
**Digital publishing — EPUB3
preservation —**

**Part 1:
Principles**

*Publications numériques — EPUB3 preservation —
Partie 1: Principes*

IECNORM.COM : Click to view the full PDF of ISO/IEC TS 22424-1:2020



IECNORM.COM : Click to view the full PDF of ISO/IEC TS 22424-1:2020



COPYRIGHT PROTECTED DOCUMENT

© ISO/IEC 2020

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

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

Published in Switzerland

Contents

	Page
Foreword	iv
Introduction	v
1 Scope	1
2 Normative references	1
3 Terms and definitions	1
4 Abbreviated terms	9
5 Packaging standards	9
6 Construction of OAIS information packages	11
6.1 Overview	11
6.2 General principles	12
6.2.1 EPUB publications shall be sent to a repository system as well-formed and complete submission information packages (SIPs)	12
6.2.2 Regardless of its type or format, it shall be possible to include any data or metadata in SIPs	14
6.2.3 It should be possible to transfer SIPs by any means, methods, or tools from the submitting organization to the repository system	16
6.2.4 The archive shall have a way to verify the identity of the submitting organization/person, no matter how the information packages are transferred	16
6.2.5 There is no 1:1 relation between OAIS information packages	16
6.2.6 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	16
6.2.7 The information package type (in this case, SIP) shall be indicated	16
6.2.8 SIP packaging method shall not restrict the application of any preservation method	17
6.2.9 The packaging method shall not limit the size of the SIP	17
6.3 Identification of information packages and their content	17
6.3.1 It shall be possible to identify any SIP uniquely both during and after the ingest process	17
6.3.2 Information objects (EPUB publications, PREMIS preservation metadata record, etc.) within SIPs shall be identified uniquely and persistently	17
6.3.3 EPUB Fragment Identifiers should not be used in EPUB publications sent to a repository system, unless the submission agreement explicitly allows their use	18
6.4 Structure of information packages	18
6.5 Generic Information package metadata	19
6.5.1 Metadata in information packages shall be based on standards	19
6.5.2 Metadata should allow (automatic) validation of the structure and content of SIPs in terms of integrity, fixity, and syntax	19
6.5.3 It shall be possible to edit metadata in information packages	19
Annex A (informative) EPUB and digital preservation: issues and recommendations	20
Bibliography	24

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.

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) or the IEC list of patent declarations received (see <http://patents.iec.ch>).

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

For an explanation of 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 www.iso.org/iso/foreword.html.

This document was prepared by Joint 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.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at www.iso.org/members.html.

Introduction

0.1 General

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 document 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 to ISO 19005-3, is related to the Portable Document Format (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;
- packaging 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 Open Archival Information System (OAIS) submission information package (SIP).

Packaging is covered in ISO/IEC TS 22424-2.

0.2 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.^[17]

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 Group¹⁾, 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).^[23]

The International Digital Publishing Forum, IDPF, has ceased operations as a membership organization in January 2017, and its website²⁾ 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 3.0.1³⁾. EPUB 3.1⁴⁾ 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 May 2019⁵⁾. Unlike 3.1, it is fully backwards compatible with 3.0.1. It will be covered in the next edition of this document.

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/>

5) <https://w3c.github.io/publ-epub-revision/epub32/spec/epub-spec.html>

Differences between EPUB specifications 2.0.1-3.2 are well documented:

- EPUB 3 Changes from EPUB 2.0.1⁶⁾
- EPUB 3.0.1 Changes from EPUB 3.0⁷⁾
- EPUB 3.2 Changes from EPUB 3.0.1⁸⁾

All EPUB specifications are available in the Web; 2.0.1 at <http://idpf.org/epub/201>, EPUB 3.0.1 at <http://idpf.org/epub/301> and 3.2 at <https://w3c.github.io/publ-epub-revision/epub32/spec/epub-spec.html>.

All EPUB publications, including ones using version 3.2, can be validated using EPUBCheck version 4.2.0, which was released in March 2019.

From long-term preservation point of view, lack of backward compatibility between successive versions of a file format would be a problem because it makes migration more challenging. In addition, EPUB 3.1 has at least one feature which would have been problematic. In EPUB 3.1 foreign resources do not require fallbacks if they are not in the spine and not embedded in EPUB Content Documents. In EPUB 3.0.1, fallback guarantees that there is a version of the document that can be rendered; in 3.1 such guarantee no longer exists.

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

- EPUB 3 Overview
- Publications 3.0.1
- Canonical fragment identifiers
- Content documents 3.0.1
- Media overlays 3.0.1
- Open Container Format 3.0.1

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

- EPUB Accessibility specification 1.0⁹⁾ addresses evaluation and certification of accessible EPUB publications, and discovery of the accessible qualities in such publications.
- EPUB Previews 1.0¹⁰⁾ describes how content previews can be included in EPUB publications.
- EPUB Distributable Objects 1.0¹¹⁾ is a draft specification that defines a method for the encapsulation, transportation, and integration of distributable objects in EPUB publications.
- EPUB Scriptable Components 1.0¹²⁾ provides an interoperable publish and subscribe (pubsub) pattern by which interactive content can be created and incorporated into EPUB publications. Same as EPUB Distributable Objects, it is as of 2019-05-13 a draft.

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/>

- EPUB Scriptable Components Packaging and Integration 1.0¹³⁾ is a draft that defines a method for the creation and inclusion of dynamic and interactive components in EPUB publications.
- EPUB Multiple-Rendition Publications 1.0¹⁴⁾ defines the creation and rendering of EPUB publications consisting of more than one rendition of the same publication.
- EPUB Dictionaries and Glossaries 1.0¹⁵⁾ provides a means for expressing dictionary and glossary semantics in EPUB publications.

