opentext[™] | Discovery

OpenText[™] Axcelerate[™] 5.15

Incoming Data Specification

Revised: 2018-Oct-12

Contents

1 Incoming Data Specification Overview	3
1.1 Delivery and Notice for OpenText Hosting	
1.2 Example Reference Files	
1.2.1 documents.csv File	4
2 Specification Details	7
2.1 Metadata File	7
2.1.1 Consistent structure	7
2.1.2 Date format	7
2.1.3 Unique key	8
2.1.4 Frequently Used Fields	
2.1.4.1 Tags assigned during review	13
2.2 Native Files	13
2.3 Text Files	13
2.4 Images	13
2.4.1 Fields in the Opticon Load File for Images	16
2.5 Filenames and Paths	17
2.6 Encoding	17
3 Changes to this Document	18
4 Terms of Use	19

1 Incoming Data Specification Overview

This document describes commonly used requirements for incoming data, for Axcelerate OnDemand and Axcelerate Cloud or licensed Axcelerate installations. The described examples and fields are typical for a third party production, including:

Structured Data:

native, imaged or mixed productions delivered in a structured format, pre-processed by another vendor

Unstructured Data:

raw native deliveries of loose files and emails or containers

 Recommind databases: Axcelerate Ingestion or Axcelerate Review & Analysis databases created by customers or vendors with licensed Axcelerate installations.

To ensure the fastest upload, incoming data should include these files:

- Metadata file
- · Opticon load file for images
- TIFF images and/or native files
- Extracted text files (one per document)
- 0

Note: Deviations from this specification and the structures outlined below may incur additional charges and not be processed in an expeditious fashion. Depending on the actual installation, there may be additional requirements, especially with regard to metadata fields, which have to be specified separately.

1.1 Delivery and Notice for OpenText Hosting

The files should be placed on either a USB hard drive, optical media such as DVD and CD, or with prior approval a secure sFTP site. Advance notice in writing that media is being delivered; cover letters detailing with which database(s) to associate the media and outlining any special load requests is appreciated. If the media is encrypted, passwords should be sent under separate cover.

1.2 Example Reference Files

This example delivery (a third party production) consists of

- A metadata file called documents.csv, in CSV format. For better display, | is used as separator.
- An Opticon file called documents.opt for references to images.
- Four documents with their respective image files.

The documents are delivered in this folder structure:

Main folder	Contained files/folders	Contained folders	Contained files
\Delivery	documents.dat		
	documents.opt		
	\nativefile	\001\	123.doc
			555.doc
			XYZ.doc
			LMN.doc
	\images	\001\	IMG_11.TIF
			IMG_12.TIF
			IMG_31.TIF
			IMG_106.TIF

The content of the reference files is shown in tables here, for better visibility.

1.2.1 documents.csv File

File content	Explanation
BEGDOC ENDDOC BEGATTACH ENDATTACH TEXTPATH NATIVELINK	Header
ABC_0000001 ABC_0000002 ABC_ 0000001 ABC0000007 \nativefile\001\123.doc	First document has 2 pages. Its attachment family starts with the first document page (ABC_ 0000001) and ends with the last page (ABC_0000007) of the attachment.

File content	Explanation
ABC_0000003 ABC_0000007	Second document has 5 pages. It
ABC_0000001 ABC_0000007	belongs to the same attachment
\nativefile\001\555.doc	anning as the first document.
ABC_0000008 ABC_0000009	Third document has 2 pages. It
ABC_0000008 ABC_0000015	belongs to the same attachment
\nativefile\001\XYZ.doc	0000015) as the fourth document.
ABC_0000010 ABC_0000015	Fourth document has 6 pages. It
ABC_0000008 ABC_0000015	belongs to the same attachment
\nativefile\001\LMN.doc	

Documents.opt file

The first reference in a line refers to the respective document listed in <code>doc-uments.csv</code>.

