---

**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:**  Final Text submitted for TR publication

**Title:**  Call for final review for publication of ISO/IEC TS 22424-1 Digital publishing — EPUB 3 Preservation — Part 1: Principles

**Status:**  This document is final text for publication of ISO/IEC TS 22424-1 Digital publishing — EPUB 3 Preservation — Part 1: Principles reflecting comments from PDTS ballots. Through four-week review in SC 34 and SC 34/JWG 7 this document will be submitted to SC 34 for publication. If you have any further comments or suggestion on this document, please let project editor (Juha Hakala, juha.hakala@helsinki.fi) by August 23, 2019.

**Date of document:**  2019-07-17

**Source:**  Project editor (Juha Hakala)

**Expected action:**  INFO

**No. of pages:**  38

**Email of convenor:**  samoh@g.skku.edu, zzosang@gmail.com

**Committee URL:**  https://isotc.iso.org/livelink/livelink/open/jtc1sc34jwg7

1 **ISO/IEC TS 22424-1:2019**

2 ISO/IEC JTC 1/SC 34/JWG 7

3 Secretariat: JISC

4 **Digital publishing — EPUB 3 Preservation — Part 1: Principles**

5 **Édition numérique — Archivage pérenne de l'EPUB 3 - Partie I :**
6 **Principes**

7

8 # TS stage

9

10 **Warning for WDs and CDs**

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

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

15

# Contents

# Foreword

ISO (the International Organization for Standardization) and IEC (the International Electrotechnical Commission) form the specialized system for worldwide standardization. National bodies that are members of ISO or IEC participate in the development of International Standards through technical committees established by the respective organization to deal with particular fields of technical activity. ISO and IEC technical committees collaborate in fields of mutual interest. Other international organizations, governmental and non-governmental, in liaison with ISO and IEC, also take part in the work. In the field of information technology, ISO and IEC have established a joint technical committee, ISO/IEC JTC 1.

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

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

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

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

This document was prepared by Technical Committee ISO/IEC JTC 1, *Information technology*, Subcommittee SC 34, *Document description and processing languages*.

A list of all parts in the ISO/IEC TS 22424 series can be found on the ISO website.

# Introduction

This document facilitates the long-term preservation of EPUB publications by specifying in general level EPUB features which are mandatory for long-term preservation (such as font embedding) and features which should be avoided if possible.

This specification can be seen as a stepping stone towards a detailed specification which would be related to EPUB in the same way as PDF/A, specified in ISO 19005-1 – 19005-3, is related to PDF. If and when the EPUB community develops detailed guidelines for the production of archivable EPUB publications, this document could be used as one of the starting points.

Long-term preservation in general requires two things:

- making the object such as EPUB publication fit for preservation – including features to be used and features to avoid;

- the packaging of the object (and any metadata related to it) together with any additional data such as other versions of the object and other documentation into an OAIS Submission Information Package (SIP).

Packaging is covered in Part 2 of this technical specification.

**EPUB**

The EPUB standard

> *defines a distribution and interchange format for digital publications and documents. The EPUB® format provides a means of representing, packaging and encoding structured and semantically enhanced Web content — including HTML, CSS, SVG and other resources — for distribution in a single-file container [EPUB 3.0.1].*

EPUB format was developed by the International Digital Publishing Forum, IDPF, which merged with the World Wide Web Consortium, W3C, in January 2017. Ongoing technical development of the standard, related extension specifications and ancillary deliverables are the responsibility of the W3C EPUB 3 Community Group[1], which published its charter in February 2017. According to the charter,

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

The International Digital Publishing Forum, IDPF, has ceased operations as a membership organization in January 2017, and its website[2] is now an archive. The latest version of the standard and information about future EPUB developments is available at the Publishing@W3C webpage, https://www.w3.org/publishing/.

The specification at hand covers EPUB 3 versions up to EPUB 3.0.1[3]. EPUB 3.1[4] was the first major revision of EPUB 3.0.1, but there are no implementations of version 3.1 and therefore it is not covered in this document. The most widely used version of the standard is still 3.0.1. EPUB 3.2, was published in

---

[1] https://www.w3.org/publishing/groups/epub3-cg/

[2] http://idpf.org/

[3] http://idpf.org/epub/301

[4] https://www.w3.org/Submission/epub31/

107    May 2019[5]. Unlike 3.1, it is fully backwards compatible with 3.0.1. It will be covered in the next edition
108    of this document.

109    Differences between EPUB specifications 2.0.1-3.2 are well documented:
110
111      •   EPUB 3 Changes from EPUB 2.0.1[6]
112      •   EPUB 3.0.1 Changes from EPUB 3.0[7]
113      •   EPUB 3.2 Changes from EPUB 3.0.1[8]

114
115    All EPUB specifications are available in the Web; 2.01 at http://idpf.org/epub/201, EPUB 3.0.1 at
116    http://idpf.org/epub/301 and 3.2 at https://w3c.github.io/publ-epub-revision/epub32/spec/epub-
117    spec.html.
118
119    All EPUB publications, including ones using version 3.2, can be validated using EPUBCheck version
120    4.2.0, which was released in March 2019.
121
122    From long-term preservation point of view, lack of backward compatibility between successive versions
123    of a file format would be a problem because it makes migration more challenging. In addition, EPUB 3.1
124    has at least one feature which would have been problematic. In EPUB 3.1 foreign resources do not
125    require fallbacks if they are not in the spine and not embedded in EPUB Content Documents. In EPUB
126    3.0.1, fallback guarantees that there is a version of the document that can be rendered; in 3.1 such
127    guarantee no longer exists.

128    EPUB 3.0.1 was prepared by the IDPF.  It consists of six interlinked documents:

129      •   EPUB 3 Overview
130      •   Publications 3.0.1
131      •   Canonical fragment identifiers
132      •   Content documents 3.0.1
133      •   Media overlays 3.0.1
134      •   Open Container Format 3.0.1

135    There are several extension specifications to these EPUB base standards. The list below is incomplete,
136    as it contains mainly specifications that are relevant from the long-term preservation point of view.
137    Some of them are still drafts:

138      •   EPUB Accessibility specification 1.0[9] addresses evaluation and certification of accessible EPUB
139         Publications, and discovery of the accessible qualities in such publications.
140      •   EPUB Previews 1.0[10] describes how content previews can be included in EPUB publications.
141      •   EPUB Distributable Objects 1.0[11] is a draft specification that defines a method for the
142         encapsulation, transportation, and integration of distributable objects in EPUB publications.
143      •   EPUB Scriptable Components 1.0[12] provides an interoperable publish and subscribe (pubsub)
144         pattern by which interactive content can be created and incorporated into EPUB publications.
145         Same as EPUB Distributable Objects, it is as of this writing (2019-05-13) a draft.
146      •   EPUB Scriptable Components Packaging and Integration 1.0[13] is a draft that defines a method
147         for the creation and inclusion of dynamic and interactive components in EPUB publications.