These extensions are not widely used and they have not been explicitly taken into account in this document. As regards accessibility, all EPUB publications are supposed to be accessible. However, accessibility features as such do not have an impact on long term preservation of EPUB publications and therefore this document does not make accessibility-related requirements.

EPUB 3 core media types have been listed at <https://www.w3.org/publishing/epub3/epub-spec.html#sec-core-media-types>. As of 2019-05-13, the latest change has been made on April 1, 2018. Starting from EPUB 3.2, core media types are part of the standard.

In 2014, EPUB 3.0 specifications were republished as ISO/IEC TS 30135-1 to ISO/IEC TS 30135-6. Each of these six ISO specifications is identical to its IDPF equivalent, for example ISO/IEC TS 30135-1 has exactly the same content as the EPUB 3.0 Overview.

ISO/IEC TS 30135-7 entitled "Part 7: EPUB3 Fixed-Layout Documents" is from EPUB 3.0.1 (EPUB 3.0 does not have fixed layout specification). ISO/IEC TS 30135 (all parts) is therefore a combination of EPUB 3.0 and Fixed-Layout Documents specification from 3.0.1.

ISO/IEC JTC 1/SC 34 is currently updating the ISO standard to match fully the version 3.0.1.

EPUB is a rich document format with a lot of features. From the digital preservation point of view this is a challenge, not least because long-term preservation has not been a priority in the development of the standard. Preserving all aspects and features of EPUB publications may be difficult, since there are features which are difficult to preserve. Moreover, EPUB reading systems usually do not support all features of the specification and finding tools supporting rare features can be difficult.

In spite of these challenges EPUB is generally regarded as a suitable format for digital archiving. For instance, the Finnish National Digital Library initiative has selected just eight archivable file formats for text, EPUB being one of them. The selection criteria were openness/transparency, adoption as a preservation standard, degree of forward/backward compatibility, degree of protection against file corruption, frequency of version releases, dependencies/interoperability, and standardization. EPUB got an A, the best grade, from everything else except the second and third criterion. For those, the grade was the second best, a B (see Reference [19], p.40). Based on these generic criteria, EPUB seems to provide a good basis for long-term preservation, although additional guidelines on how to use the standard are needed to guarantee EPUB files can be preserved efficiently.

The British Library's Digital Preservation Team has published an assessment of EPUB as a preservation format^[15]. It covers EPUB versions 3.0.1 and 2 and the overall view of EPUB is positive (Reference [15], p.2):

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

EPUB long-term preservation issues uncovered in the assessment of the British Library are discussed in [Annex A](#).

EPUB is enjoying reasonable support in the e-book market. Many suppliers, publishers, and application developers who have supported EPUB 2 have implemented version 3.0.1. According to the EPUBTest web

13) <http://www.idpf.org/epub/sc/pkg/>

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

15) <http://www.idpf.org/epub/dict/>

site¹⁶⁾, EPUB 3 support in reading systems is far from exhaustive, but market coverage is good – in January 2018, there were 59 reading systems supporting at least some of the features specified in EPUB 3.0.

E-book suppliers have produced EPUB 3 based formats that incorporate digital rights management (DRM), and EPUB modifications that may restrict using the format on other than the suppliers' own platforms. For example, the Kindle Fire eReader, released in 2015, uses a new format called Kindle Format 8 (KF8), which is partly based on EPUB 3, with Amazon's DRM. See Reference [15], 3. Publisher/supplier specific DRM often restricts the use of e-books to that publisher's/supplier's rendering devices and/or applications, and is therefore a major obstacle to digital preservation (see Reference [15], p.7).

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

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

If preserved EPUB publications are not directly accessible by the public, removing DRM, digital watermarking, and other protection mechanisms from the archived documents is not a risk. When publications are delivered to the customers as dissemination information packages (DIPs), the archive shall use a combination of administrative and technical means to protect the documents as required in the submission agreement. These means may include adding DRM protection mechanism into the DIP submitted to the user according to the requirements of the submission agreement. The agreement may also specify the customers the archive is entitled to serve; for instance, it is possible to require that the preserved documents can only be disseminated to the producer, and the producer will serve the end-users who do not have direct access the OAIS archive.

0.3 Digital preservation

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

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

- we need software and hardware to render content for human users;
- software changes over time; there are new versions from old applications, and entirely new applications;
- new or updated applications may not be able to render outdated file formats or format versions correctly
- digital preservation makes an effort to have all archived content in stable formats. Publications should also contain the smallest possible amount of features which are not commonly supported in software packages used to render the content in these formats, and also avoid adding links to external resources since then the long-term access to the publication requires also persistence of these external resources.
- when necessary, data in old formats may be migrated into more modern formats or updated versions of the same format. For instance, an e-book in EPUB 3.0.1 format may be migrated to EPUB 3.2. when version 3.0.1 is no longer widely supported by reading systems.

16) <http://epubtest.org/results>

- since the aim is to preserve the content, not the bits, the bits may change as a result of version updates and format migrations.
- Many OAIIS archives preserve successive versions of archives publications, because migration may change the look and feel of the original document, or even its intellectual content.

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

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

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

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

PDF/A forbids also the use of encryption, because encryption is generally regarded as a risk for long-term preservation. But storing unencrypted documents is a risk as well, because if they are stolen, non-authorized usage is easy. Therefore, according to the Digital preservation handbook^[25]:

Information security methods such as encryption add to the complexity of the preservation process and should be avoided if possible for archival copies. Other security approaches may therefore need to be more rigorously applied for sensitive unencrypted files; these might include restricting access to locked-down terminals in controlled locations (secure rooms), or strong user authentication requirements for remote access.

