

# PDF/A-3u as an archival format for Accessible mathematics\*

Ross Moore

Macquarie University, Sydney, Australia  
ross.moore@mq.edu.au

**Abstract.** Including  $\text{\LaTeX}$  source of mathematical expressions, within the PDF document of a text-book or research paper, has definite benefits regarding ‘Accessibility’ considerations. Here we describe three ways in which this can be done, fully compatibly with international standards ISO 32000, ISO 19005-3, and the forthcoming ISO 32000-2 (PDF 2.0). Two methods use embedded files, also known as ‘attachments’, holding information in either  $\text{\LaTeX}$  or MathML formats, but use different PDF structures to relate these attachments to regions of the document window. One uses structure, so is applicable to a fully ‘Tagged PDF’ context, while the other uses `/AF` tagging of the relevant content. The third method requires no tagging at all, instead including the source coding as the `/ActualText` replacement of a so-called ‘fake space’. Information provided this way is extracted via simple `Select/Copy/Paste` actions, and is available to existing screen-reading software and assistive technologies.

## 1 Introduction

PDF/A is being adopted by publishers and Government agencies for the long-term preservation of important documents in electronic form. There are a few variants, which pay more or less regard to Accessibility considerations; i.e., ‘a’ for *accessible*, ‘b’ for *basic*, ‘u’ for (presence of) *unicode mappings* for all font characters. Later versions [3,4] of this ISO standard [2] allow for other file attachments in various data formats. In particular, the PDF/A-3u variant allows the inclusion of *embedded files of arbitrary types*, to convey supplementary descriptions of technical portions of a document’s contents.

‘Accessibility’ is more relevant for reports and text-books than for research outputs. In fact in some countries it is a legal requirement that when a visually-impaired student enrolls in unit of study for which a text-book is mandated as ‘Required’, then a fully accessible version of the contents of that book *must* be made available. Anecdotally, visually-impaired students of mathematics and related fields much prefer mathematical material to be made available as  $\text{\LaTeX}$  source, to any other format. With a Braille reader, this is text-based and sufficiently compact that expressions can be read and re-read with ease, until a full

---

\* The final publication is available at <http://link.springer.com/>, as part of:  
S.M. Watt et al. (Eds.): CICM 2014, LNAI 8543, pp. 184–199, 2014.

understanding has been achieved. This is often preferable to having an audio version [13,14], which is less-easy to navigate. Of course having *both* a well-structured audio version, as well as textual source, is even more useful. The PDF example [12] accompanying this paper<sup>1</sup> in fact has both, though here we concentrate on how the latter is achievable within PDF documents.

Again anecdotally, the cost of reverse-engineering<sup>2</sup> all the mathematical expressions within a complete textbook is typically of the order of £10,000 or AUD 30,000 or CAD 10,000. This cost would have been dramatically reduced if the PDF had originally been created to include a  $\text{\LaTeX}$  or MathML description of each expression<sup>3</sup>, attached or embedded for recovery by the PDF reader or other assistive technology. How to do this in PDF is the purpose of this paper.

The method of *Associated Files*, which is already part of the PDF/A-3 standard [4], is set to also become part of the ISO 32000-2 (PDF 2.0) standard [6], which should appear some time in 2014 or 2015. In Sect. 3 this mechanism is discussed in more detail, showing firstly how to include the relevant information as attachments, which can be extracted using tools in the PDF browser. The second aspect is to relate the attachments to the portion of content as seen on-screen, or within an extractable text-stream. This can be specified conveniently in two different ways. One way requires structure tagging to be present (i.e., a ‘Tagged PDF’ document), while the other uses direct tagging with an `/AF` key within the content stream. In either case a PDF reader needs to be aware of the significance of this `/AF` key and its associated embedded files.

With careful use of the `/ActualText` attribute of tagged content,  $\text{\LaTeX}$  (or other) source coding of mathematical expressions can be included within a PDF document, virtually invisibly, yet extractable using normal `Select/Copy/Paste` actions. A mechanism, using very small space characters inserted before and after each mathematical expression, is discussed in Sect. 4. This is applicable with *any* PDF file, not necessarily PDF/A. It is important that these spaces not interfere with the high-quality layout of the visual content in the document, so we refer to them as ‘fake spaces’.

The various Figures in this paper illustrate the ideas and provide a look at the source coding of a PDF document<sup>1</sup> that includes all the stated methods, thus including the  $\text{\LaTeX}$  source of each piece of mathematical content. (Where explicit PDF coding is shown, the whitespace may have been massaged to conserve space within the pages of this paper.) Indeed the example document includes as many as 7 different representations of each piece of mathematical content:

- the visual form, as typically found in a PDF document;
- the  $\text{\LaTeX}$  source, in two different ways; i.e. an attachment associated with a `/Formula` structure tag and also associated directly to the (visual) content, and as the `/ActualText` replacement of a ‘fake space’.

<sup>1</sup> ... should be attached prior to the ‘References’, else downloadable online; see [12].

<sup>2</sup> ... with prior permission granted by the publisher ...

<sup>3</sup> This is distinct from including the complete  $\text{\LaTeX}$  source of the whole document. There are many reasons why an author, and hence the publisher, might not wish to share his/her manuscript; perhaps due to extra information commented-out throughout the source, not intended for general consumption.

- a MathML version as an attachment, also associated to the `/Formula structure` tag and also associated directly to the (visual) content;
- a MathML representation through the structure tagging;
- words for a phonetic audio rendering, to be spoken by ‘Read Out Loud’;
- the original  $\text{\LaTeX}$  source of the complete document, as a file attachment associated with the document as a whole.

In practice not all these views need be included to satisfy ‘Accessibility’ or other requirements. But with such an array of representations, it is up to the PDF reading software to choose those which it wants to support, or which to extract according to particular requirements of end-users. It is remarkable that a single document can be so enriched, yet still be conforming with a standard such as PDF/A-3u, see [4]. Indeed, with all content being fully tagged, this document<sup>1</sup> would also validate for the stricter PDF/A-3a standard, apart from the lack of a way to specify the proper rôle of MathML structure tagging, so that tags and their attributes are preserved under the ‘Save As Other ... XML 1.0’ export method when using Adobe’s ‘Acrobat Pro’ software. This deficiency will be addressed in PDF 2.0 [6].

