

Content Manager OnDemand for i
Version 7 Release 2

Common Server Indexing Reference



Content Manager OnDemand for i
Version 7 Release 2

Common Server Indexing Reference



Note

Before using this information and the product it supports, read the information in “Notices” on page 65.

This edition applies to version 7, release 2 of IBM Content Manager OnDemand for i (product number 5770-RD1) and to all subsequent releases and modifications until otherwise indicated in new editions.

© **Copyright IBM Corporation 2001, 2014.**

US Government Users Restricted Rights – Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

Contents

IBM Navigator for i v

Accessibility information for Content Manager OnDemand vii

System requirements ix

**Content Manager OnDemand OS/400
indexer 1**

Indexing concepts 1
Indexing parameters 2
Using BREAK=YES versus BREAK=NO in indexer
parameters 4
Controlling maximum number of pages per group. 5
Defining multi-key indexes 5
 Define a multi-key index example 5
Defining transaction fields 8
 Define a transaction report example 9
Assigning default index values 11
Defining text search fields 12
Handling SCS spooled files that have AFP overlays 13
Using a mask when defining applications fields 13
Using Tag Logical Elements (TLEs) 14

**Content Manager OnDemand PDF
indexer 17**

How OnDemand uses index information 19
Processing PDF input files with the graphical
indexer 20
Manually indexing input data 22
 Indexing concepts 23
 Coordinate system 23
 Indexing parameters 23
Indexing with Metadata Indexes 26
How to create indexing parameters 27
PDF resource collection 28
PDF indexing system requirements 29
 Specifying the location of Adobe fonts 29
 PDF indexing limitations 29
 Input data requirements 30
 National language support for indexed PDF
 documents. 30
Parameter reference 31
 BOOKMARKS 31
 COORDINATES 32
 FIELD 32
 FONTLIB 36
 HEX 36

INDEX 37
INDEXDD 38
INDEXMODE 38
INDEXSTARTBY 39
INPUTDD 40
MSGDD 40
OUTPUTDD 41
PARMDD 41
REMOVERES 42
RESOBJDD 42
RESTYPE 43
TEMPDIR 43
TRACEDD 43
TRIGGER 43
Message reference 46
ARSPDOCI reference 47
 Purpose 47
 Syntax 47
 Description 47
 Parameters 47
 IFS location 47
ARSPDUMP reference 48
 Purpose 48
 Syntax 48
 Description 48
 Parameters 48
 Examples 49
 IFS location 49
Trace facility 49

**Content Manager OnDemand generic
indexer 53**

Loading data 53
Specifying the parameter file 56
 CODEPAGE: 56
 COMMENT: 56
 GROUP_FIELD_NAME: 57
 GROUP_FIELD_VALUE: 57
 GROUP_FILENAME: 58
 GROUP_LENGTH: 59
 GROUP_OFFSET: 60
 Parameter file examples 60
 Additional indexing topics 62

Notices 65

Trademarks 67

Index 69

IBM Navigator for i

IBM Navigator for i is a powerful graphical interface for managing your IBM® i servers.

IBM Navigator for i functionality includes system navigation, configuration, planning capabilities, and online help to guide you through your tasks. IBM Navigator for i makes operation and administration of the server easier and more productive and is the only user interface to the new, advanced features of the operating system. It also includes Management Central for managing multiple servers from a central system.

Accessibility information for Content Manager OnDemand

For complete information about accessibility features that are supported by this product, see the IBM *Content Manager OnDemand for i: Common Server Administration Guide*.

System requirements

For administrative client system requirements, see: <http://publib.boulder.ibm.com/infocenter/prodguid/v1r0/clarity-reports/report/html/softwareReqsForProduct?deliverableId=1301098213267>

Content Manager OnDemand OS/400 indexer

The Content Manager OnDemand OS/400® indexer is the most common Content Manager OnDemand indexer used for IBM i spooled files. The OS/400 indexer is called by the ADDRPTOND command for SCS, SCS-extended, Advanced Function Presentation (AFP), and Line spooled files. You use the Content Manager OnDemand administrative client's graphical indexing tool to define the index criteria that the OS/400 indexer uses to locate and create index data for your spooled files.

The graphical tool can be invoked in one of two ways:

- By clicking the Select Sample Data button within the Report Wizard, or
- Selecting Sample Data and clicking the Modify button on the Indexer Information panel while creating a Content Manager OnDemand application definition

Content Manager OnDemand uses the OS/400 indexer by default for SCS, SCS-extended, AFP, and Line spooled files. See the Report Wizard section in the Introduction of the *IBM Content Manager OnDemand for i: Common Server Administration Guide* for more information on the Report Wizard. See the section on Adding the Application in the Examples chapter of the *IBM Content Manager OnDemand for i: Common Server Administration Guide* for more information on defining an application without using the Report Wizard.

Indexing concepts

Indexing parameters include information that allows Content Manager OnDemand to identify key items in the print data stream and create index elements pointing to these items. You can specify the index information that allows Content Manager OnDemand to segment the data stream into individual items called groups. A group is a collection of one or more pages. You define the bounds of the collection; for example, a bank statement, insurance policy, phone bill, or other logical segment of a report file. A group can also represent a specific number of pages in a report. For example, you might decide to segment a 10,000 page report into groups of 100 pages. Content Manager OnDemand creates indexes for each group. Groups are determined when the value of an index changes (for example, account number), or when the maximum number of pages for a group is reached.

Index data is made up of an attribute name (for example, Customer_Name) and an attribute value (for example, Frank Booth), with a defined tag that identifies the location of the data on the print page. For example, the Account_Number tag with the pointer 1,21,16 means Content Manager OnDemand can expect to find Account_Number values starting in column 21 of specific input records. Content Manager OnDemand collects 16 bytes of information starting at column 21 and adds it to a list of attribute values found in the input. Content Manager OnDemand creates an index file when you index report files. The index file includes index elements that contain the offset and length of a group. Content Manager OnDemand calculates an index element for every group found in the input file. Content Manager OnDemand then writes the attribute values extracted from the input file to the index file.

Indexing parameters

Indexing parameters can contain indexing, conversion, and resource collection parameters, options, and values. For most reports, Content Manager OnDemand requires three indexing parameters to extract or generate index data:

- **TRIGGER** Content Manager OnDemand uses triggers to determine where to locate data. A trigger instructs Content Manager OnDemand to look for certain information in a specific location in the report file. When Content Manager OnDemand finds a record in the data stream that contains the information specified in the trigger, it can begin to look for index information.
 - Content Manager OnDemand compares data in the report file with the set of characters specified in a trigger, byte for byte.
 - A maximum of 16 triggers can be specified.
 - All fixed group triggers must match before Content Manager OnDemand can generate index information. However, floating triggers can occur anywhere in the data stream. That is, index data based on a floating trigger can be collected from any record in the report file.
- **FIELD** The field parameter identifies the location, offset, and length of the data Content Manager OnDemand uses to create index values.
 - Field definitions are based on TRIGGER1 by default, but can be based on any of 16 TRIGGER parameters.
 - A maximum of 128 fields can be defined.
 - A field can also specify all or part of the actual index value stored in the database.
- **INDEX** The index parameter is where you specify the attribute name, identify the field or fields on which the index is based, and specify the type of index that Content Manager OnDemand generates. For the group-level indexes Content Manager OnDemand stores in the database, you should name the attributes the same as the application group database field names.
 - Content Manager OnDemand can create indexes for a page, group of pages, and the first and last sorted values on a page or group of pages. Content Manager OnDemand stores group-level index values in the database. Users can search for items using group-level indexes. Page-level indexes are stored with the document (for example, a statement). After retrieving a document that contains page-level indexes, you can move to a specific page by using the page-level indexes.

Content Manager OnDemand can only generate this type of page-level information when converting the input data to AFP. This type of page-level information is generated by specifying the CONVERT=YES and INDEXOBJ=ALL parameters, and by creating an index field with the TYPE=PAGE option.
 - You can concatenate field parameters to form an index.
 - A maximum of 128 index parameters can be specified.

Content Manager OnDemand creates a new group and extracts new index values when one or more of the fixed group index values change, or the GROUPMAXPAGES value is reached.

```

-----1-----2-----3-----4-----5-----6-----7-----8-----9
01                                     Page 0001
1
2                                     Jack Straw
3                                     4 Buxanchange Way
4                                     Wichitaw KS 99999-9999
5
6                                     Statement Date: 06/15/07
7                                     Account Number: 1234-5678-9876-0000
8
9                                     Balance: $2,984.17

```

Figure 1. Indexing a report

The following indexing parameters can be used to generate index data for the report shown in Figure 1. The TRIGGER definitions tell Content Manager OnDemand how to identify the beginning of a group in the input. Content Manager OnDemand requires two TRIGGER definitions to identify the beginning of a group (statement) in the sample file. For example:

- TRIGGER1 looks for a 1 in the first byte of each input record.
- TRIGGER2 looks for the string Page 0001 in column 72 of the same record.

Together, the triggers uniquely identify the start of a statement in the report.