In order to guarantee the correct processing of PDF/A files, there are specific requirements for PDF/A reading systems, such as support for embedded fonts. There are three versions of the specification: 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 all the features of PDF/A-2 as well as allows the embedding of other file formats into PDF/A conforming documents^[21].

The TI/A (Tagged Image for Archival) standard initiative intended to create an ISO recommendation to optimize the format specification for archival purposes. Unfortunately the project was disbanded in 2016, and the TI/A draft the initiative completed in September 2016 is only available in the project Intranet. However, the original TIFF/A (later TI/A) draft from February 2015 is a public document available on a PREFORMA project web site^[17]. Although this TIFF/A specification is only a draft, it is probably a good idea to use in archival TIFF images features specified mandatory in the specification, and avoid the ones which are forbidden.

The motivation behind the TI/A initiative can be applied to other image formats as well, and there are also points the EPUB community might agree with Reference ^[22]:

17) <http://www.preforma-project.eu/dpf-manager.html>

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

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

TIFF is not an EPUB core media type, but four other image types have been listed; GIF, JPEG, PNG, and SVG. It is significant from a digital preservation point of view how these formats and other core media types are used in the EPUB context. Image and audio files embedded in an EPUB publication may require migration before the EPUB publication itself has to be migrated into a more modern file format, if commonly available EPUB reading systems no longer support these file formats. This document does not provide guidelines for creating archivable files in EPUB 3 core media types, due to the magnitude of such task. But EPUB community should follow the archival file format lists of national archives or libraries (for example the Library of Congress file format list¹⁸⁾ and the U.S. National Archives list¹⁹⁾) when the core media file format list is updated. Publishers should also consider the persistence of file formats used when creating EPUBs for which the need for long-term preservation is foreseen.

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

0.4 OAIS and related standards

ISO 22424 (all parts) provides guidance on how to utilize the OAIS and current practices of OAIS archives in preservation of EPUB publications. The OAIS (ISO 14721) is equally relevant to both parts of the ISO 22424 series.

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

The model has five functional units (Ingest, Archival Storage, Access, Data management and Administration) as shown in [Figure 1](#).

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>

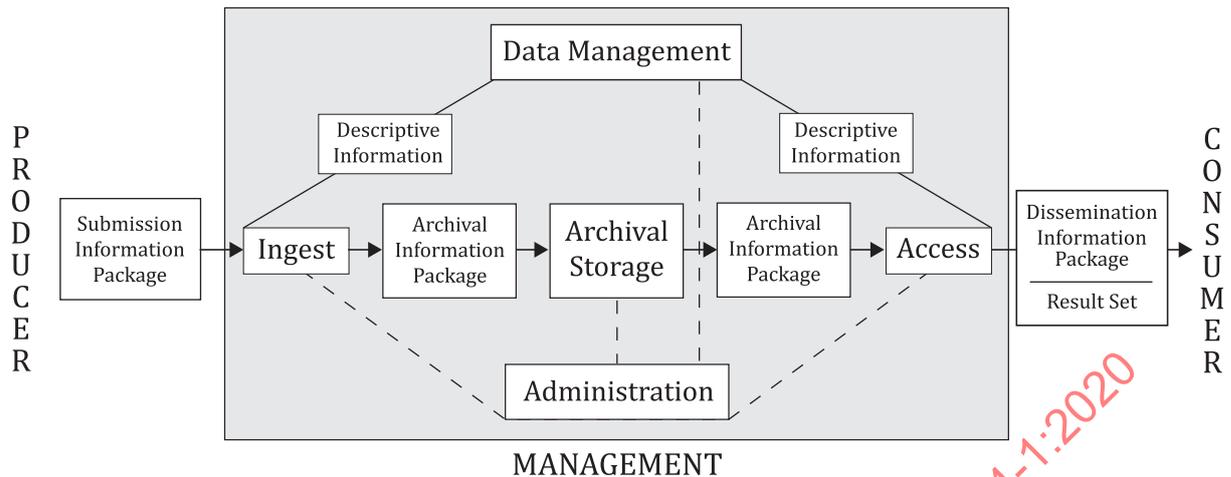


Figure 1 — OAIS model^[20]

In the model, the *ingest function* is responsible for receiving information from producers and preparing it for storage and management within the OAIS archive. The ingest accepts information – in this case, EPUB publications – from producers in the form of SIPs, performs quality assurance checks on the SIP, and generates an archival information package (AIP) from one or more SIPs (or multiple AIPs from a single SIP). Finally, the ingest function transfers the new AIPs to Archival Storage and the associated descriptive information (metadata) to Data Management.

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

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

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

OAIS submission agreements specify the principles of how documents should be prepared and submitted to the repository system. If the archive uses migration as the preservation method²¹⁾, submission agreements should specify file formats (and metadata formats) suitable for submission and/or archival, or refer to external documents listing these formats. File formats suitable for submission but not for archival are migrated during the ingest process, although the original files may be included in the AIP.

21) In this document, preservation method is assumed to be migration. In practice, emulation can 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 is not the case. Migrated document could 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.

Therefore archives often preserve also the original version of the archived resource, alongside more modern versions.

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

This document outlines the general principles for the submission of EPUB publications from digital asset management systems to repository systems. The principles of archival storage or dissemination of archived documents are not covered here, because OAIS archives may apply various methods and processes to meet the requirements of submission agreements. Bit level preservation is also out of scope; the purpose of this document is to make it easier for producers and OAIS archives to preserve access to EPUB documents.

ISO/IEC TS 22424-2 provides a technical basis to meet the principles listed in this document by specifying metadata required for long-term preservation, and a method for packaging this metadata with the original EPUB container.