---

[5] https://w3c.github.io/publ-epub-revision/epub32/spec/epub-spec.html

[6] http://www.idpf.org/epub/30/spec/epub30-changes-20111011.html

[7] http://www.idpf.org/epub/301/spec/epub-changes-20140626.html

[8] https://w3c.github.io/publ-epub-revision/epub32/spec/epub-changes.html

[9] http://www.idpf.org/epub/a11y/accessibility.html

10 http://www.idpf.org/epub/previews/epub-previews-20150826.html

11 http://www.idpf.org/epub/do/

12 http://www.idpf.org/epub/sc/api/

148   • EPUB Multiple-Rendition Publications 1.0[14] defines the creation and rendering of EPUB
149     publications consisting of more than one rendition of the same publication.
150   • EPUB Dictionaries and Glossaries 1.0[15] provides a means for expressing dictionary and glossary
151     semantics in EPUB publications.
152
153   These extensions are not widely used and they have not been explicitly taken into account in this
154   document. As regards accessibility, all EPUB publications are supposed to be accessible. However,
155   accessibility features as such do not have an impact on long term preservation of EPUB publications and
156   therefore this document does not make accessibility-related requirements.
157
158   EPUB 3 Core Media Types have been listed at https://idpf.github.io/epub-cmt/v3/. As of this writing
159   [2019-05-13], the latest change has been made on April 1, 2018. Starting from EPUB 3.2, core media
160   types are part of the standard.
161
162   In 2014, EPUB 3.0 specifications were republished as a standard, ISO/IEC TS 30135 parts 1-6, by the
163   International Standards Organization. Each of these six ISO specifications is identical to its IDPF
164   equivalent, for example TS-30135-1 has exactly the same content as the EPUB 3.0 Overview.
165
166   ISO/IEC TS 30135-7 is "Part 7: EPUB3 Fixed-Layout Documents" is from EPUB 3.0.1 (EPUB 3.0 does not
167   have fixed layout specification).  TS 30135 is therefore a combination of EPUB 3.0 and Fixed-Layout
168   Documents specification from 3.0.1.
169
170   ISO/IEC JTC 1/SC 34 is currently updating the ISO standard to match fully the version 3.0.1.
171
172   EPUB is a rich document format with a lot of features. From the digital preservation point of view this is
173   a challenge, not least because long-term preservation has not been a priority in the development of the
174   standard. Preserving all aspects and features of EPUB publications may be difficult, since there are
175   features which are difficult to preserve. Moreover, EPUB reading systems usually do not support all
176   features of the specification and finding tools supporting rare features can be difficult.
177
178   In spite of these challenges EPUB is generally regarded as a suitable format for digital archiving. For
179   instance, the Finnish National Digital Library initiative has selected just eight archivable file formats for
180   text, EPUB being one of them. The selection criteria were openness/transparency, adoption as a
181   preservation standard, degree of forward/backward compatibility, degree of protection against file
182   corruption, frequency of version releases, dependencies/interoperability, and standardization. EPUB
183   got an A, the best grade, from everything else except the second and third criterion. For those, the grade
184   was the second best, a B [File formats, p. 40]. Based on these generic criteria, EPUB seems to provide a
185   good basis for long-term preservation, although additional guidelines on how to use the standard are
186   needed to guarantee EPUB files can be preserved efficiently.
187
188   The British Library's Digital Preservation Team has published an assessment of EPUB as a preservation
189   format [Day]. It covers EPUB versions 3.0.1 and 2 and the overall view of EPUB is positive [Day, p. 2]:
190
191         *EPUB 3 is currently the closest thing available to an open standard for e-books. In 2013,*
192         *Bläsi and Rothlauf concluded that EPUB 3 had the "highest expressive power" of all formats*
193         *in the e-book ecosystem, and that it included the superset of all features used in proprietary*
194         *formats like KF8, Fixed Layout EPUB, and iBooks.*
195
196   EPUB is enjoying reasonable support in the e-book market. Many suppliers, publishers, and application
197   developers who have supported EPUB 2 have implemented version 3.0.1. According to the EPUBTest

---

13 http://www.idpf.org/epub/sc/pkg/
14 http://www.idpf.org/epub/renditions/multiple/
15 http://www.idpf.org/epub/dict/

198  web site[16], EPUB 3 support in reading systems is far from exhaustive, but market coverage is good – in
199  January 2018, there were 59 reading systems supporting at least some of the features specified in EPUB
200  3.0.
201
202  E-book suppliers have produced EPUB 3 based formats that incorporate Digital Rights Management
203  (DRM), and EPUB modifications that may restrict using the format on other than the suppliers' own
204  platforms. For example, the Kindle Fire eReader, released in 2015, uses a new format called Kindle
205  Format 8 (KF8), which is partly based on EPUB 3, with Amazon's DRM. [Day, 3]. Publisher/supplier
206  specific DRM often restricts the use of e-books to that publisher's/supplier's rendering devices and/or
207  applications, and is therefore a major obstacle to digital preservation [Day, p. 7].
208
209  The EPUB specification does not enforce a particular Digital Rights Management scheme, but DRM may
210  be layered on top of the EPUB specifications. A producer can, for instance, use one of the three major
211  rights management systems in the market (Amazon DRM, Apple FairPlay DRM for books bought from
212  iBooks, and Adobe DRM), or some other DRM system along with some additional platform-targeting.
213
214  DRM protection should be removed from EPUB publications during pre-ingest by the producer or as a
215  part of the ingest process by the OAIS archive. In practice, only national libraries may be able to do this,
216  provided that legal deposit act and / or copyright act guarantee them such privilege. If migration is the
217  chosen preservation strategy, existing EPUB publications will be converted into more modern EPUB
218  versions when rendering tools for old versions are no longer available, and (eventually) migrated into
219  other formats.
220
221  If preserved EPUB publications are not directly accessible by the public, removing DRM, digital
222  watermarking, and other protection mechanisms from the archived documents is not a risk. When
223  publications are delivered to the customers as Dissemination Information Packages (DIPs), the archive
224  shall use a combination of administrative and technical means to protect the documents as required in
225  the submission agreement. These means may include adding DRM protection mechanism into the DIP
226  submitted to the user according to the requirements of the submission agreement. The agreement may
227  also specify the customers the archive is entitled to serve; for instance, it is possible to require that the
228  preserved documents can only be disseminated to the producer, and the producer will serve the end-
229  users who do not have direct access the OAIS archive.
230
231  **Digital preservation**

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

235  Digital preservation is not about preserving just bits, but about preserving access.  The "business logic"
236  is as follows:

237  • we need software and hardware to render content for human users

238  • software changes over time; there are new versions from old applications, and entirely new
239    applications

240  • new or updated applications may not be able to render outdated file formats or format versions
241    correctly

242  • digital preservation makes an effort to have all archived content in stable formats. Publications
243    should also contain the smallest possible amount of features which are not commonly
244    supported in software packages used to render the content in these formats, and also avoid

---

[16] http://epubtest.org/testsuite/epub3/

245     adding links to external resources since then the long-term access to the publication requires
246     also persistence of these external resources.

247     • when necessary, data in old formats may be migrated into more modern formats or updated
248       versions of the same format. For instance, an e-book in EPUB 3.0.1 format may be migrated to
249       EPUB 5.2. when version 3.0.1 is no longer widely supported by reading systems.

250     • since the aim is to preserve the content, not the bits, the bits may change as a result of version
251       updates and format migrations.

252     • Many OAIS archives preserve successive versions of archives publications, because migration
253       may change the look and feel of the original document, or even its intellectual content.

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

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

264 Preservation metadata, which allows the publication to be found, rendered and authenticated correctly,
265 is a prerequisite for digital preservation. Some preservation metadata elements can or should be
266 provided by the original creator of the publication. It is also important to keep preservation
267 requirements in mind when preparing a publication, if it is known that it has to be preserved for a long
268 time. Any feature in a file format can be either essential, useful, neutral, questionable, or even
269 downright counterproductive from a long-term preservation point of view. However, publishers are
270 likely to use the features that let them achieve their own goals, and preservation may not be among
271 them.

272 There already are archivable versions of some file formats. PDF/A (ISO 19005-1:2005)[17] is probably the
273 best known example. It specifies how to use the Adobe Portable Document Format (PDF) for long-term
274 preservation. An example of a counterproductive feature for preservation in PDF is font referencing;
275 therefore in PDF/A all fonts shall be embedded in order to guarantee that the document can be
276 rendered correctly.

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

280     *Information security methods such as encryption add to the complexity of the preservation*
281     *process and should be avoided if possible for archival copies. Other security approaches may*
282     *therefore need to be more rigorously applied for sensitive unencrypted files; these might*
283     *include restricting access to locked-down terminals in controlled locations (secure rooms),*
284     *or strong user authentication requirements for remote access.*
285
286 In order to guarantee the correct processing of PDF/A files, there are specific requirements for PDF/A
287 reading systems, such as support for embedded fonts. There are three versions of the specification:
288 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

---

[17] https://www.iso.org/standard/38920.html

ix

289 all the features of PDF/A-2 as well as allows the embedding of other file formats into PDF/A conforming
290 documents [PDF/A].

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

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

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

305 TIFF is not an EPUB core media type, but four other image types have been listed; GIF, JPEG, PNG, and
306 SVG. It is significant from a digital preservation point of view how these formats and other core media
307 types are used in the EPUB context. Image and audio files embedded in an EPUB publication may
308 require migration before the EPUB publication itself has to be migrated into a more modern file format,
309 if commonly available EPUB reading systems no longer support these file formats. This specification
310 does not provide guidelines for creating archivable files in EPUB 3 core media types, due to the
311 magnitude of such task. But EPUB community SHOULD follow the archival file format lists of national
312 archives or libraries (for example the Library of Congress file format list[18] and the U.S. National
313 Archives list[19]) when the core media file format list is updated. Publishers SHOULD also consider the
314 persistence of file formats used when creating EPUBs for which the need for long-term preservation is
315 foreseen. .

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

321 **OAIS and related standards**
322 ISO 22424 provides guidance on how to utilize the Open Archival Information System (OAIS) and
323 current practices of OAIS archives in preservation of EPUB publications. The OAIS [ISO 14721] is
324 equally relevant to both parts of the ISO 22424.

325 OAIS is a reference model for long-term data storage systems. It is used by memory institutions
326 (libraries, archives, and museums) and many other organizations that need to preserve digital
327 resources in the long-term. Although an ISO standard, the OAIS was originally developed by the CCSDS,
328 The Consultative Committee for Space Data Systems[20], which still maintains the specification.

