long-term
652 preservation.

653 **3.42**
654 **package document**
655 publication resource that describes one rendition of an EPUB publication, as defined in package
656 document. The package document carries meta information about the Rendition, provides a manifest of
657 resources and defines the default reading order.

658 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

659 Note 1 to entry: It specifies all tools required to render the document, provides an exhaustive list of resources
660 belonging to the document, and defines their default reading order.

661 **3.43**
662 **PDF**
663 Portable Document Format, a set of formats and open standards maintained by the International
664 Organization for Standardization for producing and sharing electronic documents

665 Note 1 to entry: Originally developed by Adobe Systems.

666 [SOURCE: Digital preservation handbook, Glossary]

667 **3.44**
668 **PDF/A**
669 versions of the PDF standard intended for archival use

670 [SOURCE: Digital preservation handbook, Glossary]

671 **3.45**
672 **pre-ingest**
673 actions required before data can be submitted into an OAIS archive, including negotiation of data
674 acquisitions, checking rights and access criteria, licensing, and data submission

675 Note 1 to entry: This area also includes activities involving data producer support and training.

676 Note 2 to entry: Pre-ingest is not a function in the standard OAIS model, but activities in this area can form a
677 significant part of a producer's responsibilities.

678 [SOURCE: UK Data Archive. Archive training manual²⁹]

679 **3.46**
680 **preservation description information**
681 **PDI**
682 information necessary for the adequate preservation of Content Information that can be categorized as
683 provenance, reference, fixity, context, and rights information

684 **3.47**
685 **preservation metadata**
686 metadata containing information needed to archive and preserve a resource

687 [SOURCE: Understanding metadata]

688 **3.48**
689 **preservation planning functional entity**
690 OAIS functional entity that provides the services and functions for monitoring the environment of the
691 OAIS and that provides recommendations and preservation plans to ensure information stored in the

²⁹ <http://www.data-archive.ac.uk/curate/archive-training-manual/pre-ingest>

692 OAIS remains accessible to, understandable by, and sufficiently usable by the designated community
693 over the long term, even if the original computing environment becomes obsolete

694 **3.49**

695 **producer**

696 role played by those persons or client systems that provide the information to be preserved

697 Note 1 to entry: This can include other OAISs or internal OAIS persons or systems. The producer does not need to
698 be the publisher.

699 **3.50**

700 **provenance information**

701 information that documents the history of the Content Information

702 Note 1 to entry: This information states the origin or source of the Content Information, any changes that may
703 have taken place since it was generated, and who has had custody of it.

704 Note 2 to entry: The Archive is responsible for creating and preserving Provenance Information from the point of
705 ingest; however, earlier Provenance Information should be provided by the Producer. Provenance Information
706 adds to the evidence to support authenticity.

707 **3.51**

708 **publication resource**

709 resource that has the content or instructions contributing to the logic and rendering of at least one
710 rendition of an EPUB publication

711 EXAMPLE Examples of publication resources include a rendition's Package Document, EPUB
712 Content Document, EPUB style sheets, audio, video, images, and embedded fonts and
713 scripts.

714 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

715 **3.52**

716 **reading system**

717 system that processes EPUB publications for presentation to a user in a manner conformant with EPUB
718 specification

719 [SOURCE: Modified from EPUB 3.1 Recommended Specification 5 January 2017].

720 **3.53**

721 **reference information**

722 information that is used as an Identifier for the Content Information

723 Note 1 to entry: This also includes Identifiers that allow outside systems to refer unambiguously to a particular
724 Content Information.

725 EXAMPLE an ISBN is a type of Reference Information.

726 **3.54**

727 **reference model**

728 framework for understanding significant relationships among entities in an environment and for the
729 development of consistent standards or specifications supporting that environment

730 Note 1 to entry: A Reference Model is based on a small number of unifying concepts and may be used as a basis for
731 education and explaining standards to a non-specialist.

- 732 **3.55**
733 **reformatting**
734 copying information content from one storage medium to a different storage medium (media
735 reformatting) or converting from one file format to a different file format (file reformatting)
- 736 [SOURCE: Digital preservation handbook, Glossary]
- 737 **3.56**
738 **refreshing**
739 copying information content from one storage media to the same storage media
- 740 [SOURCE: Digital preservation handbook, Glossary]
- 741 **3.57**
742 **release identifier**
743 identifier that allows any instance of an EPUB publication to be compared against another to determine
744 if they are identical, different versions, or unrelated
- 745 Note 1 to entry: Release Identifiers consist of a unique identifier and the last-modified date of the document.
- 746 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]
- 747 **3.58**
748 **remotely-hosted resource**
749 objects hosted outside the EPUB Container.
- 750 Note 1 to entry: EPUB 3.1 allows fonts and resources used by scripts to be hosted externally.
- 751 **3.59**
752 **rendition**
753 one rendering of the content of an EPUB publication, as expressed by an EPUB package
- 754 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]
- 755 **3.60**
756 **repository system**
757 long term preservation system used by an archive
- 758
- 759 **3.61**
760 **rights management metadata**
761 information that identifies the access restrictions concerning the Content Information, including the
762 legal framework, licensing terms, and access control
- 763 Note 1 to entry: This contains the access and distribution conditions stated in the Submission Agreement, related
764 to both preservation (by the OAIS) and final usage (by the Consumer).
- 765 Note 2 to entry: It also includes specifications for the application of rights enforcement measures.
- 766 **3.62**
767 **spine**
768 EPUB spine element defines the default reading order of the EPUB Publication content by defining an
769 ordered list of manifest item references.
- 770
- 771 [SOURCE : EPUB Publications 3.0.1]

772 **3.63**
 773 **structural metadata**
 774 metadata that indicates how compound objects are put together, for example how the pages of a
 775 document are arranged to form chapters

776 [SOURCE: Understanding metadata]

777 **3.64**
 778 **submission agreement**
 779 agreement reached between an OAIS archive and a Producer that specifies a data model and any other
 780 arrangements needed for the data submission session

781 Note 1 to entry: This data model identifies the format/content and the logical constructs used by the Producer and
 782 how they are represented on each media delivery or in a telecommunication session.

783 **3.65**
 784 **submission information package**
 785 **SIP**
 786 information package that is delivered by a Producer to an OAIS to be used to construct or update one or
 787 more AIPs and/or the associated descriptive information.

788 **3.66**
 789 **unique identifier**
 790 primary identifier of an EPUB publication, which may be shared by one or several renditions of the
 791 same EPUB publication that conform to the EPUB standard and embody the same content.