This document is applicable to EPUB versions 3 and 3.0.1 and as such it should be used cautiously with other (previous or later) versions of the standard. If there is a need to preserve documents that are in earlier EPUB versions, they do not need to be migrated, provided that a) submission agreement specifies those EPUB versions as archivable formats, and b) there are reading systems for these EPUB versions. Additional features in future EPUB versions should be analyzed from a long-term preservation point of view. If such an analysis reveals that they may constitute a risk, they should be avoided in submitted EPUB publications, or removed during ingest.

[Annex A](#) in this document provides a summary of issues and recommendations related to the EPUB standard and its usage from long-term preservation point of view.

Digital publishing — EPUB3 preservation —

Part 1: Principles

1 Scope

The ISO/IEC TS 22424 series supports long-term preservation of EPUB publications via a dual strategy. This document considers EPUB features from a long-term preservation point of view. 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 are suitable for preservation.

ISO/IEC TS 22424-2 makes EPUB compliant with Open Archival Information System (OAIS) and current practices of OAIS archives.

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 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 terms and definitions given in ISO 14721 and the following apply.

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

- ISO Online browsing platform: available at <https://www.iso.org/obp>
- IEC Electropedia: available at <http://www.electropedia.org/>

NOTE Unless stated otherwise, the terms have been adopted from ISO 14721:2012.

3.1

administrative metadata

metadata (3.33) that provides *information* (3.28) to help manage a resource, such as when and how it was created, file type and other technical information, and access rights

Note 1 to entry: The definition is adapted from Reference [24].

3.2

archival information package

AIP

information package (3.29) consisting of *content information* (3.6) and associated *preservation description information (PDI)* (3.41), which is preserved within an *OAIS* (3.36)

3.3

archive

OAIS archive

organization that intends to preserve *information* (3.28) for access and use by a *designated community* (3.11)

**3.4
authenticity**

property that an entity is what it claims to be

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

[SOURCE: ISO/IEC 27000:2018, 3.6, modified — Note 1 to entry has been added.]

**3.5
consumer**

role played by those persons or client systems, who interact with *OAIS* (3.36) services to find preserved *information* (3.28) of interest and to access that information in detail

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

**3.6
content information**

set of *information* (3.28) that is the original target of preservation or that includes part or all of that information

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

**3.7
context information**

information (3.28) that documents the relationships of the *content information* (3.6) to its environment

Note 1 to entry: This includes reasons why the content information was created and how it relates to other content information objects.

**3.8
core media type**

set of *publication resource* (3.45) for which no *fallback* (3.23) is required

Note 1 to entry: The definition is adapted from Reference [18]. Core media types have been specified in chapter 5.1 of Reference [18].

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

**3.9
data, pl**

reinterpretable representation of *information* (3.28) in a formalized manner suitable for communication, interpretation, or processing

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

[SOURCE: ISO 5127:2017, 3.1.1.15]

**3.10
descriptive metadata
descriptive information**

metadata (3.33) about a resource for example for discovery and identification

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

Note 2 to entry: The definition is adapted from Reference [24].

3.11**designated community**

identified group of potential *consumers* (3.5) who should be able to understand a particular set of *information* (3.28)

Note 1 to entry: A designated community may be composed of multiple user communities. The community is defined by an *archive* (3.3), though this definition may change later on.

3.12**digital preservation**

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

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

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

Note 3 to entry: The definition is adapted from Reference [25].

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

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

EXAMPLE 3 **Long-term preservation** (3.31) - Access to digital materials, or at least to the *information* (3.28) contained in them, indefinitely.

3.13**digital rights management****DRM**

packaging, distributing, controlling, and tracking content based on rights and licensing *information* (3.28)

3.14**digital signature****signature**

data (3.9) appended to, or a cryptographic transformation of, a data unit that allows the recipient of the data unit to prove the source and integrity of the data unit and protect against forgery, e.g. by the recipient

[SOURCE: ISO/IEC 19784-1:2018, 4.34, modified — Note 1 to entry has been removed.]

3.15**dissemination information package****DIP**

information package (3.29), derived from one or more *AIPs* (3.2), sent by an *archive* (3.3) to a *consumer* (3.5) in response to a request in the *OAIS* (3.36)

3.16**distributable object**

component of an *EPUB publication* (3.21) that can be reused in other contexts

Note 1 to entry: A distributable object can be a complete *EPUB content document* (3.19) (e.g., a chapter of a book), a section of such a document (e.g., an exercise or a promotional excerpt), a media resource (e.g., a video or interactive feature), or a combination of such resources that are not necessarily contiguous within the parent EPUB publication but are intended to be able to be distributed as a unit.

Note 2 to entry: The definition is adapted from Reference [26].

3.17

e-book

electronic book

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

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

[SOURCE: ISO 2789:2013, 2.3.19, modified — "eBook" has been replaced by "e-book"; "(monograph)" at the end of the definition has been removed; Note 2, 3 and 4 to entry have been removed.]

3.18

EPUB container

ZIP based packaging and distribution format for *EPUB publications* (3.21)

Note 1 to entry: The definition is adapted from Reference [18].

3.19

EPUB content document

publication resource (3.45) that conforms to one of the EPUB content document definitions

Note 1 to entry: The definition is adapted from Reference [18].

3.20

EPUB navigation document

specialization of the *XHTML content document* (3.58), containing human- and machine-readable global navigation *information* (3.28)

Note 1 to entry: The definition is adapted from Reference [18].

3.21

EPUB publication

collection of one or more *renditions* (3.52) conforming to the EPUB specifications, packaged in an *EPUB container* (3.18)

Note 1 to entry: The definition is adapted from Reference [18].

3.22

EPUB reading system reading system

system that processes *EPUB publications* (3.21) for presentation to a user in a manner compliant with EPUB specifications

Note 1 to entry: The definition is adapted from Reference [18].