329 The model has five functional units:

---

[18] http://www.loc.gov/preservation/digital/formats/
[19] https://www.archives.gov/records-mgmt/policy/transfer-guidance-tables.html
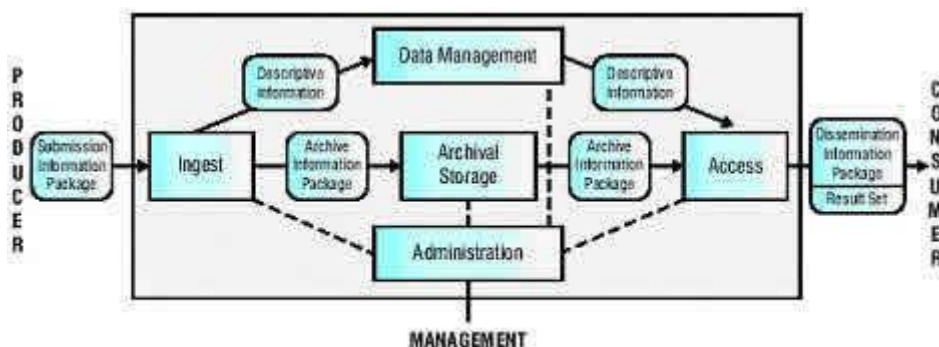[20] https://public.ccsds.org/default.aspx

Figure 3

330

**Figure 1. OAIS Model [Lavoie]**

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

338 Modifying an EPUB publication so that it is suitable for digital archiving is from the OAIS point of view a
339 part of pre-ingest and as such not a part of the OAIS model. The importance of the OAIS to the ISO
340 22424 is that the model provides a terminology, information package data model and an overall
341 framework within which digital preservation can be performed.

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

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

357 OAIS submission agreements specify the principles of how documents should be prepared and
358 submitted to the repository system. If the archive uses migration as the preservation method[21],
359 submission agreements should specify file formats (and metadata formats) suitable for submission

---

[21] 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.

360  and/or archival, or refer to external documents listing these formats. File formats suitable for
361  submission but not for archival are migrated during the ingest process, although the original files may
362  be included in the AIP.

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

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

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

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

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

389

# Digital publishing — EPUB 3 Preservation — Part 1: Principles

## 1  Scope

This document supports long-term preservation of EPUB publications via a dual strategy. First, it considers EPUB features from long-term preservation point. Some EPUB features are forbidden and some others required, depending on how they relate to a long-term preservation. EPUB publications constructed according to these guidelines should be suitable for preservation.

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

## 2  Normative references

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

ISO/IEC TS 30135 (all parts), *Information technology — Digital publishing — EPUB3*

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

## 3  Terms and definitions

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

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

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

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

**3.1**
**access functional entity**
OAIS functional entity that contains the services and functions, which make the archival information holdings and related services visible to Consumers

**3.2**
**administrative metadata**
metadata that provides information to help manage a resource, such as when and how it was created, file type and other technical information, and access rights

[SOURCE: Understanding metadata]

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

427 **3.4**
428 **archive**
429 **OAIS archive**
430 organization that intends to preserve information for access and use by a Designated Community
431 **3.5**
432 **authenticity**
433 property that an entity is what it claims to be

434 [SOURCE: ISO/IEC 27000]

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

436 **3.6**
437 **bit preservation**
438 term used to denote a very basic level of preservation of digital resource as it has been submitted
439 (literally the preservation of the **bits** forming a digital resource)

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

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

445 [SOURCE: Digital preservation handbook, Glossary]

446 **3.7**
447 **consumer**
448 role played by those persons or client systems, who interact with OAIS services to find preserved
449 information of interest and to access that information in detail

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

451 **3.8**
452 **content information**
453 set of information that is the original target of preservation or that includes part or all of that
454 information

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

457 **3.9**
458 **context information**
459 information that documents the relationships of the Content Information to its environment

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

462 **3.10**
463 **core media type**
464 a set of publication resource for which no fallback is required.

465 [SOURCE: EPUB Publications 3.0.1]

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

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

471  **3.11**
472  **data, pl**
473  reinterpretable representation of information in a formalized manner suitable for communication,
474  interpretation, or processing

475  [SOURCE: ISO 5127:2017]

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

478  **3.12**
479  **data dictionary**
480  organized and constructed (electronic data base) compilation of descriptions of data concepts that
481  provides a consistent means for documenting, storing and retrieving the syntactical form (i.e.
482  representational form) and the meaning and connotation of each data concept

483  [SOURCE: ISO 24531:2013]

484  Note 1 to entry: PREMIS[22] is a data dictionary.

485  **3.13**
486  **descriptive metadata**
487  **descriptive information**
488  metadata about a resource for example for discovery and identification

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

490  [SOURCE: Understanding metadata]

491  **3.14**
492  **designated community**
493  identified group of potential Consumers who should be able to understand a particular set of
494  information

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

497  **3.15**
498  **digital preservation**
499  series of managed activities necessary to ensure continued access to digital materials for as long as
500  necessary

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

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

---

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

**3**

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

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

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

512 [SOURCE: Digital preservation handbook, Glossary]

513 **3.16**
514 **digital rights management**
515 **DRM**
516 packaging, distributing, controlling, and tracking content based on rights and licensing information

517 [SOURCE: ISO 19153:2014]

518 **3.17**
519 **digital signature**
520 **signature**
521 data appended to, or a cryptographic transformation of, a data unit that allows the recipient of the data
522 unit to prove the source and integrity of the data unit and protect against forgery, e.g. by the recipient

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

524 **3.18**
525 **dissemination information package**
526 **DIP**
527 information package, derived from one or more AIPs, sent by an Archive to a Consumer in response to a
528 request in the OAIS

529 **3.19**
530 **distributable object**
531 component of an EPUB publication that can be reused in other contexts

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

536 [SOURCE: EPUB Distributable Objects 1.0]

537 **3.20**
538 **electronic book**
539 **e-book**
540 non-serial digital document, licensed or not, where searchable text is prevalent, and which can be seen
541 in analogy to a print book

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

544 [SOURCE: ISO 2789:2013]

545 **3.21**
546 **EPUB container**
547 ZIP based packaging and distribution format for EPUB publications

548 [SOURCE: EPUB Publications 3.0.1]

549 **3.22**
550 **EPUB content document**
551 publication resource that conforms to one of the EPUB content document definitions

552 [SOURCE: EPUB Publications 3.0.1]

553 **3.23**
554 **EPUB navigation document**
555 specialization of the XHTML content document, containing human- and machine-readable global
556 navigation information

557 [SOURCE: EPUB Publications 3.0.1]

558 **3.24**
559 **EPUB publication**
560 collection of one or more renditions conforming to the EPUB specifications, packaged in an EPUB
561 container

562 [SOURCE: EPUB Publications 3.0.1]

563 **3.25**
564 **EPUB reading system**
565 system that processes EPUB publications for presentation to a user in a manner compliant with EPUB
566 specifications

567 [SOURCE: EPUB Publications 3.0.1]

568 **3.26**
569 **fallback**
570 mechanism with which versions of the same resource in different file formats can be linked to one
571 another

572 [SOURCE: EPUB Publications 3.0.1]

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

575 **3.27**
576 **fixity information**
577 information that documents the authentication mechanisms and provides authentication keys to ensure
578 that the Content Information object has not been altered in an undocumented manner

579 [SOURCE: ISO 13527:2010]

580 **3.28**
581 **foreign resource**
582 publication resource that is not a core media type

583 [SOURCE: EPUB Publications 3.0.1]

584 **3.29**
585 **identifier**
586 data string or pointer that establishes the identity of an item, institution, or person alone or in
587 combination with other elements.

588 [SOURCE: ISO 8459:2009]

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

591 **3.30**
592 **independently understandable**
593 characteristic of information that is sufficiently complete to allow it to be interpreted, understood, and
594 used by the Designated Community without having to resort to special resources not widely available,
595 including named individuals

596 **3.31**
597 **information**
598 any type of knowledge that can be exchanged

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

600 EXAMPLE    a string of bits (the data) accompanied by a description on how to interpret the string of
601 bits as numbers representing temperature observations measured in degrees Celsius (the
602 representation information)

603 **3.32**
604 **information package**
605 logical container composed of optional content information and optional associated preservation
606 description information

607 **3.33**
608 **ingest functional entity**
609 OAIS functional entity that contains the services and functions that accept SIPs from producers,
610 prepares AIPs for storage, and ensures AIPs and their supporting descriptive information become
611 established within the OAIS

612 **3.34**
613 **long-term**
614 period of time long enough to raise concerns about the impact of changing technologies, including
615 support for new media and data formats, and of a changing designated community, on the information
616 being held in an OAIS

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

618 **3.35**
619 **long-term preservation**
620 act of maintaining information, independently understandable by a designated community, with
621 evidence supporting its authenticity over the long-term

622 **3.36**
623 **manifest**
624 EPUB manifest element provides an exhaustive list of the Publication Resources that constitute the
625 given Rendition, each represented by an item  element.