References	Explanation
ABC_0000001,CD_ 001,\IMAGES\001\IMG_ 11.TIF,Y,,,2	First image out of 2 for the first document. Y marks the first image for a document, 2 is the (optional) number of images for one document.
ABC_0000001,CD_ 001,\IMAGES\001\IMG_ 12.TIF,,,,	Second image out of 2 for the first doc- ument.
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 31.TIF,Y,,,5	First image out of 5 for the second doc- ument.
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 32.TIF,,,,	
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 33.TIF,,,,	

References	Explanation
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 34.TIF,,,,	
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 35.TIF,,,,	
ABC_0000008,CD_ 001,\IMAGES\001\IMG_ 81.TIF,Y,,,2	First image out of 2 for the third document
ABC_0000008,CD_ 001,\IMAGES\001\IMG_ 81.TIF,,,,	
ABC_0000010,,CD_ 001,\IMAGES\001\IMG_ 101.TIF,Y,,,2	First image out of 6 for the fourth document
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 102.TIF,,,,	
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 103.TIF,,,,	
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 104.TIF,,,,	
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 105.TIF,,,,	
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 106.TIF,,,,	

2 **Specification Details**

2.1 Metadata File

The metadata file contains details about the incoming records, such as Bates number, author/recipient or other fields for names, paths, and attachment information. It must be consistent in structure and delimiters.

A metadata flat data file will be provided for each data source, or CSV Merge. Multiple custodians may be aggregated into a single file, provided the file contains a custodian field containing the different values.

2.1.1 Consistent structure

- 1. Fields will match the fields specified under "Frequently Used Fields" on page 9 in name and content.
- 2. Delimiters Concordance default delimiters
 - Text delimiter "p": Hex (FE), Unicode (U+00FE), Decimal (254)
 - Field separator (not displayable or displayed as "DC4"): Hex (14), Unicode (U+0014), Decimal (20)

If other delimiters are used, this must be explicitly specified.

- 3. All rows will contain the same number of delimiters and fields.
- 4. The multi-value field delimiter is a semicolon (U+003B) and must be consistent across all fields.
- 5. The first line contains a header row with field names. The table below defines the accepted field name labels for the header row.
- 6. Extracted text is delivered separately from the metadata file as loose text files, one per document.

2.1.2 Date format

- 1. Date formats must be consistent across all fields, that is, the sent date should have the same format as the last modified date, for example.
- 2. Dates and times can be concatenated into a single field, if nothing else is specified. They may occur in two different fields, if this is required.
- 3. The default date format is configurable. If nothing else is specified, use MM/DD/YYYY and HH:MM:ss (zzz).
- 4. If a time is not available, such as the estimate date for a coded document, then 12:00 am, or 00:00 should be assigned, that is, 12/21/1999 00:00.

5. Invalid times or dates or missing times or dates in important date fields will be replaced by 01/01/1901 00:00 by default when they are loaded into an Axcelerate project.

2.1.3 Unique key

- 1. One field must contain a value unique across the Axcelerate project. Typically this is a Bates number or control number. This should be a unique value for the record across all deliveries.
- This key cannot have spaces, but any alpha-numeric character and all ASCII characters are accepted, except these: < > & /\?*"\$ |:,;
- 3. A unique key combination is also needed for attachment families. All documents of an attachment family must have the same attachment start and and attachment end number. Otherwise, attachment families cannot be identified.