3.23

fallback

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

Note 1 to entry: A *reading system* (3.22) that does not support the file format of a *foreign resource* (3.25) shall traverse the fallback chain until it finds a version it can render.

Note 2 to entry: The definition is adapted from Reference [18].

3.24

fixity information

information (3.28) that documents the authentication mechanisms and provides authentication keys to ensure that the *content information* (3.6) object has not been altered in an undocumented manner

[SOURCE: ISO 13527:2010, 1.4.2, modified — The example has been removed.]

3.25**foreign resource**

publication resource (3.45) that is not a *core media type* (3.8)

Note 1 to entry: The definition is adapted from Reference [18].

3.26**identifier**

data (3.9) string or pointer that establishes the identity of an item, institution, or person alone or in combination with other elements

Note 1 to entry: EPUB 3 specifies *unique identifiers* (3.57) and *release identifiers* (3.50); the latter is a combination of a unique identifier and the last modification data of the *rendition* (3.52) of the resource.

[SOURCE: ISO 8459:2009, 2.27, modified — Note 1 to entry has been added.]

3.27**independently understandable**

characteristic of *information* (3.28) that is sufficiently complete to allow it to be interpreted, understood, and used by the *designated community* (3.11) without having to resort to special resources not widely available, including named individuals

3.28**information**

any type of knowledge that can be exchanged

Note 1 to entry: In an exchange, this is represented by *data* (3.9).

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

3.29**information package**

logical container composed of optional *content information* (3.6) and optional associated *preservation description information* (3.41)

3.30**long-term**

period of time long enough to raise concerns about the impact of changing technologies, including support for new media and *data* (3.9) formats, and of a changing *designated community* (3.11), on the *information* (3.28) being held in an *OAIS* (3.36)

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

3.31**long-term preservation**

act of maintaining *information* (3.28), *independently understandable* (3.27) by a *designated community* (3.11), with evidence supporting its *authenticity* (3.4) over the *long-term* (3.30)

3.32**manifest**

EPUB manifest element that provides an exhaustive list of the *publication resources* (3.45) that constitute the given *rendition* (3.52), each represented by an item element

Note 1 to entry: The definition is adapted from Reference [18].

3.33

metadata

data (3.9) about other data, documents, or records that describe their content, context, structure, format, provenance, and/or rights

[SOURCE: ISO 5127:2017, 3.1.10.26.01, modified — Note 1 to entry has been removed; "describes" has been replaced by "describe"; "data format" has been replaced by "format"; the words "attached to them" at the end of the definition has been removed.]

3.34

METS

Metadata Encoding and Transmission Standard

standard for presenting *metadata* (3.33) using XML

Note 1 to entry: The definition is adapted from Reference [25].

3.35

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* (3.49) 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.

Note 3 to entry: The definition is adapted from Reference [25].

3.36

Open Archival Information System

OAIS

archive (3.3), 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* (3.28) and make it available to a *designated community* (3.11)

Note 1 to entry: It has a set of responsibilities, as defined in [Clause 4](#), which allow an OAIS archive to be distinguished from other uses of the term 'archive'.

Note 2 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 3 to entry: The OAIS abbreviation is also commonly used to refer to the Open Archival Information System *reference model* (3.47) 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* (3.31).

3.37

package document

publication resource (3.45) that describes one *rendition* (3.52) of an *EPUB publication* (3.21)

Note 1 to entry: The package document carries meta *information* (3.28) about the rendition, provides a *manifest* (3.32) of resources and defines the default reading order.

Note 2 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.

Note 3 to entry: The definition is adapted from Reference [18].

3.38**PDF****Portable Document Format**

set of formats and open standards maintained by the International Organization for Standardization for producing and sharing electronic documents

Note 1 to entry: An open format was originally developed by Adobe Systems.

Note 2 to entry: The definition is adapted from Reference [25].

3.39**PDF/A**

versions of the *PDF* (3.38) standard intended for archival use

Note 1 to entry: The definition is adapted from Reference [25].

3.40**pre-ingest**

actions required before *data* (3.9) can be submitted into an *OAIS archive* (3.3), including negotiation of data acquisitions, checking rights and access criteria, licensing, and data submission

Note 1 to entry: This area also includes activities involving data *producer* (3.43) support and training.

Note 2 to entry: Pre-ingest is not a function in the standard *OAIS* (3.36) model, but activities in this area can form a significant part of a producer's responsibilities.

Note 3 to entry: The definition is adapted from Reference [27].

3.41**preservation description information****PDI**

information (3.28) necessary for the adequate preservation of *content information* (3.6) that can be categorized as provenance, reference, fixity, context, and rights information

3.42**preservation metadata**

metadata (3.33) containing *information* (3.28) needed to archive and preserve a resource

Note 1 to entry: The definition is adapted from Reference [24].

3.43**producer**

role played by those persons or client systems that provide the *information* (3.28) to be preserved

Note 1 to entry: This can include other *OAISs* (3.36) or internal *OAIS* persons or systems. The producer does not need to be the publisher.

3.44**provenance information**

information (3.28) that documents the history of the *content information* (3.6)

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* (3.3) is responsible for creating and preserving provenance information from the point of ingest; however, earlier provenance information should be provided by the *producer* (3.43). Provenance information adds to the evidence to support *authenticity* (3.4).

**3.45
publication resource**

resource that has the content or instructions contributing to the logic and rendering of at least one *rendition* (3.52) of an *EPUB publication* (3.21)

Note 1 to entry: The definition is adapted from Reference [18].

EXAMPLE Examples of publication resources include a *rendition's package document* (3.37), *EPUB content document* (3.19), EPUB style sheets, audio, video, images, and embedded fonts and scripts.

**3.46
reference information**

information (3.28) that is used as an *identifier* (3.26) for the *content information* (3.6)

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.