The FIELD definitions determine the location of the index values in a statement. Fields are based on the location of trigger records. For example:

- FIELD1 identifies customer name index values, beginning in column 40 of the second record following the TRIGGER1 record.
- FIELD2 identifies the statement date index values, beginning in column 56 of the sixth record following the TRIGGER1 record.
- FIELD3 identifies the account number index values, beginning in column 56 of the seventh record following the TRIGGER1 record.

An INDEX definition identifies the attribute name of the index field. Indexes are based on one or more field definitions. For example:

- INDEX1 identifies the attribute name custnam, for values extracted using FIELD1.
- INDEX2 identifies the attribute name sdate, for values extracted using FIELD2.
- INDEX3 identifies the attribute name acctnum, for values extracted using FIELD3.

The following table lists the maximum values for certain indexing attributes:

| Indexing attribute | Maximum value |
|--|---------------|
| Maximum number of lines per spooled file page (greater than printer file maximum due to allowance for overprint lines) | 510 |
| Maximum page width (positions per line) | 378 |
| Maximum number of triggers per page (for documents not using multi-key) | 16 |
| Maximum number of index values per page (for documents not using multi-key) | 128 |
| Maximum number of fields per page (for documents not using multi-key) | 128 |
| Maximum number of triggers per page (for multi-key documents) | 512 |
| Maximum number of index values per page (for multi-key documents) | 1024 |

| Indexing attribute | Maximum value |
|---|------------------|
| Maximum number of fields per page (for multi-key documents) | 1024 |
| Maximum number of index values per group (document) (for multi-key documents) | 9999 |
| Maximum size of an AFP resource segment | 16,000,000 bytes |
| Maximum size of any single AFP resource | 16,000,000 bytes |

Important: For the Maximum page width indexing attribute, when using the Print Text for OnDemand (**PRTTXTOND**) command with **STMF(*NONE)**, which directs the output to a spooled file instead of a stream file, the maximum page width is 372.

Using **BREAK=YES** versus **BREAK=NO** in indexer parameters

A group is a set of pages that logically belong together. For example, all the pages in a single bank statement could comprise a group. A group is a single document, or a *segment*, as it was known in Spool File Archive. A group break is the process of closing the current group and starting a new group. In Spool File Archive, this process was known as segmentation. For a specific group index, the **BREAK** setting determines whether the OS/400 indexer begins a new document when that index's value changes.

When you specify **BREAK=YES**, the OS/400 indexer begins a new group when the value of the field on which the index is based changes. **BREAK=NO** is useful when you define two or more fields and you want the OS/400 indexer to begin a new group only when the other of the two fields' value changes. Specify **BREAK=YES** only for the index that is based on the field that you want the OS/400 indexer to use to control the group break. Specify **BREAK=NO** for all the other indexes in the group.

To expand on the bank statement example, consider storing bank statements. Each statement begins with a change in account number from the previous statement. You defined indexes for Account Number, Customer Name, and Statement Date. Most likely, you want Account Number to be set to **BREAK=YES**, Customer Name to **BREAK=NO**, and Statement Date to **BREAK=NO**. Doing this ensures that a group break occurs only when Account Number changes. The corresponding indexer parameters in the Application definition might look like this:

```
INDEX1=X'C1838396A495A3D5A494828599',FIELD1,(TYPE=GROUP,BREAK=YES) /* AccountNumber */
INDEX2=X'C3A4A2A396948599D5819485',FIELD2,(TYPE=GROUP,BREAK=NO) /* CustomerName */
INDEX3=X'E2A381A385948595A3C481A385',FIELD3,(TYPE=GROUP,BREAK=NO) /* StatementDate */
```

The Content Manager OnDemand Administrator client's Report Wizard is designed to simplify the process of defining application groups, applications, and folders. The Wizard makes the assumption that any change in an index that is defined as **TYPE=GROUP** should cause a group break. Thus, it sets all index fields to **BREAK=YES**. If the index is based on a float trigger (**TYPE=FLOAT**), then **BREAK=NO** is set by default. Indexes based on float triggers can be changed to **BREAK=YES**, if necessary.

Note that if you have selected the Allow Multiple Values option, **BREAK** is automatically set to **NO** and should not be changed.

If you already archived data with all of your indexes set to **BREAK=YES**, you can still make this change. Changing some of your indexes from **BREAK=YES** to

BREAK=NO can be done at any time. As with any change to your indexer parameters, you should verify that your reports archive correctly after the change. Any reports already archived do not need to be rearchived; however, the change will only affect reports that are archived after the change is made.

Controlling maximum number of pages per group

You might want to set a maximum number of pages for each group that is indexed. Content Manager OnDemand can use the value of the GROUPMAXPAGES indexer parameter to determine the number of pages in a group. For example, you need to index a report consisting of thousands of pages of detail. If your BREAK=YES criteria do not result in small enough groups of pages (or segments) of the report, you can use GROUPMAXPAGES=100, for example, to force Content Manager OnDemand to close the current group and begin a new group for any group that reaches 100 pages. In other words, if the GROUPMAXPAGES value is reached before the value of a group index changes, Content Manager OnDemand forces the creation of a new group. If you do not specify a value for the GROUPMAXPAGES parameter, Content Manager OnDemand does not terminate the current group and begin a new group until the value of one of the fields named by an INDEX with BREAK=YES changes.

Defining multi-key indexes

Multi-key indexes can be used when an index value occurs multiple times within a single document. For example, invoices might have invoice number, customer number, and customer name defined as the first three index fields, each occurring once within a given invoice. Then you might also want to define item number as a multi-key index, since there might be multiple item numbers within one invoice. With multi-key support, an end-user could search by item number to find any invoice for a given item number, no matter where that item number appeared in the list of invoiced items. Without the multi-key capability, only the first item number on the page would be indexed.

To enable multi-key indexing, the keyword ALLOWMULTIPLEVALUES=YES must be added to each INDEX statement that is to have multiple values captured per document. For example:

```
INDEX2=X'97969596',FIELD2,(TYPE=GROUP,BREAK=NO,ALLOWMULTIPLEVALUES=YES)
```

The keyword would be added to the Content Manager OnDemand Application definition. Go to the Indexer Information tab, then click on Keyboard and then Modify to edit the Application's Indexer Parameters. Keyword ALLOWMULTIPLEVALUES is only valid when BREAK=NO. Keyword ALLOWMULTIPLEVALUES=YES is only supported by the OS/400 indexer. It is not supported by the PDF indexer. Also note that unlike the Content Manager OnDemand Spool File Archive multi-key rule, defining an index as multi-key does not require all subsequent index fields to also be defined as multi-key. In a Common Server environment, as is shown in the example, you can define an index as multi-key and then define another one below it that is not multi-key. However, a field used for a multi-key index must be found on or below the row containing the float trigger used to locate that field.

Define a multi-key index example

The following example demonstrates how to define a multi-key index using the Report Wizard and the graphical indexer. The sample report to be archived is an AFP invoice. The following pieces of information should be used as indexes:

- Customer Number
- Invoice Number
- Invoice Date
- Item Number (this will be the multi-key index)
- Total Due

As a general rule, you should define triggers and fields from top left to bottom right of the report. This has the added benefit of making your indexer parameters easier to understand.

Figure 2 shows a page from a sample report with a multi-key index.

OnDemand - [04/16/03]

File Edit View Search Notes Options Window Help

Super Sun Seeds
A Growth Company

400 CPU Parkway
Vegetation, NJ 55090

Office: 555-499-2367
Fax: 555-415-9794

THE LAST LEAF
340 DESPERADO COURT
LONGVIEW
CA 12345-6789

CAMBIUM LAYER LIMITED
2222 SApLING CIRCLE
BARKERSVILLE
BC 47365-7290

-- Sold To --

-- Ship To --

Customer Number: 154 Invoice Number: 31354 Invoice Date: 3/16/03 Payment Date: 4/16/03

SHIP Via: TREE TRUCK Shipped Date: 3/16/03 Terms: NET 15 Salesman: MARY PINETREE

| | | | | | |
|-----|----|----------|---------------------------|------|--------|
| 4 | PK | 03698741 | STRING GRAPEFRUIT | 2.01 | 8.04 |
| 300 | EA | 11000146 | AZALIA, GIANT ROSE SEEDS | .55 | 165.00 |
| 6 | CT | 11005010 | EARLY DWARF DANISH SEEDS | 3.01 | 18.06 |
| 24 | DZ | 11005013 | MINCOR NANTES CARROT SEED | .87 | 20.88 |
| 24 | PK | 11005020 | FRENCH PICKLING SEEDS | 2.39 | 57.36 |
| 12 | CT | 12382910 | SUCCATASH SEEDS | .38 | 4.56 |
| 55 | CT | 13145340 | SOUR GRAPE SEEDS | .15 | 8.25 |
| 14 | BZ | 32165478 | BLACK EYED BANANA | 3.01 | 42.14 |
| 600 | DZ | 44646510 | PLUMP RED PLUMS | .49 | 294.00 |
| 40 | DZ | 45613712 | CRANAPPLE BERRY SEEDS | 1.28 | 51.20 |

Ready bmarshall - iSeries (7.1.1.0) Page 1 of 8 100% 24

Figure 2. Multi-key index sample report

To begin, first start the Content Manager OnDemand administrative client and log on to your instance's server. Next, click the Report Wizard toolbar button. Then select the data type; for the example, select AFP. Then select the sample input file. The graphical indexer should now display the spooled file.