File content	Explanation
BEGDOC ENDDOC BEGATTACH	Header
ENDATTACH TEXTPATH NATIVELINK	
ABC_0000001 ABC_0000002	First document has 2 pages. Its attach-
ABC_0000001 ABC_0000007	ument family starts with the first doc-
\nativefile\001\123.doc	with the last page (ABC_0000007) of the attachment.
ABC_0000003 ABC_0000007	Second document has 5 pages. It
ABC_0000001 ABC_0000007	belongs to the same attachment family
\nativefile\001\555.doc	
ABC_0000008 ABC_0000009	Third document has 2 pages. It belongs
ABC_0000008 ABC_0000015	to the same attachment family (ABC
\nativefile\001\XYZ.doc	fourth document.
ABC_0000010 ABC_0000015	Fourth document has 6 pages. It
ABC_0000008 ABC_0000015	belongs to the same attachment family
\nativefile\001\LMN.doc	

Example: Unique keys and attachments in documents.csv

2.1.4 Frequently Used Fields

Axcelerate Field	Mandatory (M) /Op- tional (O)	Multiple values possible (Y/N)	Description
Beg Doc	Μ	Ν	Beginning control number for document (the unique key used for the data) Can also be used as external Bates number for productions. For details, see the <i>Axcel</i> - <i>erate Help Center Administrator Help</i> .
End Doc	Μ	Ν	Ending control number for document Amongst others, this field is needed for val- idation of productions with external Bates numbers.
Beg Attach	Μ	Ν	Beginning control number for first page of parent document
End Attach	Μ	Ν	Ending control number for last page of last attachment
Location	0	Ν	File system path or Internet URL, either to the loose file, or to the container the doc- ument belongs to (for example a PST archive).
Custodian	0	Y	Data's custodian, owner of the files

Except for the mandatory fields, additional specifications are possible.

Axcelerate Field	Mandatory (M) /Op- tional (O)	Multiple values possible (Y/N)	Description
Document Date	0	Ν	An aggregated date field based on the fol- lowing criteria:
			 For loose files: modification date/time (or creation date/time if last modified date is not available)
			 For emails: sent date/time (or delivery date/time if sent date is not available) For attachments: inherits the date/time from the parent email. This field is commonly known as <i>Sort Date</i>. The individual dates can also be provided separately and be mapped to the date fields below.
Modification Date	Modification O Date	Ν	Last modified date (stored by host file system)
			Note: This is a file system date, not the application date.
Creation Date	0	Ν	Creation date (stored by host file system)
			Note: This is a file system date, not the application date.
Sent Date	0	Ν	Email date/time sent
Application Last Modified Date	0	Ν	Last modified date (stored by the applic- ation)
Application Create Date	0	Ν	Creation date (stored by the application)

Axcelerate Field	Mandatory (M) /Op- tional (O)	Multiple values possible (Y/N)	Description
Document	0	Ν	Aggregated field, based on this information:
1 me			 <i>Title</i> metadata field (if available) or file- name of non-email files Subject for emails
			The individual title/filename/subject fields can also be provided separately and be mapped to the fields below.
Title	0	Ν	Title metadata field within non-email file
Filename	0	Ν	Filename of non-email file
Subject	0	Ν	Email subject
Sender	0	Ν	Email sender, or sender of a chat message
Recipient	0	Y	Email recipient(s), or chat message recip- ient(s)
Email CC	0	Y	Email CC(s)
Email BCC	0	Υ	Email BCC(s)
Importance	0	Ν	Email importance flag
Read/Unread	0	Ν	Email read/unread flag
Author	0	Ν	Author metadata field within non-email file
File Name	0	Ν	Filename of non-email file
File Exten- sion	0	Ν	File extension
File Size	0	Ν	File size in Bytes
Folder Name	0	Ν	Email folder (that is, folder within a PST or NSF file)

Axcelerate Field	Mandatory (M) /Op- tional (O)	Multiple values possible (Y/N)	Description
Message ID	0	Ν	Internet MessageID for emails. Always use in combination with References field for thread detection by header analysis.
References	0	Ν	References to other items for internet mes- sages. Always use in combination with the Message ID field for thread detection by header analysis.
MD5 Hash	0	Ν	The MD5 hash is based on actual file con- tent and, for emails, on composite of metadata fields for emails. If this field is filed, duplicate detection is pos- sible in the target system.
Store Name	0	Ν	Name of container file (PST name, NSF name, Opentext database name), including extension
	0	Y	Review tags or other work product assigned to documents during a previous review. The tags to be mapped must be communicated to OpenText prior to data being loaded. See "Frequently Used Fields" on page 9.
	0	Ν	Relative path to text file, for example \textfile\001\123.txt This field is mandatory if text files are part of the incoming data.
	Ο	Ν	Relative path to native file, for example \nativefile\001\123.doc This field is mandatory if native files are part of the incoming data. Caution: The path must not con- tain spaces.