626 [SOURCE: EPUB Publications 3.0.1]

**3.37 metadata**
data about other data, documents, or records that describe their content, context, structure, format, provenance, and/or rights.

[SOURCE: ISO 5127:2017]

**3.38**
**METS**
Metadata Encoding and Transmission Standard, a standard for presenting metadata using XML.

[SOURCE: Digital preservation handbook, Glossary]

**3.39**
**migration**
means of overcoming technological obsolescence by transferring digital resources from one hardware/software generation to the next

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

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

[SOURCE: Digital preservation handbook, Glossary]

**3.40**
**Open Archival Information System**
**OAIS**
archive, consisting of an organization, which may be a part of a larger organization, of people and systems, that has accepted the responsibility to preserve Information and make it available to a Designated Community. It has a set of responsibilities, as defined in section 4, which allow an OAIS Archive to be distinguished from other uses of the term 'Archive'.

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

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

**3.41**
**package document**
publication resource that describes one rendition of an EPUB publication, as defined in package document. The package document carries meta information about the Rendition, provides a manifest of resources and defines the default reading order.

[SOURCE: EPUB Publications 3.0.1]

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

667 **3.42**
668 **PDF**
669 Portable Document Format, a set of formats and open standards maintained by the International
670 Organization for Standardization for producing and sharing electronic documents

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

672 [SOURCE: Digital preservation handbook, Glossary]

673 **3.43**
674 **PDF/A**
675 versions of the PDF standard intended for archival use

676 [SOURCE: Digital preservation handbook, Glossary]

677 **3.44**
678 **pre-ingest**
679 actions required before data can be submitted into an OAIS archive, including negotiation of data
680 acquisitions, checking rights and access criteria, licensing, and data submission

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

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

684 [SOURCE: UK Data Archive. Archive training manual[23]]

685 **3.45**
686 **preservation description information**
687 **PDI**
688 information necessary for the adequate preservation of Content Information that can be categorized as
689 provenance, reference, fixity, context, and rights information

690 **3.46**
691 **preservation metadata**
692 metadata containing information needed to archive and preserve a resource

693 [SOURCE: Understanding metadata]

694 **3.47**
695 **preservation planning functional entity**
696 OAIS functional entity that provides the services and functions for monitoring the environment of the
697 OAIS and that provides recommendations and preservation plans to ensure information stored in the
698 OAIS remains accessible to, understandable by, and sufficiently usable by the designated community
699 over the long-term, even if the original computing environment becomes obsolete

700 **3.48**
701 **producer**
702 role played by those persons or client systems that provide the information to be preserved

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

---

[23] http://www.data-archive.ac.uk/curate/archive-training-manual/pre-ingest

**3.49**
**provenance information**
information that documents the history of the Content Information

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

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

**3.50**
**publication resource**
resource that has the content or instructions contributing to the logic and rendering of at least one rendition of an EPUB publication

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

[SOURCE: EPUB Publications 3.0.1

**3.51**
**reading system**
system that processes EPUB publications for presentation to a user in a manner conformant with EPUB specification

[SOURCE: Modified from EPUB Publications 3.0.1]

**3.52**
**reference information**
information that is used as an Identifier for the Content Information

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

EXAMPLE          an ISBN is a type of Reference Information.

**3.53**
**reference model**
framework for understanding significant relationships among entities in an environment and for the development of consistent standards or specifications supporting that environment

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

**3.54**
**reformatting**
copying information content from one storage medium to a different storage medium (media reformatting) or converting from one file format to a different file format (file reformatting)

[SOURCE: Digital preservation handbook, Glossary]

**3.55**
**refreshing**
copying information content from one storage media to the same storage media

746    [SOURCE: Digital preservation handbook, Glossary]

747    **3.56**
748    **release identifier**
749    identifier that allows any instance of an EPUB publication to be compared against another to determine
750    if they are identical, different versions, or unrelated

751    Note 1 to entry: Release Identifiers consist of a unique identifier and the last-modified date of the document.

752    [SOURCE: EPUB Publications 3.0.1]

753    **3.57**
754    **remotely-hosted resource**
755    objects hosted outside the EPUB Container.

756    **3.58**
757    **rendition**
758    one rendering of the content of an EPUB publication, as expressed by an EPUB package

759    [SOURCE: EPUB Publications 3.0.1]

760    **3.59**
761    **repository system**
762    long-term preservation system used by an archive
763
764    **3.60**
765    **rights management metadata**
766    information that identifies the access restrictions concerning the Content Information, including the
767    legal framework, licensing terms, and access control

768    Note 1 to entry: This contains the access and distribution conditions stated in the Submission Agreement, related
769    to both preservation (by the OAIS) and final usage (by the Consumer).

770    Note 2 to entry: It also includes specifications for the application of rights enforcement measures.

771    **3.61**
772    **spine**
773    EPUB spine element defines the default reading order of the EPUB Publication content by defining an
774    ordered list of manifest item references.
775
776    [SOURCE : EPUB Publications 3.0.1]
777
778    **3.62**
779    **structural metadata**
780    metadata that indicates how compound objects are put together, for example how the pages of a
781    document are arranged to form chapters

782    [SOURCE: Understanding metadata]

783    **3.63**
784    **submission agreement**
785    agreement reached between an OAIS archive and a Producer that specifies a data model and any other
786    arrangements needed for the data submission session

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

789 **3.64**
790 **submission information package**
791 **SIP**
792 information package that is delivered by a Producer to an OAIS to be used to construct or update one or
793 more AIPs and/or the associated descriptive information.

794 **3.65**
795 **unique identifier**
796 primary identifier of an EPUB publication, which may be shared by one or several renditions of the
797 same EPUB publication that conform to the EPUB standard and embody the same content.

798 [SOURCE: EPUB Publications 3.0.1]

799 **3.66**
800 **XHTML content document**
801 EPUB content document that conforms to the profile for HTML defined in XHTML Content Documents

802 [SOURCE: EPUB Publications 3.0.1]

803 Note 1 to entry: see EPUB Content Documents 3.0.1, chapter 2.

804 # 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 |

805 # 5  Packaging standards

806 An archiving process includes several distinct steps. A producer – which may be the publisher or other
807 body acting on behalf of the publisher, such as the archive itself - creates a Submission Information
808 Package (SIP) and transfers it to a repository system in an OAIS archive. The archive performs a quality
809 control process to the SIP and, if the package meets the criteria set in the submission agreement,
810 accepts it, creates an Archival Information Package (AIP) and transfers the package to archival storage.
811 During ingest some of the files or metadata records within SIP may be migrated to new formats or
812 additional metadata may be added.

813 The OAIS archival storage function stores, maintains, and retrieves AIPs. Maintenance may include for
814 instance frequent error checks to protect the data against bit rot. In order to keep the documents
815 understandable it may also be necessary to migrate[24] them in new formats, or to update the AIP with
816 additional metadata related to emulation. Migration and other preservation related tasks may be
817 carried out by the producer, OAIS archive and / or third parties. The party or parties responsible should
818 be specified in the submission agreement.

---

[24] 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".

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

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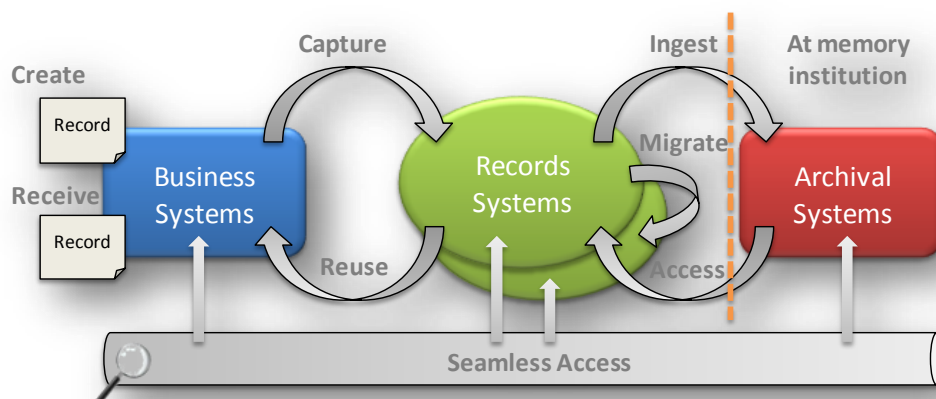832 reliable, such as other OAIS archive.