Methods used to achieve the structure tagging in the example document<sup>1</sup> have been the subject of previous talks and papers [10,9] by the author. It is not the intention here to promote those methods, but rather to present the possibilities for mathematical publishing and ‘Accessibility’ that have been opened up by the PDF/A-3 and PDF/UA standards [4,7], and the ‘fake spaces’ idea. The example document [12] is then just a ‘proof-of-concept’ to illustrate these possibilities.

Since the PDF/A-3 standard [4] is so recent, and with PDF 2.0 [6] yet to emerge, software is not yet available that best implements the ‘Associated Files’ concept. The technical content of the Figures is thus intended to assist PDF software developers in building better tools in support of accessible mathematics. It details (i) exactly what kind of information needs to be included; (ii) the kind of structures that need to be employed; and (iii) how the information and structures relate to each other. For those less familiar with PDF coding, the source snippets have been annotated with high-lighting<sup>4</sup> and extra words indicating the ideas and intentions captured within each PDF object. Lines are used to show relationships between objects within the same Figure, or ‘see Fig. Xx’ is used where the relationship extends to parts of coding shown within a different Figure. Section 2 is supplied to give an overview of the PDF file structure and language features so that the full details in the Figures can be better understood and their rôle appreciated.

## 2 Overview of the PDF file format

PDF files normally come employing a certain amount of compression, to reduce file-size, so appear to be totally intractable to reading by a human. Software

<sup>4</sup> ... with consistent use of colours, in the PDF version of this paper ...

techniques exist to undo the compression, or the PDF file may have been created without using any. The example document<sup>1</sup> was created without compression, so can be opened for reading in most editing software.

The overall structure of an uncompressed PDF file consists of:

- (a) a collection of numbered *objects*: written as `<num> 0 obj ... endobj` where the ‘...’ can represent many, many lines of textual (or binary) data starting on a new line after `obj` and with `endobj` on a line by itself. The numbering need not be sequential and objects may appear in any order. An *indirect reference* sequence of the form `<num> 0 R`<sup>5</sup> is used where data from one object is required when processing another. A *cross-reference table* (described next), allows an object and its data to be located precisely. Such indirect references are evident throughout the coding portions of Figures 1–5.
- (b) the *cross-reference table*: listing of byte-offsets to where each numbered object occurs within the uncompressed PDF file, together with a linked listing of unused object numbers. (Unused numbers are available for use by PDF editing software.)
- (c) the *trailer*, including: (i) total number of objects used; (ii) reference to the document’s `/Catalog`, see Fig. 3c; (iii) reference to the `/Info` dictionary, containing file properties (i.e., basic metadata); (iv) byte-offset to the cross-reference table; (v) encryption and decryption keys for handling compression; (vi) end-of-file marker.

Thus the data in a PDF file is contained within the collection of objects, using the cross-reference table to precisely locate those objects. A PDF browser uses the `/Catalog` object (e.g., object 2081 in Fig. 3c) to find the list of `/Page` objects (e.g., object 5 in Fig. 3b), each of which references a `/Contents` object. This provides each page’s *contents stream* of graphics commands, which give the details of how to build the visual view of the content to be displayed. A small portion of the page stream for a particular page is shown in Figures 1b, 3a, 5a.

*Character strings* are used in PDF files in various ways; most commonly for ASCII strings, in the form (...); see Figures 1a, 1b, 2b, 2c, 3a, 3c, and 5a. Alternatively, a hexadecimal representation with byte-order mark `<FEFF...>` can be used, as in Figures 1b, 3a, 5a. This is required particularly for Unicode characters above position 255, with ‘surrogate pairs’ used for characters outside the basic plane, as with the *k* variable name in those figures. Below 255 there is also the possibility of using 3-byte octal codes within the (...) string format; see footnote 11 in Sect. 4. For full details, see §7.3.4 of PDF Specifications [1,5].

*PDF names* of the form `/<name>`, usually using ordinary letters, have a variety of uses, including (i) *tag-names* in the content stream (Figures 1b, 3a, 5a); (ii) identifiers for *named resources* (Fig. 3b within object 20 and in the `/AF` tagging shown in Fig. 3a); and extensively as (iii) dictionary *keys* (in all the Figures 1, 2, 3, 5) and frequently as dictionary *values* (see below).

Other common structures used within PDF objects are as follows.

<sup>5</sup> The ‘0’ is actually a *revision number*. In a newly constructed PDF this will always be 0; but with PDF editing software, higher numbers can result from edits.

- (i) **arrays**, represented as [*⟨item⟩ ⟨item⟩ ... ⟨item⟩*], usually with similar kinds of *⟨item⟩*, (see e.g., Figures 1a, 3b, 3c) or alternating kinds (e.g., the filenames array of Fig. 2b).
- (ii) **dictionaries** of *key–value pairs*, similar to alternating arrays, but represented as <<*⟨key<sub>1</sub>⟩ ⟨value<sub>1</sub>⟩ ⟨key<sub>2</sub>⟩ ⟨value<sub>2</sub>⟩ ... >>*. The *⟨key⟩* is always a *PDFname* whereas the *⟨value⟩* may be any other element (e.g., string, number, name, array, dictionary, indirect reference). The *key–value pairs* may occur in any order, with the proviso that if the same *⟨key⟩* occurs more than once, it is the first instance whose *⟨value⟩* is used. A */Type* key, having a *PDFname* as value, is not always mandatory; but when given, one refers to the dictionary object as being of the type of this name. See Figures 1a, 2b, 2c, 3b, 3c and 5b for examples.
- (iii) **stream objects** consist of a dictionary followed by an arbitrarily-long delimited *stream* of data, having the form << ... >> **stream** ... **endstream**, with the **stream** and **endstream** keywords each being on a separate line by themselves (see objects 26 and 28 in Fig. 2c). The dictionary must include a */Length* key, whose value is the integer number of bytes within the data-stream. With the length of the data known, between the keywords on separate lines, there is no need for any escaping or special encoding of any characters, as is frequently needed in other circumstances and file-formats. See §7.3.8 of [1,5] for more details; e.g., how compression can be used.
- (iv) **graphics operators** which place font characters into the visual view occur inside a page contents stream, within portions delimited by **BT** ... **ET** (abbreviations for **B**egin/**E**nd **T**ext); see Figures 1b, 3a, 5a. These include coding */⟨fontname⟩ ⟨size⟩ Tf* for selecting the (subsetted) font, scaled to a particular size, and [*⟨string⟩*]TJ for setting the characters of the string with the previously selected font. See §9.4 of [1,5] for a complete description of the available text-showing and text-positioning operators.