opentext * | Axcelerate

2.1.4.1 Tags assigned during review

There is no specific "Tags" field, but a number of default and custom fields. If tags are required, a custom specification of these fields must be added to this standard specification.

Nested fields should be sent as individual fields with the main field as a prefix.

Example:

```
Responsive - Responsive Type
```

2.2 Native Files

Native files and references to them must meet the following requirements:

- 1. The incoming metadata file contains a relative path to the native file, the NATIVELINK field (see "Frequently Used Fields" on page 9).
- 2. Filenames matching a Bates number are acceptable, for example PROD_006789.xls.
- 3. There are no more than 1000 native files per directory.
- 4. The path to the native file has less than 255 characters, no spaces and only consists of ANSI characters.

2.3 Text Files

Extracted text files and references to them must meet the following requirements:

- 1. There is not more than one extracted text file per document, with the content of all document pages.
 - 0

Note: Multiple single-page text files for one document are not supported.

- 2. The character encoding in the text files must be consistent ideally UTF-8.
- 3. The incoming metadata file contains a relative path to the extracted text or OCR, in the TEXTPATH field. See: "Frequently Used Fields" on page 9.
- 4. There are no more than 1000 text files per directory.
- 5. The path to the text file has less than 255 characters, no spaces and only ANSI characters.
- 6. Filenames matching a Bates number are acceptable, for example PROD_006789.txt

2.4 Images

Images and the Opticon file (*.opt) must meet these requirements:

opentext * | Axcelerate

Black and white images:

single page TIFF 1bit color-depth Photometric interpretation: MinIsWhite (Black and White) Compression: CCITT – Group 4 300 dpi (default and recommended) Byte Order: little-endian Fill Order: TIFF fill order 1

0

Note: Any private tags are ignored during loading/merging.

Color images:

singe page TIFF 24bit color-depth Photometric interpretation: RGB Compression: LZW 300 dpi (default and recommended)

0

Note: Any private tags are ignored during loading/merging.

The image link is delivered separately from the metadata file in a file following the Opticon load file format specification.

The Opticon load file format is a text-delimited file containing all information necessary to link the image with the database. There is one line entry per image file.

The image file entries must be in correct order, that is, in the same order as documents occur in the metadata file. Pages must be in the same order as they occur in the documents.

The field delimiter is a comma (U+002C).

Example:

The following is a 5-image Opticon load file example. It details 4 documents with their images.

documents.opt

The first reference in a line refers to the respective document listed in documents.csv.

References	Explanation
ABC_0000001,CD_ 001,\IMAGES\001\IMG_ 11.TIF,Y,,,2	First image out of 2 for the first document. Y marks the first image for a document, 2 is the (optional) number of images for one document.
ABC_0000001,CD_ 001,\IMAGES\001\IMG_ 12.TIF,,,,	Second image out of 2 for the first doc- ument.
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 31.TIF,Y,,,5	First image out of 5 for the second doc- ument.
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 32.TIF,,,,	
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 33.TIF,,,,	
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 34.TIF,,,,	
ABC_0000003,CD_ 001,\IMAGES\001\IMG_ 35.TIF,,,,	
ABC_0000008,CD_ 001,\IMAGES\001\IMG_ 81.TIF,Y,,,2	First image out of 2 for the third document