The sample report contains AFP data, and only the text is displayed by the graphical indexer, not the AFP resources (such as special fonts, bar codes, graphics, and overlays).

Define the first trigger. Select the / (forward slash) character in the ship date as Trigger1. This trigger will be used to locate the Customer Number, Invoice Number, and Ship Date.

Define the second trigger. Select the . (period) character in the price as Trigger2. This trigger must be defined as a float trigger and will be used to locate the Item Numbers.

Define the third trigger. Select the / (forward slash) character in the payment due date as Trigger3. This trigger will be used to locate the Total Due.

After the triggers are defined, define the fields and indexes. When using the Report Wizard, the fields and indexes are defined in one step. If using the graphical indexer from within an application definition rather than the Report Wizard, the fields and indexes are defined in separate steps.

The first field and index are for the customer number. The customer number is located by using Trigger1. On the Database Field Attributes page, the customer number field is defined as a string data type.

The second field and index are for the invoice number. The invoice number is located by using Trigger1. On the Database Field Attributes page, invoice number is defined as a string data type.

The third field and index are for the invoice date. The invoice date is located by using Trigger1. On the Database Field Attributes page, invoice date is defined as a date data type, and selected as our segment field.

The fourth field and index are for the item number. The item number is located by using Trigger2. On the Database Field Attributes page, item number is defined as a string data type.

The *Mask* parameter is used to specify a pattern that the field data must match in order to be used as an index. In the example, a field must consist of eight numeric characters (each # represents one numeric character). This could be useful if the trigger (a period) could be present in row that did not contain an item number.

After defining all of the fields, you must go back and mark the item number index as multi-key (as described below).

The fifth field and index are for the total due. The total due is located by using Trigger3. On the Database Field Attributes page, total due is defined as a string data type.

That completes defining the fields and the indexes.

Now you must go back and specify the item number, which is Index4, as the multi-key. Click on the Toggle select Trigger, Index, Field Parameters toolbar button.

The administrative client opens the Select dialog box.

Click on Index 4. Then click on the Properties button to open the Update an Index dialog box.

Click on the Allow Multiple Values check box.

Click on the OK button to save the item number index as a multi-key index.

Close the Select dialog box.

To verify how the system will index the document, click on the Toggle between Display and Add Parameters toolbar button.

The defined triggers will be highlighted in red. The defined fields will be highlighted in blue.

You can now close the graphical indexer window and complete the process of using the Report Wizard to define the application group, application, and folder.

Figure 3 shows the indexer parameters that were generated for the example report.

```

TRIGGER1=*,55,X'61',(TYPE=GROUP) /* / */
TRIGGER2=*,64,X'4B',(TYPE=FLOAT) /* . */
TRIGGER3=*,31,X'61',(TYPE=FLOAT) /* / */
FIELD1=0,15,6,(TRIGGER=1,BASE=0)
FIELD2=0,33,6,(TRIGGER=1,BASE=0)
FIELD3=0,50,8,(TRIGGER=1,BASE=0)
FIELD4=0,19,8,(TRIGGER=2,BASE=0,MASK='#####')
FIELD5=0,69,12,(TRIGGER=3,BASE=0)
INDEX1=X'83A4A2A39596',FIELD1,(TYPE=GROUP,BREAK=YES) /* custno */
INDEX2=X'8995A59596',FIELD2,(TYPE=GROUP,BREAK=YES) /* invno */
INDEX3=X'8995A58481A385',FIELD3,(TYPE=GROUP,BREAK=YES) /* invdate */
INDEX4=X'89A3859495A494',FIELD4,(TYPE=GROUP,BREAK=NO,ALLOWMULTIPLEVALUES=YES)/* itemnum */
INDEX5=X'A396A3819384A485',FIELD5,(TYPE=GROUP,BREAK=NO) /* totaldue*/

```

Figure 3. Multi-key index indexer parameters

After loading the example report, you can start the Content Manager OnDemand Client, open the new folder, and search for documents.

Defining transaction fields

A transaction report contains pages of records with one or more columns of sorted data. For example, each page of a general ledger report contains up to 80 transaction records. Each record contains a unique value, such as a transaction number. The records in the report are sorted on the transaction number.

Rather than storing every transaction number in the database (perhaps hundreds of thousands of rows), you can break the report into groups of pages (say, 100 pages in a group), extract the beginning and ending transaction number for each group of pages, and store the values in the database. Then, to retrieve the group of the report that contains a specific transaction number, a user specifies a transaction number. Content Manager OnDemand compares the transaction number with the beginning and ending values stored in the database and retrieves the group that matches the query.

To define a transaction report that contains one or more columns of sorted data as described in the example, a transaction field is used. A transaction field allows Content Manager OnDemand to index a group of pages using the first index value on the first page and the last index value on the last page.

The easiest method of specifying a transaction field is to use the Report Wizard and the graphical indexer.

The indexer parameter for the transaction field will look similar to the following:
FIELD1=*,*,10,(OFFSET=(3:12),MASK='#####',ORDER=BYCOL)

The indexer parameter for the index created from the transaction field will look similar to the following:

```
INDEX1=X'D3968195',FIELD1,(TYPE=GROUPRANGE,BREAK=NO)
```

These indexer parameters would be added by the Report Wizard to the Content Manager OnDemand Application definition. To see them, go to the Indexer Information tab, then click on Keyboard and then Modify to view the Application's Indexer Parameters.

Define a transaction report example

The following example demonstrates how to define a transaction report using the Report Wizard and the graphical indexer. The sample report that we are archiving is a Loan Delinquency Report. Each page of the loan delinquency report contains loan records. Each record contains a unique value, the loan number. The records in the report are sorted on the loan number. We want to use the following pieces of information as indexes:

- Report Date
- Starting Page Number
- Loan Number (this will be the transaction field)

As a general rule you should define triggers and fields from top left to bottom right of the report. This has the added benefit of making your indexer parameters easier to understand.

Figure 4 shows a sample page of the report.

| | | | | |
|--------|-----------|-------------------------|------|----------|
| REPORT | D33313001 | ONDEMAND NATIONAL BANK | DATE | 01-15-00 |
| BANK | 001 | | TIME | 16:03:46 |
| FROM | 01/01/99 | | MODE | 9 |
| TO | 12/31/99 | LOAN DELINQUENCY REPORT | PAGE | 0001 |

| LOAN NUMBER | CUSTOMER NAME | LOAN AMOUNT | DELINQUENT 30 DAYS | DELINQUENT 60 DAYS | DELINQUENT 90 DAYS |
|-------------|------------------|---------------|--------------------|--------------------|--------------------|
| 0100000000 | AARON, ROBERT | \$10000000.00 | \$ 50.00 | \$ 50.00 | \$.00 |
| 0100000001 | ABBOTT, DAVID | \$ 11000.00 | \$ 100.00 | \$ 200.00 | \$.00 |
| 0100000002 | ABBOTT, DAVID | \$ 12000.00 | \$ 140.00 | \$.00 | \$.00 |
| 0100000003 | ABBOTT, DAVID | \$ 13000.00 | \$ 150.00 | \$.00 | \$.00 |
| 0100000005 | ROBINS, STEVEN | \$ 500.00 | \$ 50.00 | \$.00 | \$.00 |
| 0100000006 | ARNOLD, SAMUEL | \$ 1000.00 | \$ 75.00 | \$ 150.00 | \$ 225.00 |
| 0100000007 | PETERS, PAUL | \$ 650.00 | \$ 50.00 | \$.00 | \$.00 |
| 0100000008 | ROBERTS, ABRAHAM | \$ 9000.00 | \$ 120.00 | \$.00 | \$.00 |
| 0100000009 | SMITH, RANDOLPH | \$ 8000.00 | \$ 115.00 | \$.00 | \$.00 |
| 0100000010 | KLINE, PETER | \$ 8500.00 | \$ 110.00 | \$.00 | \$.00 |

Figure 4. Transaction Field sample report

To begin, first start the Content Manager OnDemand administrative client and log on to your instance's server. Next, click the Report Wizard toolbar button. Then select the data type; for the example, select SCS. Then select the sample input file. The graphical indexer should now display the spooled file.

Define the first trigger. Select the word REPORT for Trigger1. This trigger will be used to determine the start of the document, and to locate the Report Date and Starting Page Number fields.

Trigger1 is the only trigger required. Next, define fields and indexes. When using the report wizard, the fields and indexes are defined in one step. If using the graphical indexer within the application definition rather than the Report Wizard, the fields and indexes are defined in separate steps.

The first field and index are for the report date. The report date is located by using Trigger1. On the Database Field Attributes page, the report date is defined as a date data type and is selected as the segment field.

The second field and index are for the starting page number. The starting page number is located by using Trigger1. On the Database Field Attributes page, the starting page number is defined as an integer data type.

After defining all of the fields, you must change the starting page number field so that a new document group is not created each time the page number changes.

The second field and index are for the loan number. The loan number is located by using a mask. The Mask parameter is used to specify a pattern that the transaction field data must match in order to be used as an index. In the example, the field must consist of ten numeric characters (each # represents a numeric character). A transaction field does not use a trigger to locate the data, it uses the mask to define how the data must be structured, and uses any data on that page that matches that mask.