Dictionaries and arrays can be nested; that is, the *⟨value⟩* of a dictionary item's *⟨key⟩* may well be another dictionary or array, as seen in objects 20 and 90 within Fig. 3b. Similarly one or more *⟨item⟩*s in an array could well be a dictionary, another array, or an indirect reference (regarded as a 'pointer' to another object).

With the use of *PDF names*, *objects*, and *indirect references* a PDF file is like a self-contained web of interlinked information, with names chosen to indicate the kind of information referenced or how that information should be used.

The use of objects, dictionaries (with key–value pairs) and indirect references makes for a very versatile container-like file format. If PDF reader software does not recognise a particular key occurring within a particular type of dictionary, then both the key and its value are ignored. When that value is an indirect reference to another object, such as a *stream object*, then the data of that stream may never be processed, so does not contribute to the view being built. Thus PDF producing or editing software may add whatever objects it likes, for its own purposes, without affecting the views that other PDF reading software wish to construct. This should be contrasted with HTML and XML when a browser

does not recognise a custom tag. There that tag is ignored, together with its attributes, but any content of that tag *must still be handled*.

It is this feature of the PDF language which allows different reader software to support different features, and need not use all of the information contained within a PDF file. For example, some browsers support attachments; others do not. A PDF format specification now consists mostly of saying which tags and dictionary keys *must* be present, what others are allowed, and how the information attached to these keys and tags is intended to be used. Hence the proliferation of different standards: PDF/A, PDF/E, PDF/VT, PDF/UA, PDF/X, perhaps with several versions or revisions, intended for conveying different kinds of specialised information most relevant within specific contexts.

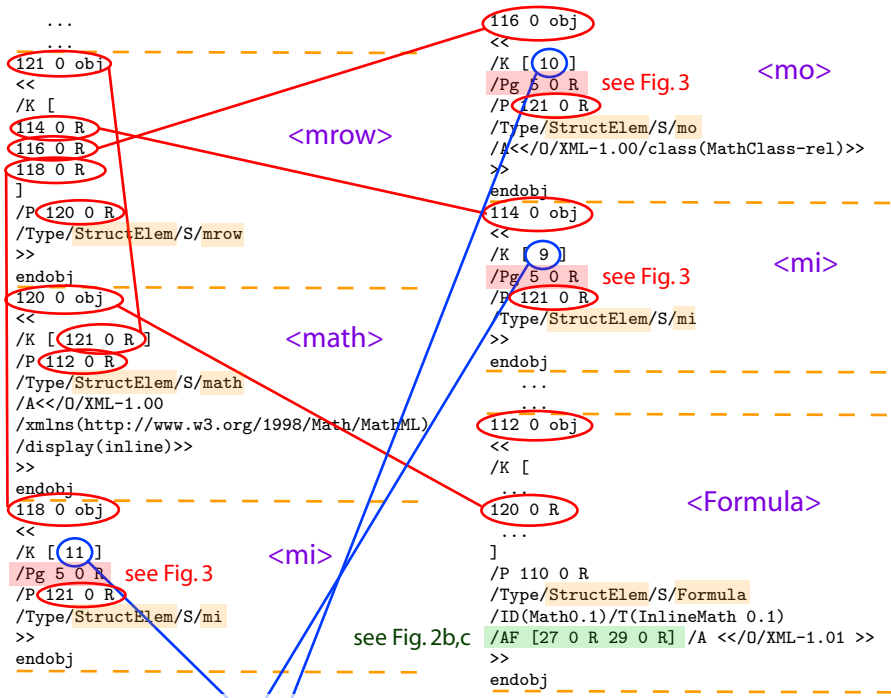
## 2.1 Tagging within PDF documents

Two types of tagging can be employed within PDF files. ‘Tagged PDF’ documents use both, with content tags connected as leaf-nodes of the structure tree.

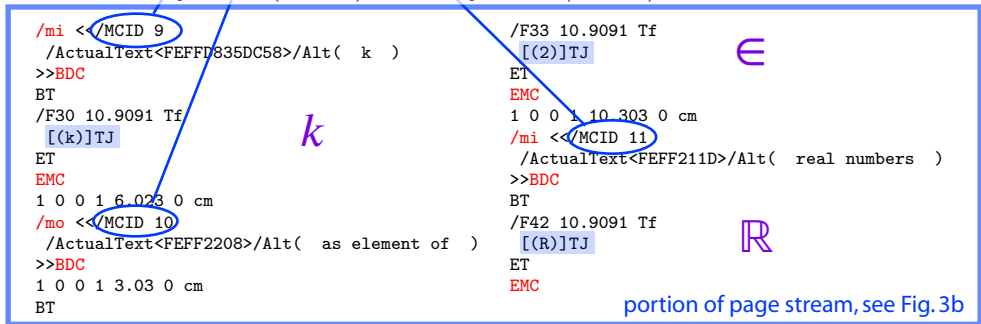
*Tagging of content* is done as  $\langle tag \rangle \langle dict \rangle$  BDC ... EMC within a contents stream. Here the BDC and EMC stand for ‘Begin Dictionary Content’ and ‘End Marked Content’ respectively, with the  $\langle dict \rangle$  providing key-value pairs that specify ‘properties’ of the *marked content*, much like ‘attributes’ in XML or HTML tagging<sup>6</sup>. The  $\langle tag \rangle$  can in principle be any *PDF name*; however, in §14.6.1 of the specifications [1,5] it stipulates that “All such tags shall be registered with Adobe Systems (see Annex E) to avoid conflicts between different applications marking the same content stream.” Thus one normally uses a standard tag, such as /Span, or in the presence of structure tagging (see below) choose the same tag name as for the parent structure node. Figures 1b, 3a, 5a show the use of Presentation-MathML content tag names, which are expected to be supported in PDF 2.0 [6]. Typical attributes are the /ActualText and /Alt strings, which allow replacement text to be used when content is extracted from the document using Copy/Paste or as ‘Accessible Text’ respectively. The /MCID attribute allows *marked content* to be linked to document structure, as discussed below. A variant of this tagging uses a *named resource* for the  $\langle dict \rangle$  element. This is illustrated with /AF content tagging in Sect. 3.2.