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

840 The content and structure of all information packages in repository systems MUST be standardized.
841 There are several packaging standards available, but the most commonly used one is the Metadata
842 Encoding and Transmission Standard (METS[25]) developed by the Library of Congress. ISO/IEC
843 JTC1/SC34 JWG 7 decided to recommend the use of METS as the container standard, although this
844 specification does allow the use of other container standards as well.

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

849 Some digital preservation projects have developed tools for creating SIPs that meet the project
850 requirements, which makes it a lot easier to submit information to the repository system. Producers
851 SHOULD nevertheless have at least basic understanding of digital preservation, since pre-ingest steps
852 from document creation to SIP submission should not risk the authenticity of the documents to be
853 submitted to the OAIS archive.

---

[25] http://www.loc.gov/standards/mets/

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

Different disciplines, even if they all use OAIS, will develop interfaces optimized for their own needs. And if the payloads are not the same, technical metadata standards will also differ. Domains have also adopted different packaging and preservation metadata standards. Almost all digital archiving projects in the library domain rely on METS and PREMIS specifications. Some libraries use BagIt[26] instead of METS for storage of ditigal objects, but BagIt specification does not require knowledge of the semantics of the resources in the container, whereas METS supports such metadata. Therefore BagIt is not an alternative to METS for long-term preservation.

Compared with libraries, the film industry started digital preservation efforts a bit later and may eventually develop different preferences[27]. And even if the same standards were used, they may be applied in a non-interoperable way even within the same domain. Therefore creating set application profiles is important in digital archiving.

## 6   Construction of OAIS information packages

According to the Open Archival Information System (OAIS) model[28], information package is "a container that contains two types of Information Objects, the Content Information and the Preservation Description Information (PDI)". Content information is the data that needs to be preserved and preservation description information is the metadata and other information that is needed in order to preserve, find and understand the data in long-term.

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

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

> *[i]t is necessary to distinguish between an Information Package that is preserved by an OAIS and the Information Packages that are submitted to, and disseminated from, an OAIS. These variant packages are needed to reflect the reality that some submissions to an OAIS will have insufficient PDI to meet final OAIS preservation requirements. In addition, they MAY be organized very differently from the way the OAIS organizes the information it is preserving. Finally, the OAIS MAY provide information to Consumers that does not include all the PDI*

---

[26] https://tools.ietf.org/html/rfc8493
[27] https://www.cen.eu/news/calls/Calls/CEN-Call_for-tender_Digitalcinema.pdf
[28] http://public.ccsds.org/publications/archive/650x0m2.pdf

884 *with the associated Content Information being disseminated. These variants are referred to*
885 *as the Submission Information Package (SIP), the Archival Information Package (AIP), and*
886 *the Dissemination Information Package (DIP). Although these are all Information Packages,*
887 *they differ in mandatory content and the multiplicity of the associations among contained*
888 *classes.*

889 The principles listed below provide SIP production guidelines for document producers (publishers or
890 third parties creating EPUB publications). The creation of the principles has been inspired by the draft
891 common requirements published by the E-ARK project (see Introduction to the Common Specification
892 for Information Packages in the E-ARK project, version 1.0[29]). Although E-ARK has served as a model for
893 this specification, these requirements have not been aligned with those of E-ARK, and therefore there
894 may be significant differences between the specifications.

## 6.1 General

### 896 6.1.1 EPUB publications SHALL be sent to a repository system as well-formed and complete
897 **Submission Information Packages (SIPs)**

898 • This specification does not assume that publishers create SIPs. The OAIS producer MAY be a
899 third party acting on behalf of the publisher, such as hosting platform or other production
900 vendor or even the OAIS archive itself.
901 • This specification and its accompanying document are mainly concerned with the structure
902 and content of SIPs. The way EPUB publications are archived and disseminated (the
903 structure of Archival Information Packages and Dissemination Information Packages, or
904 AIPs and DIPs) depends on the submission agreements made between the archive and the
905 producers, and on the operational principles of the archive, and is beyond the scope of this
906 document. It is possible that an EPUB publication is migrated into another format during
907 Ingest, and disseminated again as an EPUB publication. The archive may also preserve (in
908 bit level) the original file.
909 • Submitted EPUB publications SHALL be conformant with EPUB requirements[30] and
910 conformance SHOULD be validated.
911 • Submitted EPUB publications SHALL either contain or at least facilitate access to all the data
912 and metadata required to render the content information successfully.
913     i. Preview publications MAY be submitted, even though they are by definition not
914 complete, if the final documents are sent when ready. Depending on the submission
915 agreement, the archive MAY preserve just the final version, or both versions of the
916 resource. Identifiers SHALL be used in such a way that the OAIS archive will be able to
917 link all versions of the publication and delete preview versions, if that is the agreed
918 preservation policy.
919     ii. Distributable objects SHALL NOT be submitted individually. They MAY be embedded
920 within EPUB publications, but the archive is not obliged to deliver them as DIPs unless
921 the submission agreement mandates that.
922     iii. Fonts SHALL be embedded into the EPUB publication in full and un-obfuscated, if font
923 license allows that. If submission agreements allow submission of EPUB publications
924 with obfuscated or non-embedded fonts, there is a risk that such publications become
925 unusable in the future.
926     iv. Related resources such as audio and video SHOULD be embedded in the EPUB
927 publication.

---

[29] http://www.dasboard.eu/specifications/common-specification
[30] Conformance requirements for EPUB publications and reading systems have been specified in chapter 3.1 of EPUB Publications, version 3.0.1.

928        v.   Remotely-hosted resources SHOULD be avoided, but if used, it is necessary to ensure
929               that all remote data is available to the archive so that the data can be incorporated into
930               the AIP during ingest, and permission to do this SHALL be explicitly agreed upon in the
931               submission agreement, especially if the publisher is not in full control of remote data.
932        vi.   Descriptive and other metadata SHOULD be embedded in the SIP. METS mdRef element
933               MAY only be used if a) referred metadata is part of the same SIP, or b) the archive is
934               able to retrieve any linked external metadata and incorporate it into the AIPs in an
935               appropriate format.
936       vii.   Permission to use remote resources and metadata SHALL be specified in the
937               submission agreement[31]. The permission SHALL specify acceptable metadata and file
938               formats.

939      •   The SIP SHOULD[32] be checked for viruses and malicious software before submission to the
940         repository system.

941      •   EPUB publications in SIPs SHOULD NOT be encrypted, because that compromises long-term
942         preservation. If data is submitted in an encrypted format, the archive SHALL receive
943         necessary decryption information/details within the SIP, as agreed in the submission
944         agreement or elsewhere. When the archive disseminates the archived data to its customers,
945         it can be encrypted again.

946      •   DRM protection, if any, SHOULD be removed by the producer before the document is
947         submitted. If the content in the SIP is DRM protected, the archive SHALL receive the
948         necessary information/details to remove the DRM protection within the SIP, as agreed in
949         the submission agreement or elsewhere. Such permission may be producer-specific, based
950         on the submission agreement, or a generic permission, based on e.g. the Copyright Act.

951      •   If data is compressed, the user of the compression method SHALL be specified using the
952         Compression metadata element in the EPUB's encryption.xml file.

953      •   The submission agreement SHOULD specify at least one EPUB reading system capable of
954         rendering the submitted EPUB publications successfully. Knowing the reading system
955         requirements in advance makes it easier for the archive to design and implement the ingest
956         process. Although submitted publications will usually be validated only with automated
957         tools[33], the archive should be able to validate that the received EPUB can be presented to
958         the customers, and check for instance the look and feel of archived EPUB publications
959         before and after migration. This is possible only if the archive can operate the reading
960         systems that can render the archived publications successfully.

961      •   Each SIP SHOULD specify EPUB reading system or systems, which can render the EPUB
962         publication in the SIP. If this information is missing, reading system or systems SHALL be
963         specified in the submission agreement.

964      •   Multiple-rendition EPUB publications may be designed for multiple reading systems, in
965         which case the submission agreement may require the archive to carry out at least
966         occasional checks in all of these reading systems. If so, all these reading systems SHOULD be
967         listed in the submission agreement.

---

[31] 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. Submission agreement SHOULD specify how to handle such a situation. For instance, the agreement can require the producer either sends a new SIP with all the data and metadata embedded into it, or makes sure that the archive is allowed to access remote data and metadata.