**3.47
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.

**3.48
reformatting**

copying *information* (3.28) 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)

Note 1 to entry: The definition is adapted from Reference [25].

**3.49
refreshing**

copying *information* (3.28) content from one storage media to the same storage media

Note 1 to entry: The definition is adapted from Reference [25].

**3.50
release identifier**

identifier (3.26) that allows any instance of an *EPUB publication* (3.21) to be compared against another to determine if they are identical, different versions, or unrelated

Note 1 to entry: Release identifiers consist of a *unique identifier* (3.57) and the last-modified date of the document.

Note 2 to entry: The definition is adapted from Reference [18].

**3.51
remotely-hosted resource**

objects hosted outside the *EPUB container* (3.18)

**3.52
rendition**

rendering of the content of an *EPUB publication* (3.21), as expressed by an EPUB package

Note 1 to entry: The definition is adapted from Reference [18].

**3.53
repository system**

long-term preservation (3.31) system used by an *archive* (3.3)

3.54 spine

EPUB element that defines the default reading order of the *EPUB publication* (3.21) content by defining an ordered list of *manifest* (3.32) item references

Note 1 to entry: The definition is adapted from Reference [18].

3.55 submission agreement

agreement reached between an *OAIS archive* (3.3) and a *producer* (3.43) that specifies a *data* (3.9) model and any other arrangements needed for the data submission session

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

3.56 submission information package SIP

information package (3.29) that is delivered by a *producer* (3.43) to an *OAIS* (3.36) to be used to construct or update one or more *AIPs* (3.2) and/or the associated *descriptive information* (3.10)

3.57 unique identifier

primary identifier (3.26) of an *EPUB publication* (3.21), which may be shared by one or several *renditions* (3.52) of the same EPUB publication that conform to the EPUB standard and embody the same content

Note 1 to entry: The definition is adapted from Reference [18].

3.58 XHTML content document

EPUB content document (3.19) that conforms to the profile for HTML defined in XHTML content documents

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

Note 2 to entry: The definition is adapted from Reference [18].

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

5 Packaging standards

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

The OAIS archival storage function stores, maintains, and retrieves AIPs. Maintenance may include for instance frequent error checks to protect the data against bit rot. In order to keep the documents understandable it may also be necessary to migrate²²⁾ them in new formats, or to update the AIP with additional metadata related to emulation. Migration and other preservation related tasks may be carried out by the producer, OAIS archive and/or third parties. The party or parties responsible should be specified in the submission agreement.

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

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

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

The content and structure of all information packages in repository systems MUST be standardized. There are several packaging standards available, but the most commonly used one is the Metadata Encoding and Transmission Standard (METS²³⁾) developed by the Library of Congress. The use of METS as the container standard is recommended, although this document does allow the use of other container standards as well.

Since container standards – including METS - are rich specifications there is a need to create profiles to specify how they should be used. ISO/IEC TS 22424-2 provides a METS profile for EPUB. Other container standards are not taken into account; if METS is replaced by another container specification, profiling needs to be done separately.

Some digital preservation projects have developed tools for creating SIPs that meet the project requirements, which makes it a lot easier to submit information to the repository system. Producers should nevertheless have at least basic understanding of digital preservation, since pre-ingest steps from document creation to SIP submission should not risk the authenticity of the documents to be submitted to the OAIS archive. A simple model of such steps is shown in [Figure 2](#).

22) 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".

23) <http://www.loc.gov/standards/mets/>

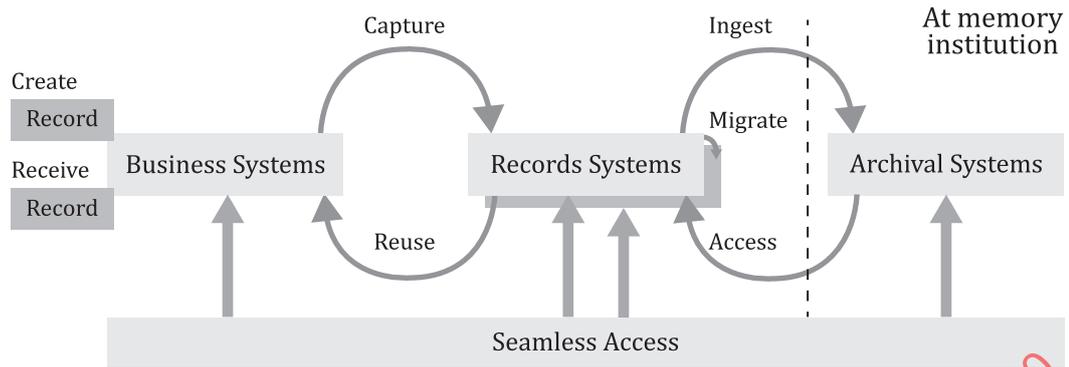


Figure 2 — Information flow between production system and OAIS archive (Reference [17], p.13)

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 BagIt²⁴⁾ instead of METS for storage of digital 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 preferences²⁵⁾. 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

6.1 Overview

According to the Open Archival Information System (OAIS) specification²⁶⁾, 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),

It 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 with the associated Content Information being disseminated. These variants are referred to as the Submission Information Package (SIP), the Archival Information Package (AIP), and the Dissemination Information Package (DIP). Although these are all

24) <https://tools.ietf.org/html/rfc8493>

25) https://www.cen.eu/news/calls/Calls/CEN-Call_for_tender_Digitalcinema.pdf

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

Information Packages, they differ in mandatory content and the multiplicity of the associations among contained classes.

The principles listed below provide SIP production guidelines for document producers (publishers or third parties creating EPUB publications). The creation of the principles has been inspired by the draft common requirements published by the E-ARK project (see Introduction to the Common Specification for Information Packages in the E-ARK project, version 2.0²⁷⁾). Although E-ARK has served as a model for this document, these requirements have not been aligned with those of E-ARK, and therefore there may be significant differences between the specifications.