*Tagging of structure* requires building a tree-like structural description of a document’s contents, in terms of Parts, Sections, Sub-sections, Paragraphs, etc. and specialised structures such as Figures, Tables, Lists, List-items, and more [1,5, §14.8.4]. Each *structure node* is a dictionary of type /StructElem having keys /S for the structure type, /K an array of links to any child nodes (or Kids) including *marked content* items, and /P an indirect reference to the parent node. Optionally there can be a /Pg key specifying an indirect reference to a /Page dictionary, when this cannot be deduced from the parent or higher ancestor. Also, the /A key can be used to specify attributes for the structure tag when the

<sup>6</sup> Henceforth we use the term ‘attribute’, rather than ‘property’.



(a) PDF coding of the /Formula structure node showing the reference to ‘Associated Files’ via the /AF key in object 112. The indirect references (27 and 29) correspond to /Filespec dictionaries, as shown in Fig. 2. (In the coding ‘...’ indicates parts omitted due to not being relevant to this structure; these portions are discussed in Sect. 4.) The corresponding *marked content* is specified via the /K [ ... ] numbers (9, 10, 11) in the child structure nodes; i.e., objects 114 (<mi>), 116 (<mo>) and 118 (<mi>), which are children of object 121 (<mrow>) under object 120 (<math>).



(b) Portion of the PDF page content stream showing the /MCID numbers (9, 10, 11) of the actual content portions of the mathematical expression. These correspond to leaf-nodes of the structure tree as presented in part (a).

**Fig. 1.** PDF coding for portions of (a) the structure tree and (b) the page content stream, corresponding to the mathematics shown as selected in Fig. 1a. It is through the /MCID numbers that the association is made to the /Formula structure tag for the corresponding piece of mathematical content.



document’s contents are exported in various formats; e.g., using ‘Save As Other ... XML 1.0’ export from Adobe’s ‘Acrobat Pro’ browser/editor. Fig. 1a shows the MathML tagging of some inline mathematical content. The tree structure is indicated with lines connecting nodes to their kids; reverse links to parents are not drawn, as this would unduly clutter the diagram. Other keys, such as `/ID` and `/T` can provide an identifier and title, for use primarily in editing software to locate specific nodes within appropriately ordered listings.

The link between structure and *marked content* (as leaf-nodes to the structure tree, say) is established using the `/MCID` number attribute. A numeric integer entry in the `/K Kids` array corresponds to an `/MCID` number occurring within the contents stream for that page specified via a `/Pg` entry, either of the structure node itself or the closest of its ancestors having such a key. Fig. 1b shows this linking via `/MCID` with lines drawn to the corresponding structure nodes shown in Fig. 1a. The interplay of structure with content was addressed in the author’s paper [10], with Figure 1 of that paper giving a schematic view of the required PDF structural objects.

### 3 ‘Associated Files’, carrying $\LaTeX$ and MathML views of mathematical content.

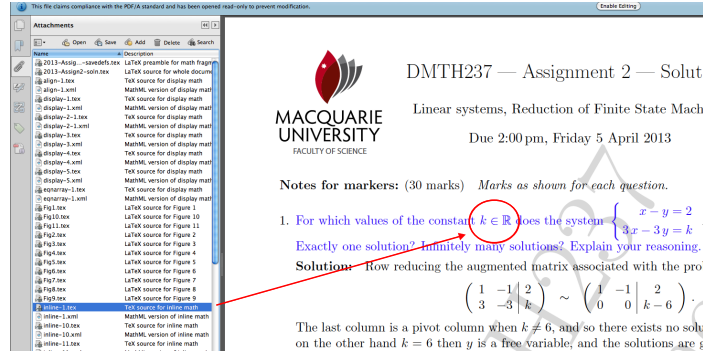
There are several ways in which file attachments may be associated with specific portions of a PDF document, using the ‘Associated Files’ technique [4, Annex E]. The file is embedded/attached and then associated, by a method, either to:

- (i) the document as a whole [4, §E.3], [6, §14.13.2] — e.g. the full  $\LaTeX$  source, or preamble file used when converting snippets of mathematical content into a MathML presentation of the same content;
- (ii) a specific page within the document [4, §E.4], [6, §14.13.3] or to a (perhaps larger) logical document part using PDF 2.0 [6, §14.13.7];
- (iii) graphic objects in a content stream [4, §E.5], [6, §14.13.4] — when structure is available, this is not the preferred method<sup>7</sup>;
- (iv) a structure node [4, §E.7], [6, §14.13.5] such as `/Figure`, `/Formula`, `/Div`, etc.
- (v) an `/XObject` [4, §E.6], [6, §14.13.6] such as an included image of a formula or other mathematical/technical/diagrammatic content;
- (vi) an annotation [4, §E.8] — but this method can be problematic with regard to validation for PDF/A [4, §6.3], and PDF/UA [7, §7.18] standards<sup>8</sup>.

<sup>7</sup> In the PDF/A-3 specifications [4, §E.5] the final paragraph explicitly states “When writing a PDF, the use of structure (and thus associating the `/AF` with the structure element, see [4, §E.7]) is preferred instead of the use of explicit marked content.” with a corresponding statement also in [6, §14.13.4].

<sup>8</sup> The method of indicating an attachment with a ‘thumb tack’ annotation located at a specific point within a document, is deprecated in the PDF/A-3 standard, as it does not provide a proper method to associate with the portion of content. Besides, the appearance of such thumb-tacks all over paragraphs containing inline mathematics is, well, ... .. downright ugly.





(a) Listing of attachments, indicating how one is associated to some inline content.

```

1860 0 obj   see Fig.3c
<</Names [ (2013-Assign2-soln.tex) 22 0 R (2013-Assign2-soln-savedefs.tex) 24 0 R (inline-1.tex)
27 0 R (inline-1.xml) 29 0 R (inline-2.tex) 31 0 R (inline-2.xml) 33 0 R (display-1.tex)
... (Fig11.tex) 1858 0 R ]>> array of filenames of embedded files
endobj

112 0 obj
<<
/K [
113 0 R
120 0 R
122 0 R
]
/P 110 0 R
/Type/StructElem/S/Formu      <Formula>
/ID(Math0.1)/T(InlineMath 0.1)
/AF [27 0 R 29 0 R] /A <</XML-1.01 >>
>>
endobj