792 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

793 **3.67**
 794 **XHTML content document**
 795 EPUB content document that conforms to the profile for HTML defined in XHTML Content Documents

796 [SOURCE: EPUB 3.1 Recommended Specification 5 January 2017]

797 Note 1 to entry: see EPUB Content Documents 3.1, chapter 2.

798 **4 Abbreviated terms**

AIP	Archival Information Package
DIP	Dissemination Information Package
DRM	Digital Rights Management
OAIS	Open Archival Information System
PDI	Preservation Description Information
SIP	Submission Information Package

799 **5 Packaging standards**

800 An archiving process includes several distinct steps. A producer – which may be the publisher or other
 801 body acting on behalf of the publisher, such as the archive itself - creates a Submission Information
 802 Package (SIP) and transfers it to a repository system in an OAIS archive. The archive performs a quality

803 control process to the SIP and, if the package meets the criteria set in the submission agreement,
804 accepts it, creates an Archival Information Package (AIP) and transfers the package to archival storage.
805 During ingest some of the files or metadata records within SIP may be migrated to new formats or
806 additional metadata may be added.

807 The OAIS archival storage function stores, maintains, and retrieves AIPs. Maintenance may include for
808 instance frequent error checks to protect the data against bit rot. In order to keep the documents
809 understandable it may also necessary be necessary to migrate³⁰ them in new formats, or to update the
810 AIP with additional metadata. Migration and other preservation related tasks may be carried out by the
811 producer, OAIS archive and / or third parties. The party or parties responsible should be specified in
812 the submission agreement.

813 The OAIS Access function allows users to retrieve information from a repository system in the form of
814 Dissemination Information Packages (DIPs) which can include all or parts of the data and metadata of
815 an AIP. Differences between SIPS, AIPs and DIPs can be substantial, depending on the preserved
816 content, requirements of submission agreement, national legislation and institutional practices. OAIS
817 does not require a 1:1 relationship between information packages, so one AIP can contain documents
818 and metadata from multiple SIPS or vice versa.

819 Transfers of package states (SIP to AIP to DIP) do NOT mean that the content SHALL change. The
820 change from SIP to AIP can be minimal, that is, the content information remains the same, but some
821 administrative metadata is added into the AIP about the actions taken during the ingest process. If an
822 EPUB publication is created according to the requirements in this document there should be no need for
823 reformatting the EPUB publication itself. During ingest it is enough to check the validity of the
824 document, and if there are no issues, it can be stored "as is". Some archives may choose to apply even
825 simpler initial ingest procedures (that is, avoid even validity checks) if the producer is well known and
826 reliable, such as other OAIS archive.

827 This specification covers only the initial stage of the archiving process, namely the creation of
828 Submission Information Package (SIP). SIP consists of data objects and representation information with
829 which the data is interpreted. Both the data (documents) and representation information (metadata)
830 MUST conform to the standards and specifications the producer and the archive have agreed upon in
831 the submission agreement. If a SIP does not meet the requirements, ingest to the repository system
832 fails. Note that a SIP MAY contain unarchivable resources, provided that they have been encoded in an
833 appropriate manner.

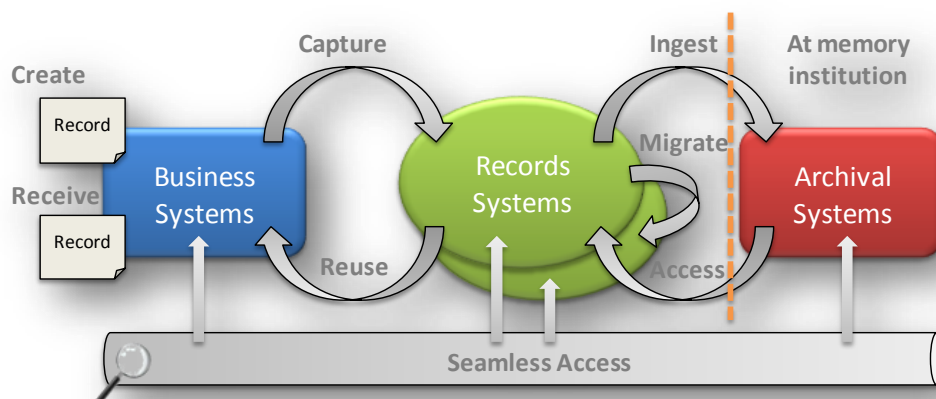
834 The content and structure of all information packages in repository systems MUST be standardized.
835 There are several packaging standards available, but the most commonly used one is the Metadata
836 Encoding and Transmission Standard (METS³¹) developed by the Library of Congress. ISO/IEC
837 JTC1/SC34 JWG 7 decided to recommend the use of METS as the container standard, although this
838 specification does allow the use of other container standards as well.

839 Since container standards – including METS - are rich specifications there is a need to create profiles to
840 specify how they should be used. This specification provides a METS profile for EPUB in Part 2. Other
841 container standards are not taken into account; if METS is replaced by another container specification,
842 profiling needs to be done separately.

843 Some digital preservation projects have developed tools for creating SIPs that meet the project
844 requirements, which makes it a lot easier to submit information to the repository system. Producers
845 SHOULD nevertheless have at least basic understanding of digital preservation, since all pre-ingest
846 steps from document creation to SIP submission should be carried out in such a way that the
847 authenticity of submitted documents can be guaranteed.

³⁰ From OAIS point of view, migration is a complex process which involves export of the document (as a migration DIP) and then migration during "ingest as new manifestation".

³¹ <http://www.loc.gov/standards/mets/>



848
849 **Figure 2. Information flow between live and archival systems**
850 **[E-ARK Common specification, p. 13]**

851 Different disciplines, even if they all use OAIS, will develop interfaces optimized for their own needs.
852 And if the payloads are not the same, technical metadata standards will also differ. Domains may even
853 adopt different packaging and preservation metadata standards. For instance, almost all digital
854 archiving projects in the library domain rely on METS and PREMIS specifications, although some
855 libraries use BagIt³² as an alternative for METS. Compared with libraries, the film industry started
856 digital preservation efforts a bit later and may eventually develop different preferences³³. And even if
857 the same standards were used, they may be applied in a non-interoperable way even within the same
858 domain. Therefore creating set application profiles is very important in digital archiving.

859 **6 Construction of OAIS information packages**

860 According to the Open Archival Information System (OAIS) model³⁴, information package is “a container
861 that contains two types of Information Objects, the Content Information and the Preservation
862 Description Information (PDI)”. Content information is the data that needs to be preserved and
863 preservation description information is the metadata and other information that is needed in order to
864 preserve, find and understand the data in long-term.