The Database Field Attributes page has specific parameters to support a transaction field. The end user of the sample report will see the folder field names. The database field names are using internally to Content Manager OnDemand and are not seen by end users.

The end user will enter search criteria (the loan number) into the field that is identified by the Query Folder Field. The document list will show two loan numbers. These are the starting and ending loan numbers of the group of the report that contains the loan number that was searched for.

The loan number is defined as a string data type.

Now you must go back and specify that the starting page number, which is Index2, should not start a new document group when the value changes. Click the Toggle select Trigger, Index, Field Parameters toolbar button.

The administrative client opens the Select dialog box.

Click on Index 2. Then click on the Properties button to open the Update an Index dialog box.

Under Break, select the No option. Click the OK button to save the starting page number index as a Break=No index. A change in the starting page number will no longer cause a new document group to be created.

Close the Select dialog box.

To verify how the system will index the document, click on the Toggle between Display and Add Parameters toolbar button.

The defined triggers will be highlighted in red. The defined fields will be highlighted in blue. The defined transactions fields will be highlighted in green.

You can now close the graphical indexer window and complete the process of using the Report Wizard to define the application group, application, and folder.

Figure 5 shows the indexer parameters that were generated for the example report.

```
TRIGGER1=*,2,X'D9C5D7D6D9E3', (TYPE=GROUP) /* REPORT */
FIELD1=0,83,8, (TRIGGER=1,BASE=0)
FIELD2=3,87,4, (TRIGGER=1,BASE=0)
FIELD3=*,*,10, (OFFSET=(3:12),MASK='#####',ORDER=BYROW)
INDEX1=X'998481A385',FIELD1, (TYPE=GROUP,BREAK=YES) /* rdate */
INDEX2=X'A297818785',FIELD2, (TYPE=GROUP,BREAK=NO) /* spage */
INDEX3=X'D396819540D5A494828599',FIELD3, (TYPE=GROUPRANGE,BREAK=NO) /* Loan Number */
```

Figure 5. Transaction Field indexer parameters

After archiving the example report, you can start the Content Manager OnDemand client, open the new folder, and search for documents.

Assigning default index values

You can create a Content Manager OnDemand application definition with an index field that does not always exist on the print page. If a value is not found for that field during indexing (in other words, if only blanks are found or the field location does not exist on the particular print page), then the DEFAULT keyword is used to determine the default value to use. The DEFAULT keyword can be placed on the FIELD indexer parameter line of the indexer parameters for a particular application definition.

The DEFAULT keyword can be specified in one of two ways. The first method allows you to specify an actual value (given in alphanumeric or hex format). The second method allows you to use the default value that you have specified on the Load Information tab of the Content Manager OnDemand application definition and index propagation (described below).

Examples of the first method:

```
DEFAULT='your_Value' (such as DEFAULT='ABC')
```

or

```
DEFAULT=x'your_Hex' (such as DEFAULT=x'C1C2C3')
```

Examples of the second method:

```
DEFAULT='_*USELOADDEFAULTORPROPAGATION'
```

or

```
DEFAULT=x'6D5CE4E2C5D3D6C1C4C4C5C6C1E4D3E3D6D9D7D9D6D7C1C7C1E3C9D6D55C6D'
```

(In this second case, the hex value specified is the hexadecimal representation of the character string `*USELOADDEFAULTORPROPAGATION*`.)

The second method (using `*USELOADDEFAULTORPROPAGATION*` or its hexadecimal representation) allows the load process to assign the default value from the Load Information tab of the application definition or for propagation to occur. To have the load process assign a default from the Load Information tab, you must specify one by using the Content Manager OnDemand Administrator Client. If you have not specified a default, propagation occurs. Propagation is the process of carrying a value over from its previously found value. This can be

useful but can also have unintended results. For example, if the field was a customer number, the value for customer number is carried from the previous document if one was not found for the current document. This might not be what you intend to happen. Exercise caution when using this second method, as propagation can occur.

Defining text search fields

The text search function is used to search for documents that contain a specified word or phrase that is not already defined as an index field for the documents. Initially, the specified index field values are used for the document search. Then, any document that matches the index fields criteria is searched for the specified text search word or phrase. For example, if the other index fields are date and account number, only documents that match the specified date and account number are searched for the specified text search word or phrase. Then, if a document contains the specified word or phrase, the document is added to the document list.

1. You can define only one text search field per folder.
2. The only valid search operator for a text search field is EQUAL.
3. Wildcards and pattern matching are not supported in a text search field.
4. The case of the specified word or phrase is ignored. For example, the phrase *customer xyz* matches *customer xyz*, *Customer Xyz*, and *CUSTOMER XYZ*.

The text search function is performed entirely on the IBM i server. Any performance impact will depend on the size and number of documents that are searched and on the performance of the system under the pre-existing workload. To limit the number of documents that are searched, users should specify criteria for some or all of the other index fields.

To create a text search field folder definition:

1. Create the application group, application, and folder by using the Report Wizard. (The Report Wizard does not include a provision for creating a text search field. However, doing so can be accomplished in just a few steps outside the Report Wizard.)
2. Copy the folder.
3. On the Field Definition tab, add a field named Full Text Search and select Text Search for the field type. Click the Add button to add the field.
4. Click OK to update the folder.

After archiving some documents into the application group, you can try the text search function.

You may want to set a number of options within the Content Manager Content Manager OnDemand Windows client to enhance the use of text search:

- From the Options menu, select the Show Search String option. This option causes the text search string that you enter to be highlighted within the document after it is opened.
- If the Autoview option is set to either First Document or Single Document, the document automatically displays with the text search string highlighted. Single Document will cause the document to automatically display if only one document meets the search criteria. First Document always causes the first document in the document list to automatically display, not matter how many documents meet the search criteria.

When you are ready to try your text search field, open the folder that contains the text search field and perform a text search. The text search string can be one or more words. Open one of the documents from the document list. The text search string should be highlighted in the document. You can use the Find Next toolbar button to find the next occurrence of the string in the document. Note that you can still perform standard searches with the folder; you do not have to specify a text search every time that you search for documents.

To use the text search function with AFP or SCS-Extended documents, you must have the Portable Application Solutions Environment (PASE; a product option of IBM i) installed. If PASE is not installed, you will receive message 161 in the Content Manager OnDemand system log when attempting to perform a text search on AFP or SCS-Extended documents. To use the text search function with SCS or Line documents, you do not need PASE.

Handling SCS spooled files that have AFP overlays

The preferred method of handling SCS spooled files that have an AFP overlay named in their associated printer file is to simply change the DEVTYPE parameter of the printer file used to create the original spooled file to *AFPDS. This will cause IBM i to put the data into spool as *AFPDS, which is the most efficient way for Content Manager OnDemand to capture (load) this type of spooled data. However, making this change will require the original, production spooled file to be printed on an AFPDS printer. In most cases, if you really are printing it with an overlay, then this should not be a problem. However, if you are printing it on a line printer with preprinted forms, this approach will not work.

If, for some reason you cannot change the original printer file's DEVTYPE parameter to *AFPDS, Content Manager OnDemand can do the conversion to AFP automatically, allowing the spooled file to be viewed and printed with fidelity. (This method is more time-consuming than letting IBM i do it using the DEVTYPE parameter of the printer file.) To enable this conversion, simply specify both the data type and the DOCTYPE indexer parameter in the Content Manager OnDemand Application definition as AFP rather than SCS. When Content Manager OnDemand encounters an *SCS spooled file that has an overlay, and the Application definition and DOCTYPE indexer parameter both specify AFP as the data type, Content Manager OnDemand will convert the *SCS data to *AFPDS and store that newly created *AFPDS spooled file. Reprints out of Content Manager OnDemand will require an AFP-capable printer, but that should be expected due to the overlay. If you specify a data type of AFP in your Content Manager OnDemand Application definitions for any other type of non-AFP spooled file, the loading of the data will fail.

Using a mask when defining applications fields

A mask specifies the pattern of symbols that the indexing program matches with data located for a particular field. With the OS/400 indexer, a mask can be used with either a trigger-based field or a transaction field. If the data matches the mask, then the indexer selects the field. If the data does not match the mask, then the field is treated as if the trigger or transaction field was not found.

You can specify the following symbols in the mask:

- @ Matches alphabetic characters
- # Matches numeric characters

- = Matches any character
- ~ Matches any non-blank character
- ^ Matches any non-blank character
- % Matches the blank character and numeric characters

For example, a mask of #####.## would cause the indexer to select the field only if the data in the field (from left to right) contains four numeric characters, followed by a decimal point, followed by two numeric characters.

An example of the indexer parameter syntax for a field with a mask is as follows:
 FIELD4=0,-24,7,(TRIGGER=3),BASE=TRIGGER,MASK='#####.##')

Note: You may need to manually add the MASK keyword to the correct field definition if you are using a group trigger-based field. Support for group trigger-based field masks may not be available with the graphical indexing tool for the version of the Content Manager OnDemand administrative client that you are using.