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

[33] One such tool is Epubcheck, available from https://github.com/idpf/epubcheck

         **15**

968      •    If a submitted EPUB publication has been optimized for a certain reading system, the
969           system SHOULD be described in the document's technical and/or preservation metadata,
970           since such information is valuable for preservation and archival access purposes.
971      •    If the optimal EPUB reading system is no longer available, the archive SHOULD, with
972           permission and support from the producer, either find another suitable reading system or
973           modify the ingest process so that the EPUB publications affected by this change can be used
974           by another EPUB reading system.[34]

975 **6.1.2   Regardless of its type or format, it SHALL be possible to include any data or metadata in**
976 **SIPs**

977      •    It SHOULD be possible to maintain the SIP and EPUB specifications independently, i.e. so
978           that any change to SIP does not automatically mean that the EPUB format needs to be
979           updated and vice versa. The exception from this rule is that any existing and future features
980           in EPUB specification which are relevant from long-term preservation point of view such as
981           font embedding SHALL be taken into account in the SIP specification.
982      •    This document does not set a priori constraints either to the current or future versions of
983           EPUB with regard to the choice of metadata and file formats or either's versions (see note 1
984           below on EPUB Core media types).
985      •    The submission agreement SHOULD specify metadata formats and file formats approved for
986           submission and archival. For EPUB publications, at least Dublin Core metadata format and
987           all EPUB core media types SHALL be supported by the archive in order to guarantee
988           efficient processing of EPUB publications.
989      •    Submission agreements SHOULD specify what kind of executables can be embedded in the
990           submitted EPUB publications (see note 2 below on interactive e-books and EPUB
991           publications).

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

996           File format lists in submission agreements may cover all EPUB core media types or – if the
997           producer does not use all the core media types  - just a subset of them. When a core media
998           type is deprecated, the producer (if it still exists) and the archive should decide whether the
999           file format in question is migrated or kept as is (and emulated). If the latter, it may be
1000         necessary to migrate the deprecated file format when DIP packages are created.

1001 NOTE 2    E-books are likely to contain more interactive features in the future. From preservation
1002           point of view it is therefore a problem that there are various ways in which EPUB 3 can
1003           support interactivity. On the other hand, some EPUB reading systems may not support
1004           interactivity at all, and even if it is supported, different reading systems may not behave
1005           identically, partly because EPUB is not specific about how support should be implemented.
1006           EPUB 3 `object` element enables the use of arbitrary embedded executables that are not
1007           inherently supported in EPUB 3 reading systems. A common use case would be to include
1008           proprietary applets or Adobe Flash applications. However, in a majority of cases, interactive
1009           publications will be created through the use of in-book source code. Because JavaScript is
1010           the de facto standard scripting language for SVG and HTML5, EPUB 3 content documents
1011           can be assumed to be scriptable only if they contain JavaScript code. The standard does not
1012           define which versions of JavaScript (ECMAScript) are required to get the support. Content
1013           creators should comply with the most commonly supported features in web browsers for

---

[34] 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.

1014 best results [Daly]. Usage of common tools and techniques will also make it easier to
1015 preserve the publication in the long-term, either via emulation (a common solution for
1016 software preservation) or migration.

1017 • Archives offering long-term preservation services for EPUB publications SHOULD keep
1018 track of EPUB core media types and consider the possibility of including them on the list of
1019 archivable formats. If this is not viable, the archives SHOULD maintain clearly defined and
1020 well tested migration pathways from non-archivable core media types into archivable
1021 formats. Then the archive would not need to migrate these images during ingest and it
1022 would be possible to preserve EPUB publications unchanged[35].

1023 • If there is a foreign resource embedded or linked to a submitted EPUB publication, a
1024 fallback chain ending in a core media type resource SHOULD be provided even if the foreign
1025 resource is in an archivable format. (Note that this requirement is stricter than those in
1026 EPUB 3.x specifications, which require a fallback only in certain situations.)

1027 • The producer MAY include foreign resources (and metadata formats) in submitted EPUB
1028 publications if they have been specified as suitable for ingest and/or archivable in the
1029 submission agreement, or if their METS encoding in SIPs makes it possible to ignore them
1030 during ingest (see below).

1031 • If foreign resources and metadata are originally in un-archivable formats (formats that have
1032 not been specified as acceptable in the submission agreement), they SHALL be migrated
1033 during pre-ingest. The SIP may contain either just migrated publications, or both the
1034 original and migrated publications. Note that the preservation method MAY be either
1035 emulation or migration, so this requirement does not mean migration-only approach.

1036 • Core media types and foreign resources not specified in the submission agreement MAY be
1037 submitted if and only if the submission agreement allows it. METS encoding of these files in
1038 SIPs SHALL make it possible to skip their validation against the generic ingest criteria
1039 during the ingest process (since otherwise the SIP shall not pass the validation) and
1040 therefore passed directly to AIP. The specifics of this type of encoding SHALL be defined in
1041 the submission agreement.

1042 • If there are alternative versions (renderings) of the publication to be included in the SIP
1043 which are not archivable, they SHOULD be migrated into acceptable file formats prior to
1044 submission by the producer or a third-party preparing a SIP on behalf of the producer. For
1045 instance, if PDF is specified as not archivable but PDF/A is, the producer should create a
1046 PDF/A version of the document, which will then be submitted to the repository system
1047 alongside the EPUB publication of the same work.

1048 • If these non-archivable originals are submitted, their METS encoding in SIPs SHALL make it
1049 possible to skip validation against the generic ingest criteria during the ingest process
1050 (since otherwise the SIP would not pass) and therefore passed directly to AIP. The specifics
1051 of this type of encoding SHALL be defined in the submission agreement[36].

1052 NOTE   EPUB 3 Fixed Layout Properties

1053 In digital preservation the usual aim is to preserve intellectual content. Preserving also
1054 the original look and feel of the document is more challenging, although that may be
1055 required for some resources or collections. Reflowable EPUB publications are designed so

---

[35] 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.

[36] 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.

that their look and feel can change with no impact on semantics, which is a good thing from the digital preservation point of view, since in these documents EPUB content presentation adapts to the user preferences and display properties, which will change in the future.

In fixed layout EPUB publications 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 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. Therefore emulation of the original hardware and software environment is likely to be the best approach for preserving such documents.

Submission agreements SHOULD specify if submission of fixed-layout EPUB publications is allowed and if so, how they are treated during ingest. One solution is to include in SIPs also reflowable versions of these publications. If this is not possible or practical, SIP SHOULD contain metadata supporting emulation of the EPUB publication or publications in the package.

### 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 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.
- A SIP MAY contain multiple EPUB publications; for instance, different renderings of the same document[37]. If so, the SIP SHALL contain descriptive and administrative metadata which allows these publications to be ingested separately.
- A SIP MAY contain alternative renderings (such as PDF or DOCX) of the publication, but if so, the SIP SHALL contain all administrative metadata required for processing of these versions, and explaining the relations between these renderings.
- A single EPUB publication MAY be split into multiple SIPs if there is a valid reason to do so, such as the complexity or large size of the document.

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

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

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

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

### 6.1.9 The packaging method SHALL NOT limit the size of the SIP

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

## 6.2 Identification of information packages and their content

### 6.2.1 It SHALL be possible to identify any SIP uniquely both during and after the ingest process

- Since multiple SIPs may be submitted to the repository system simultaneously, there is a need to identify all packages in a (globally) unique manner. Identification will also make it possible to relate the packages with appropriate submitters, earlier submissions etc. Such identification helps to streamline the whole submission process and any potential communication between the archive and the submitting organization.
- Once the ingest process has been completed and 1-n AIPs have been formed, the SIP itself is no longer needed, but sometimes it is necessary to acquire more information about submitted publications from the producer, and SIP identifier is often necessary for that. Therefore both the SIP identifier and the AIP identifier(s) which the producer receives after the SIP has been ingested SHALL be persistent.
- There are circumstances in which AIP identifiers SHOULD be not only persistent, but also globally unique. For instance, an OAIS archive can cooperate with other archives by exchanging AIPs in order to share the bit level preservation costs.

---

[37] 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.

1143 • The entire SIP or parts of it SHALL be resubmitted in a revised format if the ingest process
1144 fails due to errors in the package. To keep track of the packages, SIPs SHALL have unique
1145 identifiers.

1146 **6.2.2 Information objects (EPUB publications, PREMIS preservation metadata record, etc.)**
1147 **within SIPs SHALL be identified uniquely and persistently**

1148 • Identifiers have many vital uses in digital preservation. They are used as access keys to the
1149 archived content in repository systems and facilitate information exchange with external
1150 systems. Identifiers also enable linking different versions of an archived document to each
1151 other. Moreover, with identifiers it is possible to link documents and
1152 descriptive/administrative metadata records that describe them. These links enable the
1153 archive to e.g. create dissemination information packages with the requested content.
1154 • Submission agreements SHALL specify identifier systems used, their location (EPUB
1155 document or SIP) and who is responsible of creating them (producer, archive or a third
1156 party). For instance, if the use of EPUB release identifiers is forbidden because the
1157 repository system does not support them, another means of identifying releases is needed.
1158 • International standard identifiers, such as ISBNs for books and DOIs for articles, SHALL be
1159 used as EPUB unique identifiers whenever possible. Any exceptions (such as using other
1160 identifier systems for releases which do not have ISBNs) SHOULD be specified in the
1161 submission agreement.
1162 • It SHOULD be possible to express the identifiers (also) as actionable HTTP URIs. Usage of
1163 persistent identifiers (Handles, DOIs, URNs, or ARKs) is recommended.
1164 • If there are multiple renditions of a work in an EPUB publication, requirements in the EPUB
1165 Multiple-Rendition Publications 1.0 specification SHALL be followed. Each rendition of an
1166 EPUB publication in a SIP SHALL have its own identifier.
1167 • The SIP SHOULD contain separate descriptive and administrative metadata records for each
1168 rendition, and these records SHALL have their own identifiers.