```

(b) A portion of the PDF coding of the /Names array (upper) for embedded files, and coding of the /Formula structure element (lower), showing the /AF key with array of two 'Associated' files given as indirect references. The /Names array is used to produce the listing in (a) and provides indirect links to /Filespec entries, as shown in (c).

```

26 0 obj
<<
/Type/EmbeddedFile
/Subtype/application#2Ffx-tex
/Params <<
/ModDate (D:20140201111224+11'00')
/Size (16) /Checksum
(1EB423943E9CBB33E861CC037C8F06A)>>
/Length 16
>>
stream
\(\ k \in \mathbb{R} \)
endstream
endobj

27 0 obj
<</Type/Filespec
/F (inline-1.tex)/UF (inline-1.tex)
/Desc (TeX source for inline math)
/AFRelationship /Source
/EF<</F 26 0 R>> >>
endobj

28 0 obj
<<
/Type/EmbeddedFile
/Subtype/application#2Fmathml+xml
/Params <<
/ModDate (D:20140131152820+11'00')
/Size (164) /Checksum
(09D0F1EFDEBC32F9DB4A9F8BA05DC0E5)>>
/Length 164
>>
stream
<math
xmlns="http://www.w3.org/1998/Math/MathML"
display="inline" ><mrow
><mi
>k</mi> <mo
class="MathClass-rel">&#x2208;</mo> <mi
>&#x211D;</mi></mrow></math>
endstream
endobj

29 0 obj
<</Type/Filespec
/F (inline-1.xml)/UF (inline-1.xml)
/Desc (MathML version of inline math)
/AFRelationship /Supplement
/EF<</F 28 0 R>> >>
endobj

```

(c) PDF source of /Filespec dictionaries and /EmbeddedFile objects which hold the streams of L<sup>A</sup>T<sub>E</sub>X and MathML coding for the mathematical content indicated in (a).

Fig. 2. Embedded files associated with a /Formula structure element.

Fig. 2a shows how attachments are presented within a separate panel of a browser window, using information from an array of filenames; see Fig. 2b. This is independent of the page being displayed, so the array must be referenced from the document level. This is seen in Fig. 3c using the `/Names` key of the `/Catalog` dictionary, which references object 2080, whose `/EmbeddedFiles` key then references the filenames array (object 1860 in Fig. 2b). One can also see in Fig. 2b how each filename precedes an indirect reference to the `/Filespec` dictionary [6, §7.11.3, Tables 44 and 45] for the named file; see Fig. 2c. This dictionary contains a short description (`/Desc`) of the type of content as well as the filename to use on disk, and a link via the `/EF` key to the actual `EmbeddedFile` stream object.

That a file is ‘Associated’ is indicated by the `/AFRelationship` key, whose value is a *PDF name* indicating how the file is related to visible content. Options here are `/Source` as used with the  $\text{\LaTeX}$  source coding, or `/Supplement` as used with the MathML description. Other possibilities are `/Data` (e.g., for tabular data) and `/Alternative` for other representations such as audio, a movie, projection slides or anything else that may provide an alternative representation of the same content. `/Unspecified` is also available as a non-specific catch-all.

Not all attachments need be ‘Associated’ and conversely not all ‘Associated Files’ need be displayed in the ‘Attachments’ panel, so there is another array (object 1859) as shown in Fig. 3c, linked to the `/MarkInfo` sub-dictionary of the `/Catalog` dictionary. Files associated with the document as a whole, as in method (i) above, link via the `/AF` key in the `/Catalog` dictionary (see Fig. 3c).

For the  $\text{\LaTeX}$  source of a mathematical expression method (iv) is preferred, provided structure tagging is present within the PDF. This is discussed below in Sect. 3.1. Method (iii) also works, provided the expression is built from content confined to a single page. This is described in Sect. 3.2.

As ‘Associated Files’ have only been part of published PDF/A standards [4] since late 2012, it may be some time before PDF readers provide a good interface for ‘Associated Files’, beyond using the ‘Attachments’ pane. This ought to include interfaces to view the contents of attached files, do searching within the files, and make the file’s contents available to assistive technology. One possible way to display this association is apparent in earlier work [11], whereby a bounding rectangle appears as the mouse enters the appropriate region.

### 3.1 Embedded files associated with structure

With an understanding of how structure tagging works, as in Sect. 2.1, then associating files to structure is simply a matter of including an `/AF` key in the structure node’s dictionary, as shown in Figures 1a and 2b. The value for this key is an array of indirect references to `/Filespec` objects for the relevant files.

There is nothing in the content stream in Fig. 1b to indicate that there is a file associated with this structure node. Rather the browser, knowing that ‘Associated’ files are present, needs to have gone through some pre-processing to first locate the node (if any) to which it is associated, then trace down the structure tree to the deepest child nodes (objects 114, 116, 118). From their `/K`

please consult the cited reference.

```

/AF /inline-1 BDC
1 0 0 1 158.485 0 cm
...
/mi <</MCID 9 see Fig. 1b
/ActualText<FEFFD835DC58>/Alt( k )
>>BDC
EMC
...
1 0 0 1 10.303 0 cm

/mi <</MCID 11 see Fig. 1b
/ActualText<FEFF211D>/Alt( real numbers )
>>BDC
BT
/F42 10.9091 Tf
[(R)]TJ
ET
EMC
...
EMC
    
```

(a) A portion of the PDF content stream associating content with embedded files for the mathematical expression indicated in Fig. 2a, using a *marked content* tag /AF to refer to a *named resource* /inline-1. All content down to the final EMC is associated.

```

20 0 obj
<< /inline-1 [27 0 R 29 0 R] /inline-2 [31 0 R 33 0 R] /display-1 [35 0 R 37 0 R]
/inline-3 [39 0 R 41 0 R] /inline-4 [43 0 R 45 0 R] /inline-5 [47 0 R 49 0 R]
... /inline-31 [1475 0 R 1477 0 R] >>
-----
90 0 obj
<<
/Properties 20 0 R
/Font << /F75 97 0 R /F79 100 0 R
/F77 102 0 R /F45 105 0 R /F78 106 0 R
... >>
/XObject << /Im1 25 0 R >>
/ProcSet [ /PDF /Text ]
>>
endobj
-----
5 0 obj
<<
/Type /Page
/Contents 91 0 R
/Resources 90 0 R
/MediaBox [0 0 595.276 841.89]
/Tabs/S
/Parent 773 0 R
/StructParents 0
>>
endobj
    
```