Using Tag Logical Elements (TLEs)

Using Tag Logical Elements (TLEs) to identify index data requires no special check boxes or other special setup. The Content Manager OnDemand graphical indexer (which is invoked by the Content Manager OnDemand Administrator Client when defining an application) automatically displays TLE data at the top of each print page before displaying the data itself, allowing you to use the TLE data just as you use the print data itself to extract index information (such as a customer number or invoice number).

An example of the data you might see in the Content Manager OnDemand Administrator Client's graphical indexer when you are working with TLEs in an AFPDS spooled file is shown below. The four lines near the top, immediately following the *GROUP_START line, represent the TLE information. The AFP datastream *text* must be encoded in EBCDIC and not ASCII. This is also true of TLEs.

```
*GROUP_START          113928
Invoice Number        113928
Invoice Date          06/15/07
Customer Number       44332
Invoice Total         $  2,859.36
```

```
ABC COMPANY
101 Plagioclase Blvd.
Deva Station          VA 55564
```

```
528 555-1234
```

```
SHIP DATE    04/07/73
Dewey Cheatham & Howe
P.O. Box 47899
Ridiculous   TN 79832
```

```
CUSTOMER NUMBER 44332
```

```
PURCHASE ORDER NO. - C3050279
```

```
17 IGUANAS          3.23      0.11      77.34
93 SHOE HORNS       18.95     13.13     127.83
55 RUNCIBLE SPOONS  43.43     9.23     239.01
```

| | | | |
|------------------|-------|-------|----------|
| 55 HATRACKS | 97.00 | 43.83 | 4,721.64 |
| 93 THELMIN WIRES | 0.54 | 2.32 | 14.12 |
| 09 TOOTHPICKS | 53.00 | 19.91 | 102.43 |
| | | | 5282.37 |

Content Manager OnDemand PDF indexer

The Content Manager OnDemand PDF indexer is a program that you can use to extract index data from and generate index data about Adobe PDF input files. The index data can enhance your ability to store, retrieve, and view documents with Content Manager OnDemand. The PDF indexer supports PDF Version 1.3 or later input and output data streams. For more information about the PDF data stream, see the *Portable Document Format Reference Manual*, published by Adobe Systems Incorporated. Adobe also provides online information with the Acrobat Exchange and Acrobat Distiller products, including online guides for Adobe Capture, PDFWriter, Distiller, and Exchange.

You define and store PDF documents on the server using standard Content Manager OnDemand functions. You must define a Content Manager OnDemand application and application group. As part of the application, you must define the indexing parameters used by the PDF indexer to process input files. You can automate the indexing and loading of data by using special parameters of the ADDRPTOND (using *STMF for the INPUT parameter) or STRMONOND (using *DIR for the TYPE parameter) commands or the ARSLOAD API program. See the Command Reference appendix of the *IBM Content Manager Content Manager OnDemand for i: Common Server Administration Guide* for more information on the ADDRPTOND and STRMONOND commands. See the API Reference appendix of the *IBM Content Manager OnDemand for i: Common Server Administration Guide* for more information on the ARSLOAD API program and its parameters.

After you index and store input files in Content Manager OnDemand, you use the Content Manager OnDemand Windows client program to work with the PDF document or documents created during the indexing and loading process. See the *IBM Content Manager Content Manager OnDemand for i: Common Server Planning and Installation Guide* for more information about working with PDF documents with the Content Manager OnDemand Windows client.

Figure 6 on page 18 illustrates the process of indexing and loading PDF input files.

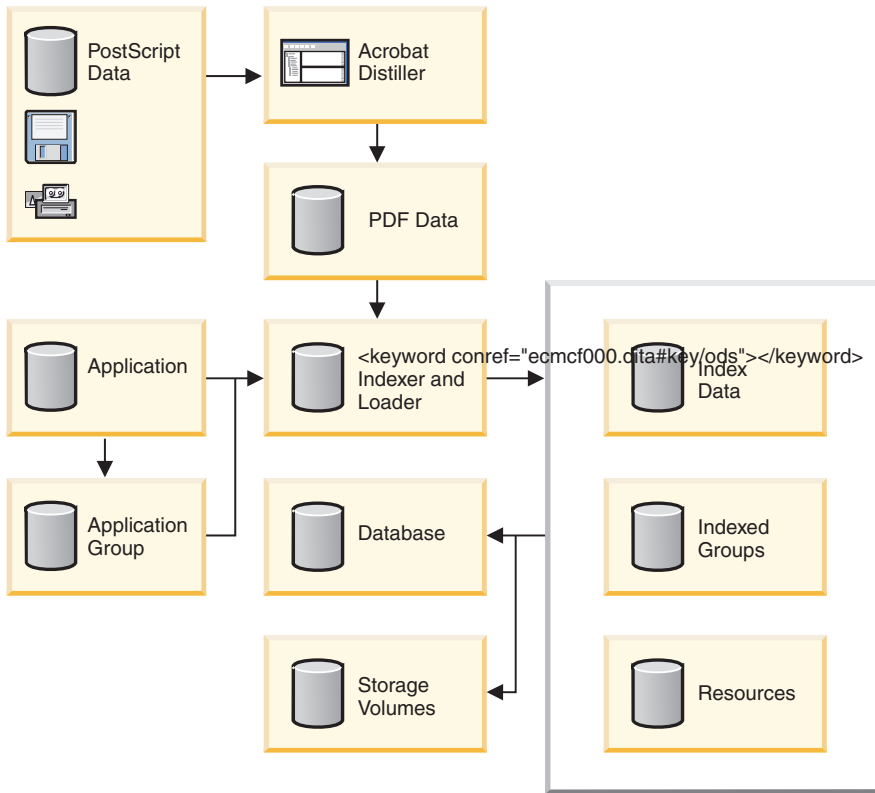


Figure 6. Processing PDF input files in Content Manager OnDemand

The PDF indexer processes PDF input files. A PDF file is a distilled version of a PostScript file, adding structure and efficiency.

Content Manager OnDemand retrieves processing information from application and application group definitions that are stored in the database. The application definition identifies the type of input data, the indexing program used to index the input files, the indexing parameters, and other information about the input data. The application group identifies the database and storage management characteristics of the data. You can use the administrative client to create the application and the indexing parameters.

When Content Manager OnDemand processes a PDF input file and the application Indexing Information page specifies PDF as the indexer, Content Manager OnDemand automatically calls the PDF indexer to process the input file. The PDF indexer processes the PDF input file with indexing parameters that determine the location and attributes of the index data. The PDF indexer extracts index data from the PDF file and generates an index file and an output file. The output file contains groups of indexed pages. A group of indexed pages can represent the entire input file or, more typically, one or more pages from the input file. If the input file contains logical groups of pages, such as statements or policies, the PDF indexer can create an indexed group for each statement or policy in the input file and users can retrieve a specific statement or set of statements, rather than the entire file.

The PDF indexer can optionally extract embedded resources from the PDF input files and store them in a resource file. The resource file is loaded into Content Manager OnDemand at the same time as the output file. After indexing the data,

Content Manager OnDemand stores the index data in the database and the indexed groups and resources on storage volumes.

How OnDemand uses index information

Every item stored in Content Manager OnDemand is indexed with one or more *group-level* indexes. Groups are determined when the value of an index changes (for example, account number). When you load a PDF file into the system, Content Manager OnDemand invokes the PDF indexer to process the indexing parameters and create the index data. Content Manager OnDemand then loads the index data into the database, storing the group-level attribute values that the PDF indexing program extracted from the data into their corresponding database fields. Figure 7 illustrates the index creation and data loading process.

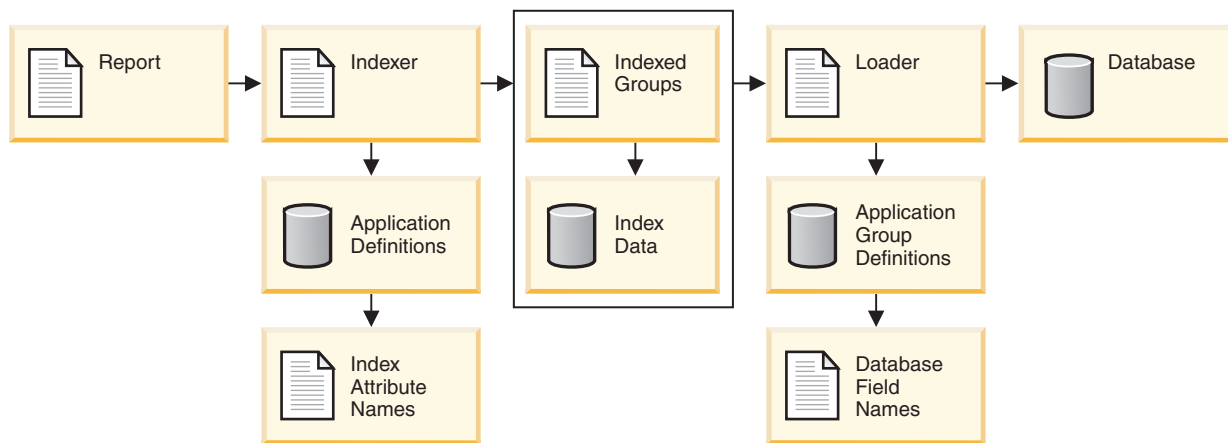


Figure 7. Indexing and loading data

You typically create an application for each report that you plan to store in Content Manager OnDemand. When you create an application, you define the indexing parameters that the indexing program uses to process the report and create the index data that is loaded into the database. For example, an INDEX parameter includes an attribute name and identifies the FIELD parameter that the indexing program uses to locate the attribute value in the input data. When you create an application, you must assign the application to an application group. The attribute name you specify on an INDEX parameter should be the same as the name of the application group database field into which you want Content Manager OnDemand to store the index values.