1169 NOTE 1 According to EPUB Multiple-Rendition Publications 1.0, the need to include more than one
1170 rendition of the content in an EPUB publication has grown as reading systems have become
1171 more sophisticated. In addition to optimizing the layout, adapting the content to specific
1172 reading systems may involve changing the content itself. Adaptation may also involve the
1173 prose of a textual work; instead of publishing several single-language EPUB publications
1174 multiple translations may be published as a single multiple-rendition EPUB publication.

1175 NOTE 2 Standard work identifier such as ISTC (International Standard Text Code) would be an ideal
1176 means of linking all manifestations to each other. Unfortunately there is no widely used
1177 standard identifier for textual works, and therefore this document does not require work
1178 level identification. However, if such identifier is available and supported in all applications
1179 involved, it is a good idea to use it. Work identifiers can be used to detect duplication of
1180 intellectual content in the archive, and if they are used in producers' and publishers' systems
1181 as well, it is possible to check overlap and possible gaps.

1182 **6.2.3 EPUB Fragment Identifiers SHOULD not be used in EPUB publications sent to a repository**
1183 **system, unless the submission agreement explicitly allows their use**

1184 • EPUB Canonical Fragment Identifiers define a standardized method of referencing content
1185 within an EPUB publication through the use of URI Fragments. From the digital
1186 preservation point of view, fragment identifiers can be problematic if the preservation
1187 strategy is not emulation, since URI fragments are media type dependent. Following
1188 migration the fragment identifiers may no longer be functional, because the new media type
1189 does not support them.

1190 • If fragment identifiers are allowed, the producer and the archive SHOULD take this into
1191 account in preservation planning, and design migrations so that the functionality provided
1192 by the fragment identifiers is preserved.

## 1193 6.3 Structure of information packages

### 1194 6.3.1 Submission information packages SHALL be built so that their components can be
### 1195 logically and physically separated from one another

1196 • For each rendition of the EPUB content document, there SHALL be a manifest file, which
1197 identifies and describes a set of resources that collectively compose a given rendition of a
1198 document, and EPUB spine, which provides a default reading order for a given rendition.
1199 • EPUB Open Container Format (OCF) defines a file format and processing model for
1200 encapsulating a set of related resources (for instance, renditions of the same resource) into
1201 a single-file (ZIP) EPUB Container[38].
1202 • The structure of each EPUB ZIP archive SHALL be described using the EPUB container.xml
1203 file (which describes the locations of root files of available renditions of the EPUB
1204 publication, and the rendition's package document and navigation document).
1205 • EPUB Package document and navigation document SHALL contain all metadata needed for
1206 rendering the publication, including the recommended reading system.
1207

## 1208 6.4 Generic Information package metadata

### 1209 6.4.1 Metadata in information packages SHALL be based on standards

1210 • METS or another agreed upon container format SHALL be used as the container standard,
1211 since this makes ingest to existing repository systems easier.
1212 • The submission agreement SHALL specify at least one mandatory metadata format for
1213 descriptive metadata. The format does not need to be Dublin Core; although EPUB
1214 publications always contain some Dublin Core metadata elements (see below), they MAY
1215 contain more complete metadata in another format, such as ONIX.
1216 • The minimum required descriptive metadata for EPUB publications are title, identifier, and
1217 language from the Dublin Core Metadata Element Set. Each rendition of a publication
1218 SHOULD also have at least the last modified date property from the DCMI Metadata Terms.
1219 Each rendition SHOULD also have the publication date encoded as DCMI Date, if the
1220 publication date is required to distinguish between publications.
1221 • SIPs submitted to a repository system MAY[39] contain preservation metadata, although such
1222 metadata will normally not be created in production systems, but in repository systems
1223 during ingest. PREMIS SHOULD be used for preservation metadata, as it is the most widely
1224 used and supported standard for this kind of metadata.
1225 • The submission agreement SHALL specify the syntax of metadata and its location (in the
1226 EPUB document, or in the SIP container), metadata formats used and metadata elements
1227 required or recommended.

---

[38] 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.

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

1228 • Since problems with text forms and encodings are common in repository systems, textual
1229 metadata SHOULD be provided also in Romanized form, using the EPUB alternate-script
1230 property to transcribe it if the metadata is originally in non-Roman script.

1231 **6.4.2 Metadata SHOULD allow (automatic) validation of the structure and content of SIPs in**
1232 **terms of integrity, fixity, and syntax**

1233 • SIPs SHALL contain message digests for all files of the SIP, and for the package itself.
1234 • File format identification and validation metadata (created with EpubCheck[40] or other
1235 validator tool) SHOULD be included in the SIP, if a validator is available.

1236 **6.4.3 It SHALL be possible to edit metadata in information packages**

1237 • If ingest has failed because of erroneous or missing metadata, the producer or a third-party
1238 responsible for the submission SHALL be able to modify the SIP so that it meets the
1239 metadata requirements in the submission agreement.
1240 • Producers and archives MAY use crowdsourcing and entity extraction activities to update
1241 descriptive metadata; an archive MAY choose to update this metadata also in the AIPs in the
1242 repository system although all the other components in the packages remain unchanged.

1243 **Annex A**

1244 (informative)

1245

1246 **EPUB and digital preservation: issues and recommendations**

1247 The British Library's EPUB Format preservation assessment includes a preservation risk summary
1248 [Whibley, p. 7-8]. The risks mentioned in the BL assessment are marked with [BL].

1249 **A.1 EPUB standard: issues**

1250 • Lack of stability in the e-book sector
1251    o EPUB does not have universally widespread support across e-book devices [BL].
1252 • Lack of EPUB format stability [BL]
1253    o Some EPUB versions have not been downward compatible. EPUB 3 is different from
1254        EPUB 3, and EPUB 3.1 from EPUB 3.0.1. The next version, 3.2, is based on 3.0.1, not on
1255        3.1.
1256    o Due to rapid technical development, future e-books (and EPUB specification) are likely
1257        to differ from the current one.
1258    o Proprietary changes and non-standard use of specifications have been used and will be
1259        used to restrict usage of EPUB publications to specific manufacturer hardware/software
1260 • Challenging EPUB features
1261    o From the long-term preservation point of view, the challenging features in EPUB include
1262        the possibility of using DRM, encryption and obfuscation, foreign resources and non-
1263        embedded resources.
1264    o While migration may be the best approach for most EPUB publications, interactive
1265        documents (containing software components) and fixed-layout documents are likely to
1266        be more suitable for preservation via emulation, so any OAIS archive preserving EPUB
1267        publications should be familiar with both preservation techniques.
1268 • Lack of archivable EPUB version
1269    o The standard is becoming richer and richer, and publishers and other users may find it
1270        more difficult to specify and avoid counterproductive features from the long-term
1271        preservation point of view. Pre-ingest (modifying the EPUB publication so that it can be
1272        preserved easily) may be difficult unless it has been taken into account from the
1273        beginning.

1274 Recommendations:

1275 • W3C should actively promote the EPUB format, because it is the only open e-book standard and
1276    it is based on open standards such as HTML5 and CSS.
1277 • EPUB community and digital preservation experts should develop a subset of EPUB ("EPUB/A")
1278    suited for long-term preservation.

1279 **A.2 EPUB usage: issues**

1280 • Ecosystem specific EPUB implementations
1281    o Major players in the e-book market (e.g. Amazon, Apple) have built EPUB based but
1282        closed (non-interoperable) ecosystems for e-books. E-books in vendor-specific formats,
1283        such as Amazon's KF8 should be migrated to EPUB before they are submitted to a
1284        repository system. Technically this is possible since EPUB is a "more or less obvious
1285        superset of what is possible in the different formats". The only exception is the fixed-

layout document specification in KF8; it is based on percentage information, not on absolute pixel positions as in EPUB 3. [Bläsi, p. 38].

- o These players have also created vendor-specific DRM solutions, which prevent the use of archived EPUB publications with other vendor's reading devices, unless the DRM protection has been removed during pre-ingest or ingest.