865 Preservation description information consists of reference information, provenance information,
866 context information, fixity information, and access rights information. See the OAIS specification for an
867 in depth explanation of these.

868 According to the OAIS specification (pages 4-35),

869 *[i]t is necessary to distinguish between an Information Package that is preserved by an OAIS*
870 *and the Information Packages that are submitted to, and disseminated from, an OAIS. These*
871 *variant packages are needed to reflect the reality that some submissions to an OAIS will*
872 *have insufficient PDI to meet final OAIS preservation requirements. In addition, they MAY be*
873 *organized very differently from the way the OAIS organizes the information it is preserving.*
874 *Finally, the OAIS MAY provide information to Consumers that does not include all the PDI*
875 *with the associated Content Information being disseminated. These variants are referred to*
876 *as the Submission Information Package (SIP), the Archival Information Package (AIP), and*
877 *the Dissemination Information Package (DIP). Although these are all Information Packages,*

³² <https://en.wikipedia.org/wiki/BagIt>

³³ https://www.cen.eu/news/calls/Calls/CEN-Call_for-tender_Digitalcinema.pdf

³⁴ <http://public.ccsds.org/publications/archive/650x0m2.pdf>

878 *they differ in mandatory content and the multiplicity of the associations among contained*
879 *classes.*

880 The principles listed below provide SIP production guidelines for document producers (publishers or
881 third parties creating EPUB publications). The creation of the principles has been inspired by the draft
882 common requirements published by the E-ARK project (see Introduction to the Common Specification
883 for Information Packages in the E-ARK project, version 1.0³⁵). Although E-ARK has served as a model for
884 this specification, these requirements have not been aligned with those of E-ARK, and therefore there
885 may be significant differences between the specifications.

886 6.1 General

887 6.1.1 EPUB publications SHALL be sent to a repository system as well-formed and complete 888 Submission Information Packages (SIPs)

- 889 • This specification does not assume that publishers create SIPs. The OAIS producer MAY be a
890 third party acting on behalf of the publisher, such as OAIS archive.
- 891 • This specification and its accompanying document are mainly concerned with the structure
892 and content of SIPs. The way EPUB publications are archived and disseminated (the
893 structure of Archival Information Packages and Dissemination Information Packages, or
894 AIPs and DIPs) depends on the submission agreements made between the archive and the
895 producers, and on the operational principles of the archive, and is beyond the scope of this
896 document. It is possible that an EPUB publication is migrated into another format during
897 Ingest, and disseminated again as an EPUB publication. The archive may also preserve (in
898 bit level) the original file.
- 899 • Submitted EPUB publications SHALL be conformant with EPUB requirements³⁶ and
900 conformance SHOULD be validated.
- 901 • Submitted EPUB publications SHALL either contain or at least facilitate access to all the data
902 and metadata required to render the content information successfully.
 - 903 i. Preview publications MAY be submitted, even though they are by definition not
904 complete, if the final documents are sent when completed. Depending on the
905 submission agreement, the archive MAY preserve just the final version, or both
906 versions of the resource.
 - 907 ii. Distributable objects SHALL NOT be submitted individually. They MAY be embedded
908 within EPUB publications, but the archive is not obliged to deliver them as DIPs unless
909 the submission agreement mandates that.
 - 910 iii. Fonts SHALL be embedded into the EPUB publication in full and un-obfuscated, if font
911 license allows that.
 - 912 iv. Related resources such as audio and video SHOULD be embedded in the EPUB
913 publication.
 - 914 v. Remotely-hosted resources SHOULD be avoided, but if used, it is necessary to ensure
915 that all remote data is available to the archive so that the data can be incorporated into
916 the AIP during ingest, and permission to do this SHALL be explicitly agreed upon in the
917 submission agreement, especially if the publisher is not in full control of remote data.
 - 918 vi. Descriptive and other metadata SHOULD be embedded in the SIP. METS mdRef element
919 MAY only be used if a) referred metadata is part of the same SIP, or b) the archive is
920 able to retrieve any linked external metadata and incorporate it into the AIPs in an
921 appropriate format.

³⁵ <http://www.dasboard.eu/specifications/common-specification>

³⁶ Conformance requirements for EPUB publications and reading systems have been specified in chapter 3.1 of EPUB Recommended specification, version 3.1.

- 922 vii. Permission to use remote resources and metadata SHALL be specified in the
 923 submission agreement³⁷. The permission SHALL specify acceptable metadata and file
 924 formats.
- 925 • The SIP SHOULD³⁸ be checked for viruses and malicious software before submission to the
 926 repository system.
 - 927 • EPUB publications in SIPs SHOULD NOT be encrypted, because that compromises long-term
 928 preservation. If data is submitted in an encrypted format, the archive SHALL receive
 929 necessary decryption information/details within the SIP, as agreed in the submission
 930 agreement or elsewhere. When the archive disseminates the archived data to its customers,
 931 it can be encrypted again.
 - 932 • DRM protection, if any, SHOULD be removed by the producer before the document is
 933 submitted. If the content in the SIP is DRM protected, the archive SHALL receive the
 934 necessary information/details to remove the DRM protection within the SIP, as agreed in
 935 the submission agreement or elsewhere. Such permission may be producer-specific, based
 936 on the submission agreement, or a generic permission, based on e.g. the Copyright Act.
 - 937 • If data is compressed, the user of the compression method SHALL be specified using the
 938 Compression metadata element in the EPUB's encryption.xml file.
 - 939 • The submission agreement SHOULD specify at least one EPUB reading system capable of
 940 rendering the submitted EPUB publications successfully. Knowing the reading system
 941 requirements in advance makes it easier for the archive to design and implement the ingest
 942 process. Although submitted publications will usually be validated only with automated
 943 tools³⁹, the archive should be able to validate that the received EPUB can be presented to
 944 the customers, and check for instance the look and feel of archived EPUB publications
 945 before and after migration. This is possible only if the archive can operate the reading
 946 systems that can render the archived publications successfully.
 - 947 • Each SIP SHOULD specify EPUB reading system or systems, which can render the EPUB
 948 publication in the SIP. If this information is missing, reading system or systems SHALL be
 949 specified in the submission agreement.
 - 950 • Multiple-rendition EPUB publications may be designed for multiple reading systems, in
 951 which case the submission agreement may require the archive to carry out at least
 952 occasional checks in all of these reading systems. If so, all these reading systems SHOULD be
 953 listed in the submission agreement.
 - 954 • If a submitted EPUB publication has been optimized for a certain reading system, the
 955 system SHOULD be described in the document's technical and/or preservation metadata,
 956 since such information is valuable for preservation and archival access purposes.
 - 957 • If the optimal EPUB reading system is no longer available, the archive SHOULD, with
 958 permission and support from the producer, either find another suitable reading system or
 959 modify the ingest process so that the EPUB publications affected by this change can be used
 960 by another EPUB reading system.⁴⁰

³⁷ If there are remote resources or associated metadata linked to the SIP with a LINK element, these external resources will be retrieved as part of the ingest process and included in the AIP. If external resources cannot be retrieved, the ingest process fails. The producer SHALL send either a new SIP with all the data and metadata embedded into it, or make sure that the archive is allowed to access remote data and metadata.