You define database fields when you create an application group. OnDemand creates a column in the application group table for each database field that you define. When you index a report, you create index data that contains index field names and index values extracted from the report. Content Manager OnDemand stores the index data into the database fields.

To search for reports stored in OnDemand, the user opens a folder. The search fields that appear when the user opens the folder are mapped to database fields in an application group (which, in turn, represent index attribute names). The user constructs a query by entering values in one or more search fields. OnDemand searches the database for items that contain the values (index attribute values) that match the search values entered by the user. Each item contains group-level index

information. OnDemand lists the items that match the query. When the user selects an item for viewing, the OnDemand client program retrieves the selected item from disk or archive storage.

Processing PDF input files with the graphical indexer

This section describes how to use the graphical indexer to create indexing information for PDF input files.

Important: If you plan to use the Report Wizard or the graphical indexer to process PDF input files, then you must first install Adobe Acrobat on the system from which you plan to run the administrative client. You must purchase Adobe Acrobat from Adobe or some other software vendor.

Content Manager OnDemand provides the ARSPDF32.API file to enable PDF viewing from the client. If you install the client after you install Adobe Acrobat, then the installation program will copy the API file to the Acrobat plug-in directory. If you install the client before you install Adobe Acrobat, then you must copy the API file to the Acrobat plug-in directory. Also, if you upgrade to a new version of Acrobat, then you must copy the API file to the new Acrobat plug-in directory. The default location of the API file is `\Program Files\IBM\Content Manager OnDemand32\PDF`. The default Acrobat plug-in directory is `\Program Files\Adobe\Acrobat x.y\Acrobat\Plug_ins`, where `x.y` is the version of Acrobat, for example, 4.0, 5.0, and so forth.

You can define indexing information in a visual environment. You begin by opening a sample input file with the graphical indexer. (**Note:** The input file is limited to a PC file when using the graphical PDF indexer. The graphical PDF indexer is designed to work with workstation PDF files, not PDF spooled files in an output queue on the IBM i server.) You can run the graphical indexer from the report wizard or by choosing the sample data option from the Indexing Information page of the application. After you open an input file in the graphical indexer, you define triggers, fields, and indexes. The PDF indexer uses the triggers, fields, and indexes to locate the beginning of a document in the input data and extract index values from the input data. Once you have defined the triggers, fields, and indexes, you can save them in the application so that Content Manager OnDemand can use them later on to process the input files that you load into the system.

You define a trigger, field, or index by drawing a box around a text string with the mouse and then specifying properties. For example, to define a trigger that identifies the beginning of a document, you could draw a box around the text string `Account Number` on the first page of a statement in the input file. Then, on the Add a Trigger dialog box, you would accept the default values provided, such as the location of the text string on the page. When processing an input file, the PDF indexer attempts to locate the specified string in the specified location. When a match occurs, the PDF indexer knows that it has found the beginning of a document. The fields and indexes are based on the location of the trigger.

The PDF file that you open with the graphical indexer should contain a representative sample of the type of input data that you plan to load into the system. For example, the sample input file must contain at least one document. A good sample should contain several documents so that you can verify the location of the triggers, fields, and indexes on more than one document. The sample input file must contain the information that you need to identify the beginning of a document in the input file. The sample input file should also contain the

information that you need to define the indexes. When you load an input file into the system, the PDF indexer will use the indexing information that you create to locate and extract index values for each document in the input file.

The following example describes how to use the graphical indexer from the Report Wizard to create indexing information for an input file. The indexing information consists of a trigger that uniquely identifies the beginning of a document in the input file and the fields and indexes for each document.

1. To begin, start the administrative client.
2. Log on to a server.
3. Start the report wizard by clicking the Report Wizard icon on the toolbar. The report wizard opens the Sample Data dialog box.
4. Click Select Sample Data to open the Open dialog box. **Note:** The Sample Data is limited to a PC file when using the graphical PDF indexer. The graphical PDF indexer is designed to work with workstation PDF files, not PDF spooled files in an output queue on the IBM i server.
5. Type the name or full path name of a file in the space provided or use the Look in or Browse commands to locate a file.
6. Click Open. The graphical indexer opens the input file in the report window.
7. Press F1 to open the main help topic for the report window. The main help topic contains general information about the report window and contains links to other topics that describe how to add triggers, fields, and indexes. Under Options and Commands, click Indexer Information page to open the Indexing Commands topic. (You can also use the content help tool to display information about the icons on the toolbar.) Under Tasks, Indexer Information page, click Adding a trigger (PDF).
8. Close any open help topics and return to the report window.
9. Define a trigger.
 - Find a text string that uniquely identifies the beginning of a document. For example, Account Number, Invoice Number, Customer Name, and so forth.
 - Using the mouse, draw a box around the text string. Start just outside of the upper left corner of the string. Click and hold mouse button one. Drag the mouse towards the lower right corner of the string. As you drag the mouse, the graphical indexer uses a dotted line to draw a box. When you have enclosed the text string completely inside of a box, release the mouse button. The graphical indexer highlights the text string inside of a box.
 - Click the Define a Trigger icon on the toolbar to open the Add a Trigger dialog box. Verify the attributes of the trigger. For example, the text string that you selected in the report window should be displayed under Value; for Trigger1, the Pages to Search should be set to Every Page. Click Help for assistance with the other options and values that you can specify.
 - Click OK to define the trigger.
 - To verify that the trigger uniquely identifies the beginning of a document, first put the report window in display mode. Then click the Select tool to open the Select dialog box. Under Triggers, double click the trigger. The graphical indexer highlights the text string in the current document. Double click the trigger again. The graphical indexer should highlight the text string on the first page of the next document. Use the Select dialog box to move forward to the first page of each document and return to the first document in the input file.
 - Put the report window in add mode.
10. Define a field and an index.

- Find a text string that can be used to identify the location of the field. The text string should contain a sample index value. For example, if you want to extract account number values from the input file, then find where the account number is printed on the page.
 - Using the mouse, draw a box around the text string. Start just outside of the upper left corner of the string. Click and hold mouse button one. Drag the mouse towards the lower right corner of the string. As you drag the mouse, the graphical indexer uses a dotted line to draw a box. When you have enclosed the text string completely inside of a box, release the mouse button. The graphical indexer highlights the text string inside of a box.
 - Click the Define a Field icon on the toolbar to open the Add a Field dialog box.
 - On the Field Information page, verify the attributes of the index field. For example, the text string that you selected in the report window should be displayed under Reference String; the Trigger should identify the trigger on which the field is based. Click Help for assistance with the options and values that you can specify.
 - On the Database Field Attributes page, verify the attributes of the database field. In the Database Field Name space, enter the name of the application group field into which you want Content Manager OnDemand to store the index value. In the Folder Field Name space, enter the name of the folder field that will appear on the client search screen. Click Help for assistance with the other options and values that you can specify.
 - Click OK to define the field and index.
 - To verify the locations of the fields, first put the report window in display mode. The fields should have a blue box drawn around them. Next, click the Select tool to open the Select dialog box. Under Fields, double-click Field 1. The graphical indexer highlights the text string in the current document. Double click Field 1 again. The graphical indexer should move to the next document and highlight the text string. Use the Select dialog box to move forward to each document and display the field. Then return to the first document in the input file.
 - Put the report window in add mode.
11. Click the Display Indexer Parameters tool to open the Display Indexer Parameters dialog box. The Display Indexer Parameters dialog box lists the indexing parameters that the PDF indexer will use to process the input files that you load into the application. At a minimum, you need one trigger, one field, and one index. See "Parameter reference" on page 31 for details about the indexing parameters.
 12. When you have finished defining all of the triggers, fields, and indexes, close the report window.
 13. Click Yes to save the changes to the indexer parameters.
 14. On the Sample Data window, click Next to continue with the report wizard.

Manually indexing input data

Note: If you prefer creating your own PDF indexing parameters manually rather than using the graphical PDF indexer, you can use the instructions in the remainder of this chapter to do so.

Indexing concepts

Indexing parameters include information that allow the PDF indexer to identify key items in the print data stream, *tag* these items, and create *index elements* pointing to the tagged items. Content Manager OnDemand uses the tag and index data for efficient, structured search and retrieval. You specify the index information that allows the PDF indexer to segment the data stream into individual items, called *groups*. A group is a collection of one or more pages, such as a bank statement, insurance policy, phone bill, or other logical segment of a report. The PDF indexer creates indexes for each group when the value of an index changes (for example, account number).

A tag is made up of an *attribute name*, for example, Customer Name, and an *attribute value*, for example, Earl Hawkins. Tags also include information that tell the PDF indexer where to locate the attribute value on a page. For example, a tag used to collect customer name index values provides the PDF indexer with the starting and ending position on the page where the customer name index values appear. The PDF indexer generates index data and stores it in a generic index file.