- Encryption and obfuscation [BL]
  - o Encryption may prevent the rendering of documents.
  - o Where not easily substituted, obfuscated fonts may lead to loss of critical information.

- Incomplete support in EPUB viewers [BL]
  - o Support for all aspects of the EPUB standard appears to be mixed, although impact of this is unclear. In the short-term, if the EPUB publication has been optimized for a specific reading system or systems, metadata embedded in the SIP should specify these systems. In the long-term, functionalities that are not widely supported may be lost.

- Losing information
  - o Where not easily substituted, non-embedded fonts may lead to loss of critical information.
  - o Metadata (and data) may not be embedded, but just linked to the SIP. During ingest, retrieval of linked information may fail.

- Invalid or badly formed EPUB files [BL]
  - o May affect the ability to render files now or in the future.

- Documents relying on EPUB features that may be difficult to preserve
  - o Fixed-layout documents: digital preservation usually concentrates on preserving the intellectual content, not the original look and feel of the document since that is regarded as difficult in the long-term. Preserving fixed-layout EPUB publications for the long-term may therefore be more demanding and require emulation instead of migration.
  - o Interactive documents that contain embedded applications supporting the required functionality may require a combination of migration and emulation methods for preservation; the former for intellectual content, the latter for software components.

- Legal issues [BL]
  - o It may be illegal to remove DRM, de-obfuscate embedded fonts, or to migrate the document to some other e-book format.

- Interactivity and animations
  - o With EPUB 3, there are two possibilities to realize built-in animations and interactive features. One is to use a CSS construct for transformations; another, more versatile approach is to use embedded JavaScript, Adobe Flash, or other software components that may enable complex interactive behaviour [Bläsi, p. 32]. Although EPUB 3 allows the use of JavaScript, it does not standardize the use of JavaScript elements in e-books. This can easily lead to proprietary extensions as well as incompatible EPUB 3 reading systems that support a different or incompatible subset of scripting elements [Bläsi, p. 17].
  - o Different e-book formats support interactivity in different ways, and apart from EPUB, features may be undocumented. Therefore migrating interactivity features between e-book formats is either difficult or impossible.

- Non-archivable core media types
  - o Depending on the chosen preservation strategy, some current or future core media types may be regarded as unsuitable for digital preservation. For instance, GIF is not an archivable format according to the requirements of the Finnish National Digital Library initiative. [File formats, p.25].

- Non-archivable foreign resources
  - o Foreign content may be both non-archivable and unsupported by EPUB reading systems the archive is able to use.

1337       o  For the time being, there are no video codecs among the core media types. There is a
1338         recommendation that reading systems should support either H.264 or VP8. Neither of
1339         these are archivable or even ingestible formats in the Finnish National Digital Library
1340         specification, which approves JPEG 2000 sequence and MPEG-4 AVC as archive formats
1341         and DV (Digital Video), MPEG-1, MPEG-2, and WMV (Windows Media Video) as
1342         ingestible formats.

1343    •  External references [BL]
1344       o  Externally referenced content (metadata, core media types, or foreign resources) SHALL
1345         be retrieved during pre-ingest and embedded into the SIP, during ingest and embedded
1346         into the AIP. If retrieval fails, the AIP is incomplete. If the submission agreement allows
1347         such policy, the archive can store the incomplete AIP and try to retrieve the missing
1348         content post-ingest. If the second attempt is successful, the AIP is ingested again into the
1349         repository system, and the missing content is added.

1350    •  Missing or poor fallback documents
1351       o  If a foreign resource cannot be rendered, there SHOULD be a core media type fallback
1352         document. However, even if a fallback resource is present it may not produce the same
1353         rendition than the original resource and there is no guarantee either that the original
1354         semantics will be preserved.

1355  Recommendations:

1356    •  EPUB 3 covers the superset of the expressive abilities of all the other e-book formats. Therefore
1357      there is no technical or functional reason not to use and establish EPUB 3 as an interoperable
1358      open e-book format standard [Bläsi, p. 8]. Having a universally supported e-book format would
1359      benefit current e-book users and make long-term preservation of e-books easier.

1360    •  Readium[41] project is developing a robust and efficient reader for EPUB publications. Such tools
1361      will make it easier to use rich EPUB documents, and EPUB community should continue
1362      investments on Readium and similar initiatives.

1363    •  The EPUB community should create EPUB/A, a subset of EPUB 3 with features suitable for long-
1364      term preservation. The specification should be complemented by an explanation why the EPUB
1365      3 features not included in the EPUB/A format may jeopardize digital preservation, and a
1366      justification for those featureas that are required.

1367    •  When new EPUB core media types are added, the archivability of these file formats should be
1368      taken into account. EPUB community could co-operate with the digital preservation community
1369      to achieve this goal.

1370    •  Legal aspects of long-term preservation of EPUB 3 documents should be investigated.

1371    •  Open source licenses such as SIL Open Font License[42] should be used when possible.

1372    •  Foreign resources should be used with caution, until the archivability of the utilized file formats
1373      has been verified.

1374    •  Core media types that are considered to be non-archivable should be avoided whenever
1375      possible. For instance, it is better to use a JPEG or a PNG than a GIF image.

1376

1377

---

[41] http://readium.org/
[42] http://scripts.sil.org/OFL_web

# Bibliography

[1]     BLÄSI, Christoph and Franz Rothlauf: *On the interoperability of eBook formats*. [online]. Brussels: European and International Booksellers Federation, 2013. Available from: http://wi.bwl.uni-mainz.de/publikationen/InteroperabilityReportGutenbergfinal07052013.pdf [viewed 2017-06-12].

[2]     DALY, Liza: *EPUB 3 and interactivity*. [online]. EPUBZone, 2014. Available from: http://epubzone.org/news/epub-3-and-interactivity [viewed 2017-06-21].

[3]     WHIBLEY, Simon: *EPUB format preservation assessment*. Version 1.2. [online]. British Library, 2015. Available from: http://wiki.dpconline.org/images/a/a9/EPUB_Assessment_v1.2.pdf [viewed 2018-03-14].

[4]     *Digital Preservation Handbook*. 2nd, revised ed. [online]. Digital Preservation Coalition, 2017. Available from: http://www.dpconline.org/handbook. [viewed 2017-07-25].

[5]     E-ARK: *Common specification for information packages. Version 1.0* [online]. E-ARK Project, 2016. Available from: http://www.dasboard.eu/specifications/common-specification. [viewed 2017-08-23].

[6]     *EPUB Publications 3.0.1.* International Digital Publishing Forum, 2014 [online]. Available from: http://www.idpf.org/epub/301/spec/epub-publications.html [viewed 2018-10-30].

[7]     *File formats. Version 1.6.1.* [online]. The National Digital Library, 2017. Available from: http://digitalpreservation.fi/files/File-Formats-1.6.1-en.pdf. [viewed 2018-09-10].

[8]     ISO 14721. *Space data and information transfer systems – Open archival information system (OAIS) – Reference model*. [online]. ISO, 2012. Available from: https://www.iso.org/standard/57284.html [paywall]. [viewed 2017-07-10]. Also available as CSSDS specification from: https://public.ccsds.org/pubs/650x0m2.pdf [viewed 2017-07-17].

[9]     ISO/FDIS 20614. *Information and documentation – Data exchange protocol for interoperability and preservation*. [online]. ISO, 2017. The standard, when completed, will be available from: https://www.iso.org/standard/68562.html. [viewed 2017-07-17].

[10]    ISO 20652. *Space data and information systems - Producer-archive interface - Methodology abstract standard [PAIMAS]*. [online]. ISO, 2006. Available from: https://www.iso.org/standard/39577.html. [paywall]. [viewed 2017-07-01]. Also available as CCSDS 651.0-M-1:2004 from: https://public.ccsds.org/publications/archive/651x0m1.pdf [viewed 2017-06-22].

[11]    LAVOIE, Brian: *Meeting the challenges of digital preservation: The OAIS reference model*. [online]. Dublin, OH: OCLC Research, 2000. Available from: http://www.oclc.org/research/publications/library/2000/lavoie-oais.html [viewed 2017-07-17].

[12]    *PDF/A* [online]. Wikipedia, 2017-05.17. Available from: https://en.wikipedia.org/wiki/PDF/A [viewed 2017-07-18]. Archived version available from: https://web.archive.org/web/20170713180152/https://en.wikipedia.org/wiki/PDF/A

[13]    *TI/A Standard Initiative homepage*. [online]. TI/A Standard Initiative, 2016. Available from: http://ti-a.org/ . [viewed 2017-07-18].

[14]    W3C. *EPUB 3 Community Group Charter*. [online]. W3C, 2017. Available from: https://www.w3.org/2017/02/EPUB3CGcharter. [viewed 2017-07-24].