³⁸ Some producers may not be able to make virus checks, but all OAIS archives SHALL be. Virus checks are commonly done during ingest.

³⁹ One such tool is Epubcheck, available from <https://github.com/idpf/epubcheck>

⁴⁰ While this standard is about the "state" in which the EPUB publication itself shall be in order to be archivable, the SIP may include a lot of other information (metadata, executables, other renditions of the EPUB publication, additional documentation etc) which may make it easier to preserve the intellectual content in the long term.

961 **6.1.2 Regardless of its type or format, it SHALL be possible to include any data or metadata in**
 962 **SIPs**

- 963 • It SHOULD be possible to maintain the SIP and EPUB specifications independently, i.e. so
 964 that any change to SIP does not automatically mean that the EPUB format needs to be
 965 updated and vice versa. The exception from this rule is that any existing and future features
 966 in EPUB specification which are relevant from long-term preservation point of view such as
 967 font embedding SHALL be taken into account in the SIP specification.
- 968 • This document does not set a priori constraints either to the current or future versions of
 969 EPUB with regard to the choice of metadata and file formats or either's versions (see note 1
 970 below on EPUB Core media types).
- 971 • The submission agreement SHOULD specify metadata formats and file formats approved for
 972 submission and archival. For EPUB publications, at least Dublin Core metadata format and
 973 all EPUB core media types SHALL be supported by the archive in order to guarantee
 974 efficient processing of EPUB publications.
- 975 • Submission agreements may specify what kind of executables may be embedded in the
 976 submitted EPUB publications, or forbid their use entirely (see note 2 below on interactive e-
 977 books and EPUB publications).

978 NOTE 1 EPUB community may change the list of EPUB Core Media Types any time, independent of
 979 the EPUB specification updates. New core media types may be approved and old ones
 980 deprecated. If core media types are not checked from long-term preservation point of view,
 981 some new EPUB core media types may turn out to be non-archivable.

982 File format lists in submission agreements may cover all EPUB core media types or – if the
 983 producer does not use all the core media types – just a subset of them. When a core media
 984 type is deprecated, the producer (if it still exists) and the archive should decide whether the
 985 file format in question is migrated or kept as is (and emulated). If the latter, it may be
 986 necessary to migrate the deprecated file format when DIP packages are created.

987 NOTE 2 E-books are likely to become more interactive in the future, which is why there are various
 988 ways EPUB 3 can support interactivity. However, some EPUB reading systems may not
 989 support interactivity, and even if it is supported, different reading systems may not behave
 990 identically, partly because EPUB is not specific about how support should be implemented.
 991 EPUB 3 `object` element enables the use of arbitrary embedded executables that are not
 992 inherently supported in EPUB 3 reading systems. A common use case would be to include
 993 proprietary applets or Adobe Flash applications. However, in a majority of cases, interactive
 994 publications will be created through the use of in-book source code. Because JavaScript is
 995 the de facto standard scripting language for SVG and HTML5, EPUB 3 content documents
 996 can be assumed to be scriptable only if they contain JavaScript code. The standard does not
 997 define which versions of JavaScript (ECMAScript) are required to get the support. Content
 998 creators should comply with the most commonly supported features in web browsers for
 999 best results [Daly]. Whatever the chosen approach, interactivity will be difficult to preserve
 1000 in the long-term, since applications should be adapted to new hardware and software
 1001 environments, which can be difficult if not impossible.

- 1002 • Archives offering long-term preservation services for EPUB publications SHOULD keep
 1003 track of EPUB core media types and consider the possibility of including them on the list of
 1004 archivable formats. If this is not viable, the archives SHOULD maintain clearly defined and
 1005 well tested migration pathways from non-archivable core media types into archivable

formats. Then the archive would not need to migrate these images during ingest and it would be possible to preserve EPUB publications unchanged⁴¹.

- If there is a foreign resource embedded or linked to a submitted EPUB publication, a fallback chain ending in a core media type resource SHOULD be provided even if the foreign resource is in an archivable format. (Note that this requirement is stricter than those in EPUB 3.x specifications, which require a fallback only in certain situations.)
- The producer MAY include foreign resources (and metadata formats) in submitted EPUB publications if they have been specified as suitable for ingest and/or archivable in the submission agreement, or if they are encoded in SIPs in such a way that they will be ignored during ingest (see below).
- If foreign resources and metadata are originally in un-archivable formats, they SHALL be migrated when the SIP is formed. The AIP may contain either just migrated publications, or both the original and migrated publications.
- Core media types and foreign resources not specified in the submission agreement MAY be submitted if and only if the submission agreement allows it. These files SHALL be encoded in the SIP in such a way that they are not validated against the generic ingest criteria during the ingest process (since otherwise the SIP shall not pass the validation) and therefore passed directly to AIP. The specifics of this type of encoding SHALL be defined in the submission agreement.
- If there are alternative versions (renderings) of the publication to be included in the SIP which are not archivable, they SHOULD be migrated into acceptable file formats prior to submission by the producer or a third-party preparing a SIP on behalf of the producer. For instance, if PDF is specified as not archivable but PDF/A is, the producer should create a PDF/A version of the document, which will then be submitted to the repository system alongside the EPUB publication of the same work.
- If these non-archivable originals are included in the SIP, they SHALL be encoded in such a way that they are not validated against the generic ingest criteria during the ingest process (since otherwise the SIP would not pass) and therefore passed directly to AIP. The specifics of this type of encoding SHALL be defined in the submission agreement⁴².

NOTE EPUB 3 Fixed Layout Properties

In digital preservation the usual aim is to preserve intellectual content, not the original look and feel of the document, because trying to preserve the original layout for decades and even centuries may either be difficult or impossible. EPUB publications are generally designed so that their look and feel can change with no impact on semantics, which is a good thing from the digital preservation point of view. EPUB content presentation adapts to the user preferences and display properties, which are both likely to change radically in the future.