(b) A portion of the /Properties dictionary (upper, object 20) which is linked to a /Page object (lower right, object 5) via its /Resources key (see lower left, object 90). Thus a name (such as /inline-1) is associated with an array of /Filespec references (viz. [27 0 R 29 0 R]), which lead to the L<sup>A</sup>T<sub>E</sub>X and MathML files seen in Fig. 2c.

```

2080 0 obj
<<
/Dests 2079 0 R
/EmbeddedFiles 1860 0 R
>>
endobj
2081 0 obj
<<
/Type /Catalog
/Pages 773 0 R
/Names 2080 0 R
/ViewerPreferences << /DisplayDocTitle true >>
/OutputIntents [ << /Type /OutputIntent
/S/GTS_PDFa1 /DestOutputProfile 1 0 R
/OutputConditionIdentifier
(sRGB_IEC61966-2-1_no_black_scaling) /Info
(sRGB IEC61966 v2.1 without black scaling) >> ]
/Metadata 2 0 R/Lang (en-US)
/PageMode/UseOutlines
/MarkInfo << /Marked true /AF 1859 0 R >>
/AF [ 22 0 R 24 0 R ]
/PageLabels << /Nums [0 << /P(1) >> ... ] >>
/OpenAction 4 0 R
/StructTreeRoot 95 0 R
>>
endobj
-----
1859 0 obj
[ 27 0 R 29 0 R 31 0 R 33 0 R 35 0 R 37 0 R 39 0 R 41 0 R 43 0 R 45 0 R 47 0 R 49 0 R 51 0 R
53 0 R 55 0 R 57 0 R 59 0 R 61 0 R 63 0 R 65 0 R 67 0 R 69 0 R 71 0 R 73 0 R 75 0 R 77 0 R
... 1850 0 R ]
    
```

(c) The document's /Catalog (object 2081) indicates presence of embedded files via the /Names key (object 2080). This references the array (object 1860 in Fig. 2b), to establish the correspondence between filenames and /Filespec dictionaries. Embedded files which are 'Associated' to content portions are also listed in an array (object 1859) referenced from the /AF key in the /MarkInfo dictionary.

**Fig. 3.** Embedded files associated with specific content.

entries (viz., 9, 10, 11 resp.), the relevant *marked content* in the page’s contents stream is located using these /MCID numbers.

### 3.2 Embedded files associated with content

With an understanding of how content tagging works, as in Sect. 2.1, and the fact that *marked content* operators may be nested, then associating files to content is also quite simple. One simply uses an /AF tag within the page’s content stream with BDC ... EMC surrounding the content to be marked, as shown in Fig. 3a. This employs the *named resource* variant (here /inline-1) to indicate the array of ‘Associated’ files. Fig. 3b shows how this name is used as a key (in dictionary object 20) having as value an array of indirect references to /Filespec objects (27 and 29). These resources can be specific to a particular page dictionary (object 5), but in the example document<sup>1</sup> the *named resources* are actually made available to all pages, since this accords with not including multiple copies of files when a mathematical expression is used repeatedly.

Finally Fig. 3c shows the coding required when embedded files, some of which may also be associated to content or structure, are present within a PDF document. One sees that the array (object 1859) of indirect object references in the lower part of Fig. 3c refer to the same /Filespec objects (27 and 29) as the *named resources* (object 20) in the upper part of Fig. 3b. These are the same references using /AF keys seen in Fig. 1b and Fig. 2b to the objects themselves in Fig. 2c.

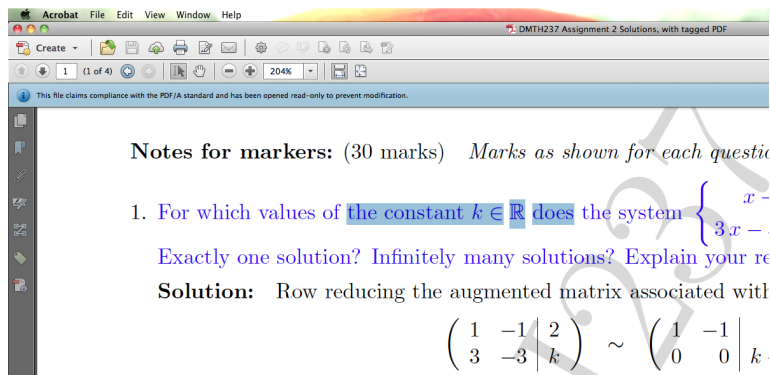
This mechanism makes it easier for a PDF reader to determine that there are files associated to a particular piece of content, by simply encountering the /AF tag linked with a *named resource*. This should work perfectly well with a PDF file that is not fully tagged for structure. However, if the content is extended (e.g., crosses a page-boundary) then it may be harder for a PDF writer to construct the correct content stream, properly tagging two or more portions.

## 4 Access-tags: attaching L<sup>A</sup>T<sub>E</sub>X source to ‘fake’ spaces.

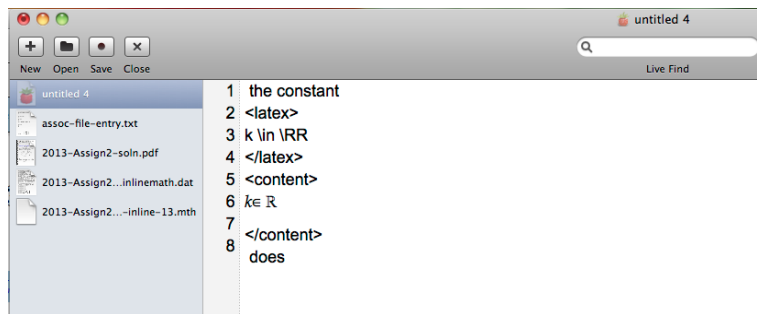
A third method allows inclusion of the L<sup>A</sup>T<sub>E</sub>X source of mathematics so that it may be readily extracted, using just the usual Select/Copy/Paste actions. This works with some existing PDF reader applications, including the freely available ‘Adobe Reader’. It is achieved by making use of the /ActualText attribute [5, §14.9.4] for a piece of *marked content*, whether or not structure tagging is present. It can be done by existing PDF-writing software that supports tagging of content, as in Sect. 2.1, and specification of a value for the /ActualText attribute.

Fig. 4 shows how this works, by tagging a ‘fake space’ character immediately before mathematical content, and another immediately afterwards. By selecting (see Fig. 4a) then Copy/Paste the content into text-editing software, the result should be similar to Fig. 4b. The PDF reader must recognise<sup>9</sup> /ActualText and replace the copied content (e.g. a single font character) with its value.