References	Explanation
ABC_0000008,CD_ 001,\IMAGES\001\IMG_ 81.TIF,,,,	
ABC_0000010,,CD_ 001,\IMAGES\001\IMG_ 101.TIF,Y,,,2	First image out of 6 for the fourth doc- ument
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 102.TIF,,,,	
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 103.TIF,,,,	
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 104.TIF,,,,	
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 105.TIF,,,,	
ABC_0000010,CD_ 001,\IMAGES\001\IMG_ 106.TIF,,,,	

2.4.1 Fields in the Opticon Load File for Images

Field	Mandatory (M) /Op- tional (O)	Description
ALIAS	Μ	Should match your BEGDOC field (see "Frequently Used Fields" on page 9) for the first page of a record, subsequent lines are the interior pages of the doc- ument, up to the next Unique key

opentext ** | Axcelerate

Field	Mandatory (M) /Op- tional (O)	Description
VOLUME	0	This entry is the name of the volume where the image resides. This is typically the volume name of a CD or server.
PATH	Μ	This is the full path and file name (and extension) of the image. File name and path should only consist of ANSI characters and must have no spaces. They have less than 255 characters.
DOC_BREAK	Μ	Enter a 'Y' to denote whether this image marks the beginning of a document.
FOLDER_ BREAK	0	A 'Y' denotes that this image marks the beginning of a folder. (not used)
BOX_BREAK	0	A 'Y' denotes that this image marks the beginning of a box. (not used)
PAGES	0	This entry is the number of document pages. (not used)

2.5 Filenames and Paths

Filenames and paths can have any character allowed for filenames in Windows, but must not contain spaces. They should only consist of ANSI characters. The paths must not have more than 255 characters.

2.6 Encoding

Opticon files are ANSI encoded. For all other files, UTF-* encoding is expected.

3 Changes to this Document

Date	Topic title	Text before change	Text after change	Remarks
2015- 03-18	"Frequently Used Fields" on page 9		Field list was reworked.	
2015- 03-27	"Frequently Used Fields" on page 9		Added Application Create Date and Application Last Modified Date to field list.	
2015- 10-13	"Frequently Used Fields" on page 9	Email From Email To	Sender Recipient	
2015- "Text Files" - 11-05 on page 13	"Text Files" on page 13	-	There is not more than one extracted text file per document, with the content of all document pages.	
	• Note: Multiple single-page text files for one document are not supported.			
2017- 01-12	"Frequently Used Fields" on page 9	Internet MessageID for emails. -	Internet MessageID for emails. Always use in combination with References field for thread detec- tion by header analysis.	
2017- 01-12	"Frequently Used Fields" on page 9	-	Added <i>References</i> field to the list.	
2017- 05-2	"Frequently Used Fields" on page 9	-	Added hint on external Bates number usage for the BEC DOC and END DOC fields.	

4 Terms of Use

OpenText[™] Axcelerate[™] 5.15

Incoming Data Specification Rev.: 2018-Oct-12

This documentation has been created for OpenText[™] Axcelerate[™] 5.15.

It is also valid for subsequent software versions as long as no new document version is shipped with the product or is published at <u>https://knowledge.opentext.com</u>.

Open Text Corporation

275 Frank Tompa Drive, Waterloo, Ontario, Canada, N2L 0A1 Tel: +1-519-888-7111 Toll Free Canada/USA: 1-800-499-6544 International: +800-4996-5440 Fax: +1-519-888-0677 Support: <u>https://support.opentext.com</u> For more information, visit <u>https://www.opentext.com</u>

Copyright © 2018 Open Text. All Rights Reserved.

Trademarks owned by Open Text.

One or more patents may cover this product. For more information, please visit, <u>https://www.opentext.com/patents</u>.

Disclaimer

No Warranties and Limitation of Liability

Every effort has been made to ensure the accuracy of the features and techniques presented in this publication. However, Open Text Corporation and its affiliates accept no responsibility and offer no warranty whether expressed or implied, for the accuracy of this publication.