Coordinate system

The location of the text strings the PDF indexer uses to determine the beginning of a group and index values are described as *x* and *y* pairs in a coordinate system imposed on the page. For each text string, you identify its upper left and lower right position on the page. The upper left corner and lower right corner form a string box. The string box is the smallest rectangle that completely encloses the text string. The origin is in the upper left hand corner of the page. The *x* coordinate increases to the right and *y* increases down the page. You also identify the page on which the text string appears. For example, the text string Customer Name, that starts 4 inches to the right and 1 inch down and ends 5.5 inches to the right and 1.5 inches down on the first page in the input file can be located as follows:

```
ul(4,1),lr(5.5,1.5),1,'Customer Name'
```

Content Manager OnDemand provides the ARSPDUMP command to help you identify the locations of text strings on the page. See “ARSPDUMP reference” on page 48 for more information about ARSPDUMP.

Indexing parameters

Processing parameters can contain index and conversion parameters, options, and values. For most reports, the PDF indexer requires at least three indexing parameters to generate index data:

- TRIGGER

The PDF Indexer supports the following types of triggers:

- GROUP TRIGGERS

The PDF Indexer compares words in the input file with the text string specified in a trigger. The location of the trigger string value must be identified using the *x,y* coordinate system and page offsets. A maximum of 16 triggers (group or float) can be specified. Group triggers are used in conjunction. For example, all the group triggers must match before the PDF Indexer can begin to locate index information. The group triggers and the fields based on them are used to define the extent of the groups. The indexer must find all the group triggers at least once within the document or it will stop processing and issue an error message.

- FLOAT TRIGGERS

Float triggers are used to locate fields which might occur more than once within a group, or might not occur at all. The PDF Indexer compares words in the input file with the text string specified in a trigger. The location of the trigger string value must be identified using the x,y coordinate system and page offsets. A maximum of 16 triggers (group or float) can be specified. A single float trigger must match before the PDF Indexer can begin to locate index information. The fields based on floating triggers do not define the extent of the groups. If a floating trigger is not found, the indexer continues processing with no error.

The following rules apply when using floating triggers:

1. Trigger1 must be a group trigger.
2. Fields based on floating trigger must contain a default value.
3. Fields based on floating triggers cannot be combined with any other field in an index.
4. At least one index must contain a field (or fields) based on a group trigger.

- FIELD

The field parameter specifies the location of the data that the PDF indexer uses to create index values.

- Field definitions are based on TRIGGER1 by default, but can be based on any of 16 TRIGGER parameters.
- The location of the field must be identified using the x,y coordinate system and page offsets.
- A maximum of 128 fields can be defined.
- A field parameter can also specify all or part of the actual index value stored in the database.

- INDEX

The index parameter is where you specify the attribute name and identify the field or fields on which the index is based. We strongly encourage you to name the attribute the same as the application group database field name.

- The PDF indexer creates indexes for a group of one or more pages.
- You can concatenate field parameters to form an index, unless any of the fields was based on a floating trigger. Fields based on floating triggers cannot be combined with any other field in an index.
- A maximum of 128 index parameters can be specified.

The PDF Indexer creates a new group and extracts new index values when one or more of the index values change, unless the index contains a field based on a floating trigger. Fields based on floating triggers cannot be used to create a new group.

Figure 8 on page 25 depicts a portion of a page from a sample input file. The text strings that determine the beginning of a group and the index values are enclosed in rectangles.

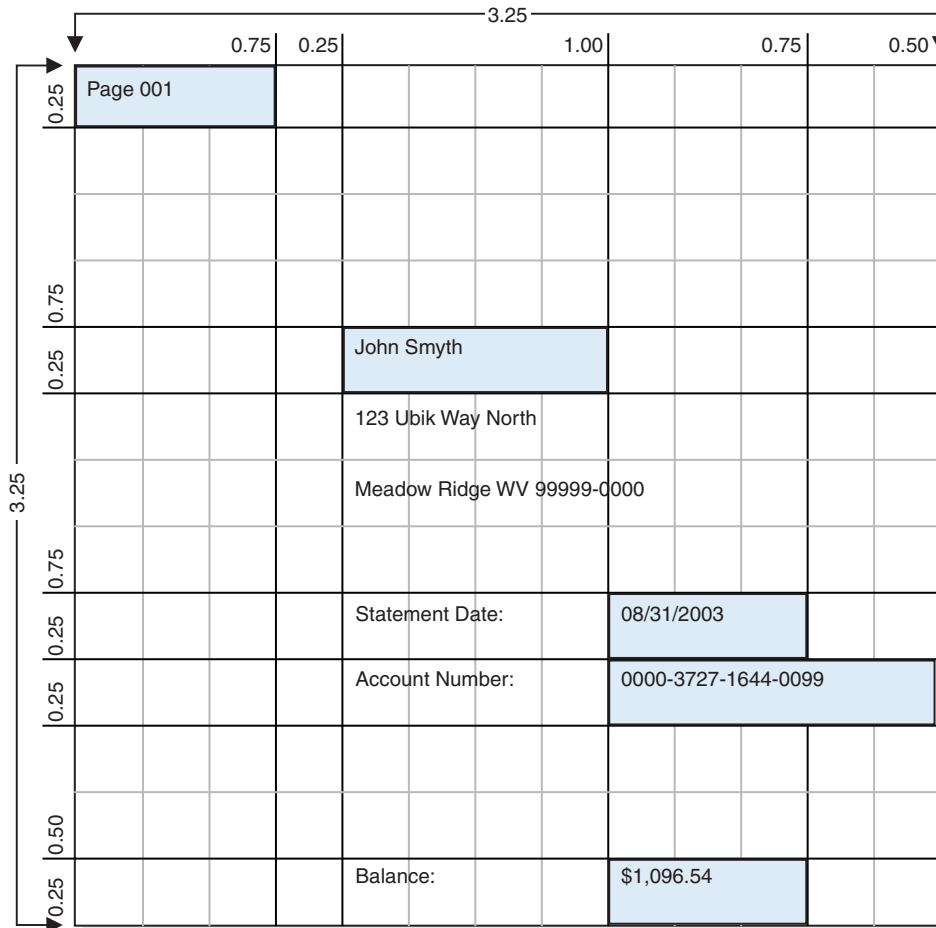


Figure 8. Indexing data with the PDF indexer

TRIGGER parameters tell the PDF indexer how to identify the beginning of a group in the input. The PDF indexer requires one TRIGGER parameter to identify the beginning of a group (statement) in the sample file. FIELD parameters determine the location of index values in a statement. A field is collected when the trigger (or triggers) associated with the field is found. INDEX parameters identify the attribute names of the index fields. Indexes are based on one or more field parameters. The following parameters could be used to index the report depicted in Figure 8. See “Parameter reference” on page 31 for details about the parameter syntax.

- Define a trigger to search each page in the input data for the text string that identifies the start of a group (statement):


```
TRIGGER1=u1(0,0),lr(.75,.25),*,'Page 001'
```
- Define fields to identify the location of index data. For the sample report, we might define four fields:
 - FIELD1 identifies the location of customer name index values.


```
FIELD1=u1(1,1),lr(2,1.25),0
```
 - FIELD2 identifies the location of statement date index values.


```
FIELD2=u1(2,2),lr(2.75,2.25),0
```
 - FIELD3 identifies the location of account number index values.


```
FIELD3=u1(2,2.25),lr(3.25,2.5),0
```
 - FIELD4 identifies the location of the balance index values.

FIELD4=u1(2,3),1r(2.75,3.25),0

- Define indexes to identify the attribute name for an index value and the field parameter used to locate the index value.
 - INDEX1 identifies the customer name, for values extracted using FIELD1.
INDEX1='cust_name',FIELD1
 - INDEX2 identifies the statement date, for values extracted using FIELD2.
INDEX2='sdate',FIELD2
 - INDEX3 identifies the account number, for values extracted using FIELD3.
INDEX3='acct_num',FIELD3
 - INDEX4 identifies the balance, for values extracted using FIELD4.
INDEX4='balance',FIELD4

Indexing with Metadata Indexes

An Adobe PDF document can contain metadata, which is general information such as title and author that applies to the entire document.

You typically create the document's metadata when the document is created and can modify the metadata at any time.

When INDEXMODE=METADATA is specified, the IBM Content Manager OnDemand PDF Indexer extracts fields from the Document Information Dictionary that correspond to the following metadata keywords, if they exist, and place their values into the .ind file:

- Title
- Author
- Subject
- Keywords
- Creator
- Producer
- CreationDate
- ModDate
- Trapped

The metadata keywords are the group field names within the .ind file and can be mapped to the application group fields in the application. You can opt not to map any group field names. Because the metadata keywords apply to the entire document, you can index the document only as one group. If TRIGGER, FIELD, or INDEX parameters are specified, they are ignored. Metadata indexing cannot be combined with indexing using a TRIGGER. If the document contains none of these metadata fields, the IBM Content Manager OnDemand PDF Indexer issues the following error message and stops processing: ARS4940 Index not found by page *page number*

where *page number* is the number specified in the INDEXSTARTBY parameter.