Fixed layout EPUB publications are an exception. In them, the intellectual content and the design of the document cannot be separated: any change in the appearance of the document may cause significant changes in the meaning or even lose it completely. Therefore fixed-layout EPUB publications give the content creators greater control over presentation. This control is based on a set of metadata properties with which the

⁴¹ An OAIS archive does not need to migrate non-archivable file formats during the ingest process. Depending on the preservation strategy, migration may only happen when a real risk to the format emerges – such as the loss of applications capable of rendering it - or when the document is disseminated for the first time.

⁴² Ideally, a well-designed and built repository system should be able to validate any file format. In practice, there are file formats validation tools cannot process. If there is a need to preserve these files in bit-level, they have to be ignored during validation.

intended rendering dimensions can be specified [EPUB 3 Fixed]. However, if the document is migrated, these metadata properties may be lost, and even if that does not happen changes in hardware (e.g. display technologies), operating systems, and middleware may change the original look and feel of the document.

Submission agreements SHOULD specify if submission of fixed-layout EPUB publications is allowed and if so, how they are treated during ingest. The best solution is to encode them in such a way that they are not validated, and include in SIPs also reflowable versions of publications. If the producer is not able to create such versions, submission agreement MAY require the archive or a third party to do that as a part of the ingest process. Submission agreements may also allow submission of fixed layout EPUB publications, with mutual understanding that preserving semantics which is tied to the original look and feel will in the long term be impossible.

6.1.3 It SHOULD be possible to transfer SIPs by any means, methods, or tools from the submitting organization to the repository system

- Although there are no general limitations (it is possible to use e.g. FTP or UPS), submission agreements MAY limit the options available by specifying the protocols to be used during submission.
- SIPs SHALL be composed so that their structure and content does not limit the use of any particular transfer method.

6.1.4 The archive SHALL have a way to verify the identity of the submitting organization/person, no matter how the information packages are transferred

- If submission is taken care of by a third party service and the producer is a different organization of person, the archive SHALL be able to verify the identity of both of them.
- There are various ways to implement this requirement, including digital signatures, secure channels, recording relevant information within the SIP as metadata, or even manual exchange of data on secure media.
- Part 2 of this specification provides an example of how a digital signature can be used for verification.

6.1.5 There is no 1:1 relation between OAIS information packages

- SIPs SHALL be composed so that their structure and content SHALL NOT prescribe or limit SIP -> AIP -> DIP conversions.
- During ingest, it SHALL be possible to transform one SIP into 1-n AIPs, or many SIPs into 1-n AIPs. For instance, a SIP might consist of all yearbooks of a publisher (e.g. 15 EPUBs) which are then archived in separate AIPs. Relevant data and metadata SHALL always be archived; number of AIPs created during ingest depends on the internal practices and processes of the archive, which are not within the scope of this specification.

6.1.6 A SIP MAY contain 0-n EPUB 3 publications, and one EPUB 3 publication MAY be submitted to the repository system in 1-n SIPs

- A SIP MAY contain only metadata about EPUB publication, not the publication itself.

- 1091 • A SIP SHOULD contain multiple EPUB publications only if they are interrelated; for instance,
1092 different renderings of the same document⁴³.
- 1093 • A SIP MAY contain alternative renderings (such as PDF or DOCX) of the publication, but the
1094 SIP SHALL contain all administrative metadata required for processing of these versions,
1095 and explaining the relations between these renderings.
- 1096 • A single EPUB publication MAY be split into multiple SIPs if there is a valid reason to do so,
1097 such as the complexity or large size of the document.

1098 **6.1.7 The information package type (in this case, SIP) SHALL be indicated**

- 1099 • Only packages which are marked to be SIPs will be ingested. AIPs, DIPs and unlabeled
1100 packages are not suitable for ingest.

1101 **6.1.8 SIP packaging method SHALL not restrict the application of any preservation method**

- 1102 • Although the most common preservation method is migration, some archives MAY choose
1103 emulation as the primary approach, which will have an impact on the OAIS Preservation
1104 Description Information required.
- 1105 • Some information objects (such as programs) are not suitable for migration. Submission
1106 agreements SHOULD specify a preservation strategy for such resources.

1107 **6.1.9 The packaging method SHALL NOT limit the size of the SIP**

- 1108 • Some archives can have problems in e.g. validating and ingesting very large data objects. If
1109 there is a risk that the SIPs are becoming too large for the submission method used or the
1110 ingest process used by the archive, an appropriate splitting mechanism SHOULD be applied.
1111 Describing such mechanisms is beyond the scope of this specification.

1112 **6.2 Identification of information packages and their content**

1113 **6.2.1 It SHALL be possible to identify any SIP uniquely during the ingest process**

- 1114 • Since multiple SIPs may be submitted to the repository system simultaneously, there is a
1115 need to identify all packages in a (globally) unique manner. Identification will also make it
1116 possible to relate the packages with appropriate submitters, earlier submissions etc. Such
1117 identification helps to streamline the whole submission process and any potential
1118 communication between the archive and the submitting organization.
- 1119 • Once the ingest process has been completed and 1-n AIPs have been formed, the SIP is no
1120 longer needed. The producer receives persistent identifier(s) for the AIP(s) containing the
1121 submitted data. The SIP identifier is no longer needed after this.
- 1122 • There are circumstances in which AIP identifiers SHOULD be not only persistent, but also
1123 globally unique. For instance, an OAIS archive can cooperate with other archives by
1124 exchanging AIPs in order to share the bit level preservation costs.
- 1125 • The entire SIP or parts of it SHALL be resubmitted in a revised format if the ingest process
1126 fails due to errors in the package. To keep track of the packages, SIPs SHALL have unique
1127 identifiers.

⁴³ OAIS archives may have different ideas of what “interrelated” means. For instance, archives tend to prefer large SIPs which may contain large number of documents gathered for years, while libraries archive publications on an individual basis.