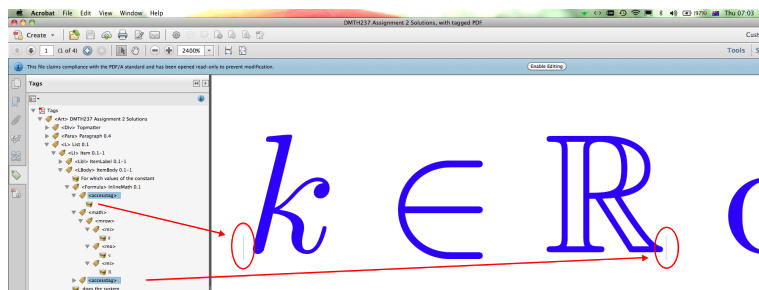
<sup>9</sup> Adobe’s ‘Reader’ and ‘Acrobat Pro’ certainly do, along with other software applications, but Apple’s standard PDF viewers currently do not support /ActualText.



(a) Selection across mathematical content



(b) Pasted text from the selection



(c) Access-tags selected in the 'Tags' tree, to show 'fake spaces'.

**Fig. 4.** This shows how the selection in (a), when copied and pasted into a text file, recovers the L<sup>A</sup>T<sub>E</sub>X source (b) that was used to specify the visual appearance of the mathematical content. In (c) we see structure within a `/Formula` node, (see also Fig. 5b) with leaf-nodes of `/accesstag` structure nodes being *marked content* of type `/AccessTag`. This consists of a single space character carrying an `/ActualText` attribute which holds the replacement text; as seen explicitly in the coding shown in Fig. 5a. The 'fake spaces' are very narrow; when selected they can be seen very faintly in (c) within the ovals indicated, at the outer edge of the the bounding rectangles of the outermost math symbols.

```

/AF /inline-1 BDC
1 0 0 1 51.508 0 cm
/AccessTag <</MCID 8 /ActualText
(\015<latex>\015k\134in\134RR
\015</latex>
\015<content>\015)
>>BDC
BT
/F79 1 Tf
[( )TJ
ET
EMC
/mi <</MCID 9 /ActualText<FEFFD835DC58>
/Alt( k )
>>BDC
BT
/F30 10.9091 Tf
[(k)TJ
ET
EMC
/mo <</MCID 10 /ActualText<FEFF2208>
/Alt( as element of )
>>BDC
1 0 0 1 3.03 0 cm
BT
/F33 10.9091 Tf
[(2)TJ
ET
EMC
/mi <</MCID 11 /ActualText<FEFF211D>
/Alt( real numbers )
>>BDC
BT
/F42 10.9091 Tf
[(R)TJ
ET
EMC
/mi <</MCID 12
/ActualText (\015</content>\015)
>>BDC
BT
/F79 1 Tf
[( )TJ
ET
EMC
EMC

```

fake space

see Fig. 1b

see Fig. 1b

see Fig. 1b

see Fig. 1b

fake space

portion of page stream, see Fig. 3b

(a) Complete portion of the content stream corresponding to the mathematics shown as selected in Fig. 4a. This is the same content as in Figures 1b and 3a but with the ‘...’ parts there now showing the /AccessTag coding of a ‘fake space’ with /ActualText attribute. The /AF ... BDC ... EMC wrapping of Fig. 3a is also shown. Being part of the document’s content, these space characters are also assigned /MCID numbers to be linked to structure nodes, as in (b) below.

```

122 0 obj
<<
/K [ 12 ]
/Pg 5 0 R
/Type/StructElem/S/accesstag
>>
endobj
...
120 0 obj
<<
/K [ 121 0 R ]
/P 112 0 R
/Type/StructElem/S/math
/A<</0/XML-1.00
/xmlns(http://www.w3.org/1998/Math/MathML)
/display(inline)>>
>>
endobj
...
113 0 obj
<<
/K [ 8 ]
/Pg 5 0 R
/Type/StructElem/S/accesstag
>>
endobj
112 0 obj
<<
/K [
113 0 R
120 0 R
122 0 R
]
/P 109 0 R
/Type/StructElem/S/Formula
/ID(Math0.1)/T(InlineMath 0.1)
/AF [ 27 0 R 29 0 R ] /A <</0/XML-1.01 >>
>>
endobj

```

<accesstag>

see Fig. 3b

see Fig. 3b

<math>

Kids

<Formula>

see Fig. 2a

see Fig. 2b,c

(b) Portion of the structure tree as in Fig. 1a, but now showing how the ‘fake spaces’ can be linked to structure nodes, here /accesstag. The missing portions of Fig. 1a, indicated there by ‘...’ are now filled-in, but leaving out other parts whose purpose has already been explained. Fig. 4c, shows the tagging opened out within the ‘Tags’ navigation panel, with the /accesstag structure nodes selected.

**Fig. 5.** File content included as /ActualText for a ‘fake space’, which itself can be tagged as *marked content* linked to an /accesstag structure node.

The *PDF name* `/AccessTag` tags a single ‘space’ character `[( )]Tf` as *marked content*, having replacement text in the `/ActualText` attribute; see Fig. 5a. An `/MCID` number is not needed for this technique to work; these *are* included in the example document<sup>1</sup> which is fully tagged<sup>10</sup>.

We refer to these tagged ‘fake spaces’ as ‘Access-tags’, since a motivation for their use is to allow Assistive Technology (e.g., a Braille reader) access to the  $\LaTeX$  source of mathematical content. The spaces are ‘fake’ in the sense that they are just 1pt in height and have nearly zero width. This makes them hard to select by themselves, but nearly impossible to separate from the mathematical content which they accompany; see Fig. 4c. They act as ordinary spaces when copied, but this is substituted with the `/ActualText` replacement, if supported. Another aspect of their ‘fakeness’ is that they take no part in the typesetting, when using pdf- $\LaTeX$  (post-2014). Of course the same idea could be implemented in different ways with different PDF-producing software.

One places into the initial ‘Access-tag’ text of the  $\LaTeX$  source — with care given to encode special characters<sup>11</sup> — as the value of its `/ActualText` attribute. The source coding is preceded by the string `<latex>` and followed by `</latex>` and `<content>`, with return-characters (in octal `\015`) to allow these ‘delimiters’ to ultimately copy onto lines by themselves. The trailing ‘Access-tag’ just takes `</content>`. As a result, the eventual Paste gives text as shown in Fig. 4b.