6.2 General principles

6.2.1 EPUB publications shall be sent to a repository system as well-formed and complete submission information packages (SIPs)

- This document does not assume that publishers create SIPs. The OAI producer may be a third party acting on behalf of the publisher, such as hosting platform or other production vendor or even the OAI archive itself.
- This document and its accompanying document are mainly concerned with the structure and content of SIPs. The way EPUB publications are archived and disseminated (the structure of archival information packages and dissemination information packages, or AIPs and DIPs) depends on the submission agreements made between the archive and the producers, and on the operational principles of the archive, and is beyond the scope of this document. It is possible that an EPUB publication is migrated into another format during ingest, and disseminated again as an EPUB publication. If so, the OAI archive may also preserve (in bit level) the original file.
- Submitted EPUB publications shall be conformant with EPUB requirements²⁸⁾ and conformance should be validated.
- Submitted EPUB publications shall either contain or at least facilitate access to all the data and metadata required to render the content information successfully.
 - i. Preview publications may be submitted, even though they are by definition not complete, if the final documents are sent when ready. Depending on the submission agreement, the archive may preserve just the final version, or both versions of the resource. Identifiers shall be used in such a way that the OAI archive will be able to link all versions of the publication and delete preview versions, if that is the agreed preservation policy.
 - ii. Distributable objects shall not be submitted individually. They may be embedded within EPUB publications, but the archive is not obliged to deliver them as DIPs unless the submission agreement mandates that.
 - iii. Fonts shall be embedded into the EPUB publication in full and un-obfuscated, if font license allows that. If submission agreements allow submission of EPUB publications with obfuscated or non-embedded fonts, there is a risk that such publications become unusable in the future.
 - iv. Related resources such as audio and video should be embedded in the EPUB publication.
 - v. Remotely-hosted resources should be avoided, but if used, it is necessary to ensure that all remote data is available to the OAI archive so that the data can be incorporated into the AIP during ingest, and permission to do this shall be explicitly agreed upon in the submission agreement, especially if the publisher is not in full control of remote data.

27) <https://www.dilcis.eu/specifications/common-specification>

28) Conformance requirements for EPUB publications and reading systems have been specified in chapter 3.1 of EPUB Publications, version 3.0.1.

- vi. Descriptive and other metadata should be embedded in the SIP. METS mdRef element may only be used if a) referred metadata is part of the same SIP, or b) the archive is able to retrieve any linked external metadata and incorporate it into the AIPs in an appropriate format.
 - vii. Permission to use remote resources and metadata shall be specified in the submission agreement. 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. The permission shall specify acceptable metadata and file formats.
- The SIP should be checked for viruses and malicious software before submission to the OAIS archive. Some producers may not be able to make virus checks, but all OAIS archives shall be. Virus checks are commonly done during ingest.
 - EPUB publications in SIPs should not be encrypted, because that compromises long-term preservation. If data is submitted in an encrypted format, the archive shall receive necessary decryption information/details within the SIP, as agreed in the submission agreement or elsewhere. When the archive disseminates the archived data to its customers, it can be encrypted again.
 - DRM protection, if any, should be removed by the producer before the document is submitted. If the content in the SIP is DRM protected, the archive shall receive the necessary information/details to remove the DRM protection within the SIP, as agreed in the submission agreement or elsewhere. Such permission may be producer-specific, based on the submission agreement, or a generic permission, based on e.g. the Copyright Act.
 - If data is compressed, the user of the compression method shall be specified using the Compression metadata element in the EPUB's encryption.xml file.
 - The submission agreement should specify at least one EPUB reading system capable of rendering the submitted EPUB publications successfully. Knowing the reading system requirements in advance makes it easier for the archive to design and implement the ingest process. Although submitted publications will usually be validated only with automated tools²⁹⁾, the archive should be able to validate that the received EPUB can be presented to the customers, and check for instance the look and feel of archived EPUB publications before and after migration. This is possible only if the archive can operate the reading systems that can render the archived publications successfully.
 - Each SIP should specify EPUB reading system or systems, which can render the EPUB publication in the SIP. If this information is missing, reading system or systems shall be specified in the submission agreement.
 - Multiple rendition EPUB publications may be designed for multiple reading systems, in which case the submission agreement may require the archive to carry out at least occasional checks in all of these reading systems. If so, all these reading systems should be listed in the submission agreement.
 - If a submitted EPUB publication has been optimized for a certain reading system, the system should be described in the document's technical and/or preservation metadata, since such information is valuable for preservation and archival access purposes.
 - If the optimal EPUB reading system is no longer available, the archive should, with permission and support from the producer, either find another suitable reading system or modify the ingest process so that the EPUB publications affected by this change can be used by another EPUB reading system. While this document 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.

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

6.2.2 Regardless of its type or format, it shall be possible to include any data or metadata in SIPs

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

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

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

E-books are likely to contain more interactive features in the future. From preservation point of view it is therefore a problem that there are various ways in which EPUB 3 can support interactivity. On the other hand, some EPUB reading systems do not support interactivity at all, and even if it is supported, it is possible that different reading systems do not behave identically, partly because EPUB is not specific about how support should be implemented. EPUB 3 `object` element enables the use of arbitrary embedded executables that are not inherently supported in EPUB 3 reading systems. A common use case would be to include proprietary applets or Adobe Flash applications. However, in a majority of cases, interactive publications will be created through the use of in-book source code. Because JavaScript is the de facto standard scripting language for SVG and HTML5, EPUB 3 content documents can be assumed to be scriptable only if they contain JavaScript code. The standard does not define which versions of JavaScript (ECMAScript) are required to get the support. Content creators should comply with the most commonly supported features in web browsers for best results^[4]. Usage of common tools and techniques will also make it easier to preserve the publication in the long-term, either via emulation (a common solution for software preservation) or migration.