1128 **6.2.2 Information objects (EPUB publications, PREMIS preservation metadata record, etc.)**
 1129 **within SIPs SHALL be identified uniquely and persistently**

- 1130 • Identifiers have many vital uses in digital preservation. They are used as access keys to the
 1131 archived content in repository systems and facilitate information exchange with external
 1132 systems. Identifiers also enable linking different versions of an archived document to each
 1133 other. Moreover, with identifiers it is possible to link documents and
 1134 descriptive/administrative metadata records that describe them. These links enable the
 1135 archive to e.g. create dissemination information packages with the requested content.
- 1136 • Submission agreements SHALL specify identifier systems used, their location (EPUB
 1137 document or SIP) and who is responsible of creating them (producer, archive or a third
 1138 party). For instance, if the use of EPUB release identifiers is forbidden because the
 1139 repository system does not support them, another means of identifying releases is needed.
- 1140 • International standard identifiers, such as ISBNs for books and DOIs for articles, SHALL be
 1141 used as EPUB unique identifiers whenever possible. Any exceptions (such as using other
 1142 identifier systems for releases which do not have ISBNs) SHOULD be specified in the
 1143 submission agreement.
- 1144 • It SHOULD be possible to express the identifiers (also) as actionable HTTP URIs. Usage of
 1145 persistent identifiers (Handles, DOIs, URNs, or ARKs) is recommended.
- 1146 • If there are multiple renditions of a work in an EPUB publication, requirements in the EPUB
 1147 Multiple-Rendition Publications 1.0 specification SHALL be followed. Each rendition of an
 1148 EPUB publication in a SIP SHALL have its own identifier.
- 1149 • The SIP SHOULD contain separate descriptive and administrative metadata records for each
 1150 rendition, and these records SHALL have their own identifiers.

1151 NOTE According to EPUB Multiple-Rendition Publications 1.0, the need to include more than one
 1152 rendition of the content in an EPUB publication has grown as reading systems have become
 1153 more sophisticated. In addition to optimizing the layout, adapting the content to specific
 1154 reading systems may involve changing the content itself. Adaptation may also involve the
 1155 prose of a textual work; instead of publishing several single-language EPUB publications
 1156 multiple translations may be published as a single multiple-rendition EPUB publication.

1157 **6.2.3 EPUB Fragment Identifiers SHOULD not be used in EPUB publications sent to a repository**
 1158 **system, unless the submission agreement explicitly allows their use**

- 1159 • EPUB Canonical Fragment Identifiers define a standardized method of referencing content
 1160 within an EPUB publication through the use of URI Fragments. From the digital
 1161 preservation point of view, fragment identifiers can be problematic if the preservation
 1162 strategy is not emulation, since URI fragments are media type dependent. Following
 1163 migration the fragment identifiers may no longer be functional, because the new media type
 1164 does not support them.
- 1165 • If fragment identifiers are allowed, the producer and the archive SHOULD take this into
 1166 account in preservation planning, and design migrations so that the functionality provided
 1167 by the fragment identifiers is preserved.

1168 **6.3 Structure of information packages**

1169 **6.3.1 Submission information packages SHALL be built in such a way that their components**
 1170 **can be logically and physically separated from one another**

- 1171 • For each rendition of the EPUB content document, there SHALL be a manifest file, which
 1172 identifies and describes a set of resources that collectively compose a given rendition of a
 1173 document, and EPUB spine, which provides a default reading order for a given rendition.

- 1174
- 1175
- 1176
- 1177
- 1178
- 1179
- 1180
- 1181
- 1182
- EPUB Open Container Format (OCF) defines a file format and processing model for encapsulating a set of related resources (for instance, renditions of the same resource) into a single-file (ZIP) EPUB Container⁴⁴.
 - The structure of each EPUB ZIP archive SHALL be described using the EPUB container.xml file (which describes the locations of root files of available renditions of the EPUB publication, and the rendition's package document and navigation document).
 - EPUB Package document and navigation document SHALL contain all metadata needed for rendering the publication, including the recommended reading system.

1183 6.4 Generic Information package metadata

1184 6.4.1 Metadata in information packages SHALL be based on standards

- 1185
- 1186
- 1187
- 1188
- 1189
- 1190
- 1191
- 1192
- 1193
- 1194
- 1195
- 1196
- 1197
- 1198
- 1199
- 1200
- 1201
- 1202
- 1203
- 1204
- 1205
- METS or another agreed upon container format SHALL be used as the container standard, since this makes ingest to existing repository systems easier.
 - The submission agreement SHALL specify at least one mandatory metadata format for descriptive metadata. The format does not need to be Dublin Core; although EPUB publications always contain some Dublin Core metadata elements (see below), they MAY contain more complete metadata in another format, such as ONIX.
 - The minimum required descriptive metadata for EPUB publications are title, identifier, and language from the Dublin Core Metadata Element Set. Each rendition of a publication SHOULD also have at least the last modified date property from the DCMI Metadata Terms. Each rendition SHOULD also have the publication date encoded as DCMI Date, if the publication date is required to distinguish between publications.
 - SIPs submitted to a repository system MAY⁴⁵ contain preservation metadata, although such metadata will normally not be created in production systems, but in repository systems during ingest. PREMIS SHOULD be used for preservation metadata, as it is the most widely used and supported standard for this kind of metadata.
 - The submission agreement SHALL specify the syntax of metadata and its location (in the EPUB document, or in the SIP container), metadata formats used and metadata elements required or recommended.
 - Since problems with text forms and encodings are common in repository systems, text metadata SHOULD be provided in Romanized form, using the EPUB alternate-script property to transcribe it if the metadata is originally in some other script.

1206 6.4.2 Metadata SHOULD allow (automatic) validation of the structure and content of SIPs in 1207 terms of integrity, fixity, and syntax

- 1208
- 1209
- 1210
- SIPs SHALL contain message digests for all files of the SIP, and for the package itself.
 - File format identification and validation metadata (created with EpubCheck⁴⁶ or other validator tool) SHOULD be included in the SIP, if a validator is available.

⁴⁴ EPUB specifications do not require or recommend any specific ZIP tool. It is possible to use for instance ePubPack (<https://sourceforge.net/projects/epubpack/>) to create EPUB ZIP containers from a folder.

⁴⁵ Adding preservation metadata during pre-ingest might be tricky since preservation metadata is the core of any preservation system and its use is highly regulated within repository systems. Errors in preservation metadata prepared by the submitter may cause serious problems in the preservation process.

⁴⁶ As of this writing (2018-01-12) EpubCheck does not support EPUB 3.1.

1211 **6.4.3 It SHALL be possible to edit metadata in information packages**

- 1212
- 1213
- 1214
- 1215
- 1216
- 1217
- If ingest has failed because of erroneous or missing metadata, the producer or a third-party responsible for the submission SHALL be able to modify the SIP so that it meets the metadata requirements in the submission agreement.
 - Producers and archives MAY use crowdsourcing and entity extraction activities to update descriptive metadata; an archive MAY choose to update this metadata also in the AIPs in the repository system although all the other components in the packages remain unchanged.

Annex A (informative)

EPUB and digital preservation: issues and recommendations

1218
1219
1220
1221

1222 The British Library's EPUB Format preservation assessment includes a preservation risk summary
1223 [Whibley, p. 7-8]. The risks mentioned in the BL assessment are marked with [BL].