Assistive Technology (e.g., a Braille reader) works either by (a) emulating Copy/Paste of on-screen portions of the document’s window; or (b) by directly accessing the ‘Accessible Text’ view of the PDF document’s contents. In both cases the `/ActualText` replacements are extracted. (The ‘Accessible Text’ view can be exported directly using Adobe’s ‘Acrobat Pro’ software, see [12].) For mathematical symbols, where Figures 1b, 3a and 5a show the presence of both `/Alt` and `/ActualText` attributes, then the `/Alt` contributes to the ‘Accessible Text’, whereas `/ActualText` supplies what is copied to the Clipboard for Copy/Paste. In either case, a human reader when encountering `<latex>` on a line by itself can choose whether to read (or listen to) the following lines of  $\LaTeX$  coding, else use a Find action to skip down to where the next `</latex>` occurs. This is followed by a line containing `<content>`. Similarly the human can read/listen or skip down to where `</content>` occurs.

*Acknowledgements:* The author wishes to acknowledge James Davenport and Emma Cliffe, at the University of Bath, for valuable discussions regarding the needs of mathematics students having a visual disability. The ‘fake space’ idea emerged as a result. Thanks also to reviewers for suggesting, among other ideas, inclusion of the PDF language overview as in Sect. 2, allowing later sections to be written more succinctly. Finally, I wish to thank members of ISO TC 171 for encouragement and support regarding ‘fake spaces’ and other aspects.

“The committee is really happy to have someone actually implementing math accessibility in PDF . . .” [15] — Neil Soiffer, Senior Scientist, Design Science Inc., 12 April 2014.

<sup>10</sup> In PDF 2.0 this will also need an association of `/accessTag` to `/Custom` within the `/RoleMap` dictionary. The `/AccessTag PDF name` can be replaced by `/Span`.

<sup>11</sup> Octal codes: `\134` for backslash, `\050` for ‘(’ and `\051` for ‘)’, `\015` for line-end.





## References

1. Adobe Systems Inc.; PDF Reference 1.7, November 2006. Also available as [5].  
[http://www.adobe.com/devnet/pdf/pdf\\_reference.html](http://www.adobe.com/devnet/pdf/pdf_reference.html).
2. ISO 19005-1:2005; Document Management — Electronic document file format for long term preservation — Part 1: Use of PDF 1.4 (PDF/A-1); Technical Committee ISO/TC 171/SC 2 (Sept. 2005). Revisions via Corrigenda: ISO 19005-1:2005/Cor 1:2007 (March 2007); ISO 19005-1:2005/Cor 2:2011 (Dec. 2011).  
[http://www.iso.org/iso/catalogue\\_detail?csnumber=38920](http://www.iso.org/iso/catalogue_detail?csnumber=38920).
3. ISO 19005-2:2011; Document Management — Electronic document file format for long term preservation — Part 2: Use of ISO 32000-1 (PDF/A-2); Technical Committee ISO/TC 171/SC 2 (June 2011).  
[http://www.iso.org/iso/catalogue\\_detail?csnumber=50655](http://www.iso.org/iso/catalogue_detail?csnumber=50655).
4. ISO 19005-3:2012; Document Management — Electronic document file format for long term preservation — Part 3: Use of ISO 32000-1 with support for embedded files (PDF/A-3); Technical Committee ISO/TC 171/SC 2 (October 2012).  
[http://www.iso.org/iso/catalogue\\_detail?csnumber=57229](http://www.iso.org/iso/catalogue_detail?csnumber=57229).
5. ISO 32000-1:2008; Document management — Portable document format (PDF 1.7); Technical Committee ISO/TC 171/SC 2 (July 2008). Also available as [1].  
[http://www.iso.org/iso/catalogue\\_detail?csnumber=51502](http://www.iso.org/iso/catalogue_detail?csnumber=51502).
6. ISO 32000-2-20140220; Document management — Portable document format — Part 2: PDF 2.0; Technical Committee ISO/TC 171/SC 2, in draft form (Feb. 2014).
7. Document management applications — Electronic document file format enhancement for accessibility — Part 1: Use of ISO 32000-1 (PDF/UA-1); Technical Committee ISO/TC 171/SC 2 (July 2012). [http://www.iso.org/iso/home/store/catalogue\\_tc/catalogue\\_detail.htm?csnumber=54564](http://www.iso.org/iso/home/store/catalogue_tc/catalogue_detail.htm?csnumber=54564).
8. Technical Implementation Guide; AIIM Global Community of Information Professionals. <http://www.aiim.org/Research-and-Publications/standards/committees/PDFUA/Technical-Implementation-Guide>. Also available as [7].
9. Moore, Ross R.; Tagged Mathematics in PDFs for Accessibility and other purposes, in *CICM-WS-WiP 2013, Workshops and Work in Progress at CICM*, CEUR Workshops Proceedings. <http://ceur-ws.org/Vol-1010/paper-01.pdf>.
10. Moore, Ross R.; Ongoing efforts to generate “tagged PDF” using pdfI<sub>E</sub>X, in *DML 2009, Towards a digital Mathematics Library, Proceedings*, Petr Sojka (editor), Muni Press, Masaryk University, 2009. ISBN 978-80-20-4781-5. Reprinted as: TUGboat, Vol.30, No.2 (2009), pp. 170–175.  
<http://www.tug.org/TUGboat/tb30-2/tb95moore.pdf>.
11. Moore, Ross R.; serendiPDF, with searchable math-fields in PDF documents; TUGboat, Vol.23, No.1 (2002), pp. 65–69.  
<http://www.tug.org/TUGboat/tb23-1/moore.pdf>.
12. Moore, Ross R.; DMTH237 Assignment 2 Solutions, with tagged PDF; web address for downloading the example document having attachments and ‘fake spaces’; also other views of this document’s contents exported using ‘Acrobat Pro’ software.  
<http://rutherglen.science.mq.edu.au/~maths/CICM/>
13. Raman, T.V.; An Audio View of (L<sup>A</sup>)T<sub>E</sub>X Documents; TUGboat, Vol.13, No.3 (1992), pp. 372–379; <http://tugboat.tug.org/TUGboat/tb13-3/raman.pdf>.
14. Raman, T.V.; An Audio View of (L<sup>A</sup>)T<sub>E</sub>X Documents—Part II; TUGboat, Vol.16, No.3 (1995), pp. 310–314; <http://tugboat.tug.org/TUGboat/tb16-3/tb48rama.pdf>.
15. Soiffer, N.; posting to ‘PDF-Access’ mailing list server; 12 April 2014.  
<http://listserv.aiim.org/scripts/wa.exe?A0=PDF-ACCESS>