- Archives offering long-term preservation services for EPUB publications should keep track of EPUB core media types and consider the possibility of including them on the list of archivable formats. If this is not viable, the archives should maintain clearly defined and well tested migration pathways from non-archivable core media types into archivable formats. Then the archive would not need to migrate these images during ingest and it would be possible to preserve EPUB publications unchanged³⁰⁾.
- If there is a foreign resource embedded or linked to a submitted EPUB publication, a fallback chain ending in a core media type resource should be provided even if the foreign resource is in an archivable

30) An OAIS archive does not need to migrate non-archivable file formats during the ingest process. Depending on the preservation strategy, migration can 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.

format. (Note that this requirement is stricter than those in EPUB 3 and 3.0.1 specifications, which require a fallback only in certain situations.)

- The producer may include foreign resources (and metadata formats) in submitted EPUB publications if they have been specified as suitable for ingest and/or archivable in the submission agreement, or if their METS encoding in SIPs makes it possible to ignore them during ingest (see below).
- If foreign resources and metadata are originally in un-archivable formats (formats that have not been specified as acceptable in the submission agreement), they shall be migrated during pre-ingest. The SIP may contain either just migrated publications, or both the original and migrated publications. Note that the preservation method may be either emulation or migration, so this requirement does not mean migration-only approach.
- Core media types and foreign resources not specified in the submission agreement may be submitted if and only if the submission agreement allows it. METS encoding of these files in SIPs shall make it possible to skip their validation against the generic ingest criteria during the ingest process (since otherwise the SIP shall not pass the validation) and therefore passed directly to AIP. The specifics of this type of encoding shall be defined in the submission agreement.
- If there are alternative versions (renderings) of the publication to be included in the SIP which are not archivable, they should be migrated into acceptable file formats prior to submission by the producer or a third-party preparing a SIP on behalf of the producer. For instance, if PDF is specified as not archivable but PDF/A is, the producer should create a PDF/A version of the document, which will then be submitted to the repository system alongside the EPUB publication of the same work.
- If these non-archivable originals are submitted, their METS encoding in SIPs shall make it possible to skip validation against the generic ingest criteria during the ingest process (since otherwise the SIP would not pass) and therefore passed directly to AIP. The specifics of this type of encoding shall be defined in the submission agreement. 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.

NOTE 2 EPUB 3 Fixed Layout Properties.

In digital preservation the usual aim is to preserve intellectual content. Preserving also the original look and feel of the document is more challenging, although that may be required for some resources or collections. Reflowable EPUB publications are designed so that their look and feel can change with no impact on semantics, which is a good thing from the digital preservation point of view, 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 (ISO/IEC TS 30135-7). 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.

6.2.3 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.

6.2.4 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 or 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.
- ISO/IEC TS 22424-2 provides an example of how a digital signature can be used for verification.

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

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

6.2.6 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 document³¹⁾. 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.

6.2.7 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.

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

6.2.8 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.

6.2.9 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 document.

6.3 Identification of information packages and their content

6.3.1 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.
- The entire SIP or parts of it shall be resubmitted in a revised format if the ingest process fails due to errors in the package. To keep track of the packages, SIPs shall have unique identifiers.

6.3.2 Information objects (EPUB publications, PREMIS preservation metadata record, etc.) within SIPs shall be identified uniquely and persistently

- Identifiers have many vital uses in digital preservation. They are used as access keys to the archived content in repository systems and facilitate information exchange with external systems. Identifiers also enable linking different versions of an archived document to each other. Moreover, with identifiers it is possible to link documents and descriptive/administrative metadata records that describe them. These links enable the archive to e.g. create dissemination information packages with the requested content.
- Submission agreements shall specify identifier systems used, their location (EPUB document or SIP) and who is responsible of creating them (producer, archive or a third party). For instance, if the use of EPUB release identifiers is forbidden because the repository system does not support them, another means of identifying releases is needed.
- International standard identifiers, such as ISBNs for books and DOIs for articles, shall be used as EPUB unique identifiers whenever possible. Any exceptions (such as using other identifier systems for releases which do not have ISBNs) should be specified in the submission agreement.

- It should be possible to express the identifiers (also) as actionable HTTP URIs. Usage of persistent identifiers (Handles, DOIs, URNs, or ARKs) is recommended.
- If there are multiple renditions of a work in an EPUB publication, requirements in the EPUB Multiple-Rendition Publications 1.0 specification shall be followed. Each rendition of an EPUB publication in a SIP shall have its own identifier.
- The SIP should contain separate descriptive and administrative metadata records for each rendition, and these records shall have their own identifiers.

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

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

6.3.3 EPUB Fragment Identifiers should not be used in EPUB publications sent to a repository system, unless the submission agreement explicitly allows their use

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

6.4 Structure of information packages

Submission information packages shall be built so that their components can be logically and physically separated from one another

- For each rendition of the EPUB content document, there shall be a manifest file, which identifies and describes a set of resources that collectively compose a given rendition of a document, and EPUB spine, which provides a default reading order for a given rendition.
- EPUB Open Container Format (OCF) defines a file format and processing model for encapsulating a set of related resources (for instance, renditions of the same resource) into a single-file (ZIP) EPUB container³²⁾.
- The structure of each EPUB ZIP archive shall be described using the EPUB container.xml file (which describes the locations of root files of available renditions of the EPUB publication, and the rendition's package document and navigation document).
- EPUB package document and EPUB navigation document shall contain all metadata needed for rendering the publication, including the recommended reading system.

32) 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.