1224 **A.1 EPUB standard: issues**

- 1225 • Lack of stability in the e-book sector
 - 1226 ○ EPUB does not have universally widespread support across e-book devices [BL].
- 1227 • Lack of EPUB format stability [BL]
 - 1228 ○ Evolving standards, context and the format itself places uncertainty on the future
 - 1229 preservation situation
 - 1230 ○ Proprietary changes and non-standard use of specifications may be used to restrict
 - 1231 access to specific manufacturer hardware/software
- 1232 • Challenging EPUB features
 - 1233 ○ From the long-term preservation point of view, the challenging features in EPUB include
 - 1234 the possibility of using DRM, encryption and obfuscation, foreign resources, non-
 - 1235 embedded resources, interactive documents (containing software components), and
 - 1236 fixed-layout documents.
- 1237 • Lack of archivable EPUB version
 - 1238 ○ The standard is becoming richer and richer, and publishers and other users may find it
 - 1239 more difficult to specify and avoid counterproductive features from the long-term
 - 1240 preservation point of view. Pre-ingest (modifying the EPUB publication so that it can be
 - 1241 preserved easily) may be difficult unless it has been taken into account from the
 - 1242 beginning.

1243 Recommendations:

- 1244 • W3C should actively promote the EPUB format, because it is the only open e-book standard and
- 1245 it is based on open standards such as HTML5 and CSS.
- 1246 • EPUB community and digital preservation experts should develop a subset of EPUB ("EPUB/A")
- 1247 suited for long term preservation.

1248 **A.2 EPUB usage: issues**

- 1249 • Ecosystem specific EPUB implementations
 - 1250 ○ Major players in the e-book market (e.g. Amazon, Apple) have built EPUB based but
 - 1251 closed (non-interoperable) ecosystems for e-books. E-books in vendor-specific formats,
 - 1252 such as Amazon's KF8 should be migrated to EPUB before they are submitted to a
 - 1253 repository system. Technically this is possible since EPUB is a "more or less obvious
 - 1254 superset of what is possible in the different formats". The only exception is the fixed-
 - 1255 layout document specification in KF8; it is based on percentage information, not on
 - 1256 absolute pixel positions as in EPUB 3. [Bläsi, p. 38].
 - 1257 ○ These players have also created vendor-specific DRM solutions, which prevent the use of
 - 1258 archived EPUB publications with other vendor's reading devices, unless the DRM
 - 1259 protection has been removed during pre-ingest or ingest.
- 1260 • Encryption and obfuscation [BL]

- 1261 ○ Encryption may prevent the rendering of documents.
- 1262 ○ Where not easily substituted, obfuscated fonts may lead to loss of critical information.
- 1263 • Incomplete support in EPUB viewers [BL]
- 1264 ○ Support for all aspects of the EPUB standard appears to be mixed, although impact of
1265 this is unclear. In the short-term, if the EPUB publication has been optimized for a
1266 specific reading system or systems, metadata embedded in the SIP should specify these
1267 systems. In the long-term, functionalities that are not widely supported may be lost.
- 1268 • Losing information
- 1269 ○ Where not easily substituted, non-embedded fonts may lead to loss of critical
1270 information.
- 1271 ○ Metadata (and data) may not be embedded, but just linked to the SIP. During ingest,
1272 retrieval of linked information may fail.
- 1273 • Invalid or badly formed EPUB files [BL]
- 1274 ○ May affect the ability to render files now or in the future.
- 1275 • Documents relying on EPUB features that may be difficult to preserve
- 1276 ○ Fixed-layout documents: digital preservation usually concentrates on preserving the
1277 intellectual content, not the original look and feel of the document since that is regarded
1278 as difficult in the long-term. Preserving fixed-layout EPUB publications for the long-term
1279 may therefore be very demanding if not impossible.
- 1280 ○ Interactive documents that contain embedded applications supporting the required
1281 functionality are a challenge, because it may be necessary to modify or rewrite these
1282 applications when hardware or software platforms change or when the documents are
1283 migrated.
- 1284 • Legal issues [BL]
- 1285 ○ It may be illegal to remove DRM, de-obfuscate embedded fonts, or to migrate the
1286 document to some other e-book format.
- 1287 • Interactivity and animations
- 1288 ○ With EPUB 3, there are two possibilities to realize built-in animations and interactive
1289 features. One is to use a CSS construct for transformations; another, more versatile
1290 approach is to use embedded JavaScript, Adobe Flash, or other software components
1291 that may enable complex interactive behaviour [Bläsi, p. 32]. Although EPUB 3 allows
1292 the use of JavaScript, it does not standardize the use of JavaScript elements in e-books.
1293 This can easily lead to proprietary extensions as well as incompatible EPUB 3 reading
1294 systems that support a different or incompatible subset of scripting elements [Bläsi, p.
1295 17]. Maintaining the functionality provided by scripts after migrations and after
1296 hardware and software platform changes may be difficult.
- 1297 ○ Amazon's KF8 does not support interactive features, and iBooks supports them in a
1298 different and undocumented way. Therefore migrating this functionality between
1299 different e-book formats is not possible.
- 1300 • Non-archivable core media types
- 1301 ○ Depending on the chosen preservation strategy, some current or future core media
1302 types may be regarded as unsuitable for digital preservation. For instance, GIF is not an
1303 archivable format according to the requirements of the Finnish National Digital Library
1304 initiative. [File formats, p.25].
- 1305 • Non-archivable foreign resources
- 1306 ○ Foreign content may be both non-archivable and unsupported by EPUB reading systems
1307 the archive is able to use.
- 1308 ○ For the time being, there are no video codecs among the core media types. There is a
1309 recommendation that reading systems should support either H.264 or VP8. Neither of
1310 these are archivable or even ingestible formats in the Finnish National Digital Library
1311 specification, which approves JPEG 2000 sequence and MPEG-4 AVC as archive formats

- 1312 and DV (Digital Video), MPEG-1, MPEG-2, and WMV (Windows Media Video) as
 1313 ingestible formats.
- 1314 • External references [BL]
 - 1315 ○ Externally referenced content (metadata, core media types, or foreign resources) SHALL
 1316 be retrieved during pre-ingest and embedded into the SIP, during ingest and embedded
 1317 into the AIP. If retrieval fails, the AIP is incomplete. If the submission agreement allows
 1318 such a policy, the archive can store the incomplete AIP and try to retrieve the missing
 1319 content post-ingest. If the second attempt is successful, the AIP is ingested again into the
 1320 repository system, and the missing content is added.
 - 1321 • Missing or poor fallback documents
 - 1322 ○ If a foreign resource cannot be rendered, there SHOULD be a core media type fallback
 1323 document. But according to EPUB 3.1 in some circumstances fallback can be omitted. In
 1324 all EPUB versions there is a risk that even if a fallback resource is present it may not
 1325 produce the same rendition than the original resource and there is no guarantee either
 1326 that the original semantics will be preserved.

1327 Recommendations:

- 1328 • EPUB 3 covers the superset of the expressive abilities of all the other e-book formats. Therefore
 1329 there is no technical or functional reason not to use and establish EPUB 3 as an interoperable
 1330 open e-book format standard [Bläsi, p. 8]. Having a universally supported e-book format would
 1331 benefit current e-book users and make long-term preservation of e-books easier.
- 1332 • Radium⁴⁷ project is developing a robust and efficient reader for EPUB publications. Such tools
 1333 will make it easier to use rich EPUB documents, and EPUB community should continue
 1334 investments on Radium and similar initiatives.
- 1335 • The EPUB community should create EPUB/A, a subset of EPUB 3 with features suitable for long-
 1336 term preservation. The specification should be complemented by an explanation why the EPUB
 1337 3 features not included in the EPUB/A format may jeopardize digital preservation, and a
 1338 justification for those features as that are required.
- 1339 • When new EPUB core media types are added, the archivability of these file formats should be
 1340 taken into account. EPUB community could co-operate with the digital preservation community
 1341 to achieve this goal.
- 1342 • Legal aspects of long-term preservation of EPUB 3 documents should be investigated.
- 1343 • Open source licenses such as SIL Open Font License⁴⁸ should be used when possible.
- 1344 • Foreign resources should be used with caution, until the archivability of the utilized file formats
 1345 has been verified.
- 1346 • Core media types that are considered to be non-archivable should be avoided whenever
 1347 possible. For instance, it is better to use a JPEG or a PNG than a GIF image.

1348

1349

⁴⁷ <http://readium.org/>

⁴⁸ http://scripts.sil.org/OFL_web

1350 **Bibliography**

- 1351 [1] Bläsi, Christoph and Franz Rothlauf: *On the interoperability of eBook formats*. [online]. Brussels:
 1352 European and International Booksellers Federation, 2013. Available from: [http://wi.bwl.uni-](http://wi.bwl.uni-mainz.de/publikationen/InteroperabilityReportGutenbergfinal07052013.pdf)
 1353 [mainz.de/publikationen/InteroperabilityReportGutenbergfinal07052013.pdf](http://wi.bwl.uni-mainz.de/publikationen/InteroperabilityReportGutenbergfinal07052013.pdf) [viewed 2017-06-
 1354 12].
- 1355 [2] Daly, Liza: *EPUB 3 and interactivity*. [online]. EPUBZone, 2014. Available from:
 1356 <http://epubzone.org/news/epub-3-and-interactivity> [viewed 2017-06-21].
- 1357 [3] Whibley, Simon: *EPUB format preservation assessment*. Version 1.2. [online]. British Library,
 1358 2015. Available from: http://wiki.dpconline.org/images/a/a9/EPUB_Assessment_v1.2.pdf
 1359 [viewed 2018-03-14].
- 1360 [4] Digital Preservation Handbook. 2nd, revised ed. [online]. Digital Preservation Coalition, 2017.
 1361 Available from: <http://www.dpconline.org/handbook>. [viewed 2017-07-25].
- 1362 [5] E-ARK: *Common specification for information packages. Version 1.0* [online]. E-ARK Project, 2016.
 1363 Available from: <http://www.dashboard.eu/specifications/common-specification>. [viewed 2017-
 1364 08-23].
- 1365 [6] EPUB 3 Fixed-Layout Documents [online]. International Digital Publishing Forum, 2012.
 1366 Available from: <http://www.idpf.org/epub/fxl/>. [viewed 2017-07-11].
- 1367 [7] EPUB 3.1. W3C Member submission 25 January 2017 [online]. Ed. by Markus Gylling, Tzviya
 1368 Siegman and Matt Garrish. W3C, 2017. Available from:
 1369 <https://www.w3.org/Submission/2017/SUBM-epub31-20170125/> [viewed 2017-07-11].
- 1370 [8] File formats. Version 1.5.1. [online]. The National Digital Library, 2017. Available from:
 1371 <http://www.kdk.fi/images/tiedostot/NDL-File-Formats-v1.5.1-en.pdf>. [viewed 2018-03-13].
- 1372 [9] ISO 14721. *Space data and information transfer systems – Open archival information system*
 1373 *(OAIS) – Reference model*. [online]. ISO, 2012. Available from:
 1374 <https://www.iso.org/standard/57284.html> [paywall]. [viewed 2017-07-10]. Also available as
 1375 CCSDS specification from: <https://public.ccsds.org/pubs/650x0m2.pdf> [viewed 2017-07-17].
- 1376 [10] ISO/FDIS 20614. Information and documentation – Data exchange protocol for interoperability
 1377 and preservation. [online]. ISO, 2017. The standard, when completed, will be available from:
 1378 <https://www.iso.org/standard/68562.html>. [viewed 2017-07-17].
- 1379 [11] ISO 20652. Space data and information systems - Producer-archive interface - Methodology
 1380 abstract standard [PAIMAS] . [online]. ISO, 2006. Available from:
 1381 <https://www.iso.org/standard/39577.html>. [paywall]. [viewed 2017-07-01]. Also available as
 1382 CCSDS 651.0-M-1:2004 from: <https://public.ccsds.org/publications/archive/651x0m1.pdf>
 1383 [viewed 2017-06-22].
- 1384 [12] LAVOIE, Brian: *Meeting the challenges of digital preservation: The OAIS reference model*. [online].
 1385 Dublin, OH: OCLC Research, 2000. Available from:
 1386 <http://www.oclc.org/research/publications/library/2000/lavoie-oais.html> [viewed 2017-07-
 1387 17].
- 1388 [13] PDF/A [online]. Wikipedia, 2017-05.17. Available from: <https://en.wikipedia.org/wiki/PDF/A>
 1389 [viewed 2017-07-18]. Archived version available from:
 1390 <https://web.archive.org/web/20170713180152/https://en.wikipedia.org/wiki/PDF/A>

- 1391 [14] TI/A Standard Initiative homepage. [online]. TI/A Standard Initiative, 2016. Available from:
1392 <http://ti-a.org/>. [viewed 2017-07-18].
- 1393 [15] W3C. EPUB 3 Community Group Charter. [online]. W3C, 2017. Available from:
1394 <https://www.w3.org/2017/02/EPUB3CGcharter>. [viewed 2017-07-24].
- 1395