The PDF indexer converts dates that are specified in the PDF format of (D:YYYYMMDDHHmmSSOHH'mm) to a format of YYYYMMDDHHmmSS. The index values CreationDate and ModDate contain the date formatted with the local time. If the time zone information is specified in the PDF date (the OHH'mm section) the PDF indexer creates another index value named CreationDateTZ or

ModDateTZ which contains the date formatted with the time adjusted to Universal Time. For more information on Adobe date formats, see the Adobe PDF Reference, 5th Edition, ISBN-10: 0321304748.

The only parameter required for metadata indexing is: `indexmode=metadata`

Here is an example of an index file created by Metadata indexing:

```
COMMENT:
COMMENT: Content Manager OnDemand Generic Index File Format
COMMENT:
COMMENT:
COMMENT:Code Page of the Index Data
CODEPAGE:1208
COMMENT:Index Field(s)
GROUP_FIELD_NAME:Title
GROUP_FIELD_VALUE:Content Manager OnDemand: Administrator's Guide
GROUP_FIELD_NAME:Author
GROUP_FIELD_VALUE:IBM
GROUP_FIELD_NAME:Creator
GROUP_FIELD_VALUE:XPP
GROUP_FIELD_NAME:Producer
GROUP_FIELD_VALUE:IBM (ID Workbench)
GROUP_FIELD_NAME:CreationDate
GROUP_FIELD_VALUE:20090408173745
GROUP_FIELD_NAME:CreationDateTZ
GROUP_FIELD_VALUE:20090408233745
GROUP_FIELD_NAME:ModDate
GROUP_FIELD_VALUE:20090408173745
GROUP_FIELD_NAME:ModDateTZ
GROUP_FIELD_VALUE:20090408233745
COMMENT:Index Offsets and Length
GROUP_OFFSET:0
GROUP_LENGTH:748641
GROUP_PAGES:387
GROUP_FILENAME:\pdf\pdfoutput\admin.pdf
COMMENT:
COMMENT:
COMMENT:
COMMENT:End Generic Indexing File
```

How to create indexing parameters

About this task

There are two parts to creating indexing parameters. First, process sample input data to determine the x,y coordinates of the text strings the PDF indexer uses to identify groups and locate index data. Then, create the indexing parameters using the administrative client.

Content Manager OnDemand provides the ARSPDUMP command to help you determine the location of trigger and field string values in the input data. The ARSPDUMP command processes one or more pages of sample report data and generates an output file. The output file contains one record for each text string on a page. Each record contains the x,y coordinates for a box imposed over the text string (upper left, lower right).

The process works as follows:

- Obtain a printed copy of the sample report.
- Identify the string values that you want to use to locate triggers and fields

- Identify the number of the page where each string value appears. The number is the *sheet number*, not the page identifier. The sheet number is the order of the page as it appears in the file, beginning with the number 1 (one), for the first page in the file. A page identifier is user-defined information that identifies each page (for example, iv, 5, and 17-3).
- Process one or more pages of the report with the ARSPDUMP command.
- In the output file, locate the records that contain the string values and make a note of the x,y coordinates.
- Create TRIGGER and FIELD parameters using the x,y coordinates, page number, and string value.

Indexing parameters are part of the Content Manager OnDemand application. The administration client provides an edit window you can use to maintain indexing parameters for the application.

PDF resource collection

The PDF reports that you store in Content Manager OnDemand might contain embedded resources such as fonts and images. When the report is indexed, the report is usually broken up into smaller pieces, and the resources are placed into each new report. Reports contain their own resources, so the size of the indexed reports can become much larger than the original PDF reports.

In order to decrease the size of the indexed reports, the PDF indexer can optionally extract these resources from the PDF reports and place them in a resource file. Content Manager OnDemand loads the resource file at the same time as Content Manager OnDemand loads the indexed report files. When a report is retrieved for viewing or printing, the resources are reinserted into the report, and then the report is sent to the client.

A PDF report might contain no resources if the report uses only the fourteen standard fonts that are listed in the PDF reference. These fonts are guaranteed to be available on the client and are not embedded in the report.

The resources that the PDF indexer collects are based on the value of the RESTYPE parameter. The following table lists values for this parameter. For a complete list, see “RESTYPE” on page 43.

Table 1. Available values for the RESTYPE parameter

| RESTYPE | Function | Description |
|-------------|---------------------------|---|
| NONE | Do not collect resources. | Report does not contain resources, or the resources are small. |
| ALL | Collect fonts and images. | To save space that is used to store the reports. |
| FONT | Collect fonts only. | To save space that is used to store the reports. Report contains fonts only. |
| IMAGE | Collect images only. | To save space that is used to store the reports. Report contains images only. |
| FONT, IMAGE | Collect fonts and images. | To save space that is used to store the reports. |

There is no resource exit for the PDF indexer.

PDF indexing system requirements

Specifying the location of Adobe fonts

The PDF indexer must be able to access fonts to insert appropriate information in a PDF output file. If a font is referenced in an input file but is not available on the system, the PDF indexer might be unable to index the document.

If you installed fonts for use with the PDF indexer, you should verify the location of the fonts. If necessary, you must move your existing fonts to the directory structure specified in Table 2. If you do not move your existing fonts, documents that indexed correctly with earlier versions of the PDF indexer might fail with this version.

The directories specified in Table 2 must be used for any Adobe font sets that you install (including but not limited to DBCS fonts).

Table 2. Location of Adobe font files

| Directory type | Directory | Directory contents |
|------------------------|--|---|
| Font directory | /QIBM/ProdData/Content Manager OnDemand/Adobe/Fonts | contains font resources ending in .PFM, .PFB, .TT, .TTF, .MMM |
| CMap file directory | /QIBM/ProdData/Content Manager OnDemand/Adobe/Resource/CMap | contains font CMap resources; files have no extension |
| CIDFont file directory | /QIBM/ProdData/Content Manager OnDemand/Adobe/Resource/CIDFont | contains font resources ending in .oft |

The ACRO_RES_DIR and PSRESOURCEPATH environment variables used in previous versions of Content Manager OnDemand for i are no longer required.

PDF indexing limitations

If you are using the PDF indexer to generate index data for PostScript and PDF files that are created by user-defined programs, remember:

- The PDF Indexer was tested using documents containing up to 100,000 pages. However, there are many factors that affect the number of pages that can be successfully indexed and stored on your system. Those factors include:
 - the system resources available such as CPU, memory, and disk
 - the size of the PDF input file
 - the type and number of resources such as fonts and images used in the PDF input file

If your PDF file does not store successfully, consider:

- splitting the file into a number of separate, smaller files
- reducing the number of different fonts used
- changing the type of fonts used
- reducing the number or size of the images included in the file

- IBM recommends that the CCSID of the PDF input file be 1252 (WinAnsiEncoding). Using another CCSID might cause unexpected results.
- The PDF indexer supports DBCS languages. However, IBM does not provide any DBCS fonts. You can purchase DBCS fonts from Adobe. The PDF indexer supports all DBCS fonts, except encrypted Japanese fonts.
- Input data delimited with PostScript Passthrough markers cannot be indexed.
- The PDF Indexer does not support documents containing Digital Signatures, or that are password protected.
- The Adobe Toolkit does not validate link destinations or bookmarks to other pages in a document or to other documents. Links or bookmarks may or may not resolve correctly, depending on how you segment your documents.

Tip: We recommend using the PDF Indexer parameter BOOKMARKS=NO so that bookmarks are not copied to the documents created by the indexing process.

- If a font is referenced in an input file but not embedded in the file and the PDF indexer cannot locate the font, the PDF Indexer may be unable to index the document. If the fonts are available on the system, they can be accessed at indexing time if they are referenced in an input file and the location is specified on the FONTLIB parameter.

Input data requirements

The PDF Indexer processes PDF input data. The Content Manager OnDemand directory monitor (started with the Start Monitor for OnDemand (STRMONOND) command with *DIR specified for the Type parameter) and the Add Report to OnDemand (ADDRPTOND) command are the two most common ways to invoke the PDF Indexer to index and load PDF data into Content Manager OnDemand on IBM i. You can also use the ARSLOAD API.

The PDF Indexer generates the index data and then adds the index information to the database and loads the input data on to the storage media defined for the particular Content Manager OnDemand application group to which the data belongs.

If you plan to automate the data indexing and loading process on the Content Manager OnDemand server, either the input file name, specific parameters on the command used to load the data, or a monitor user exit program must identify the application group and application to load. The PDF file name extension is required to initiate a load process. The case (uppercase or lowercase) of the extension (.pdf) is ignored. Application group and application names are case sensitive. Application group and application names may include special characters such as the blank character when using ADDRPTOND or ARSLOAD with a specific application group and application name provided. However, STRMONOND and ARSLOAD when using the MVS naming convention (-A and -G parameters) do not support archiving PDF files that have spaces in the file name. See the *IBM Content Manager OnDemand for i: Common Server Administration Guide* for more information about using the STRMONOND and ADDRPTOND commands and the ARSLOAD API to load data into Content Manager OnDemand.

National language support for indexed PDF documents

Consider the following when using the PDF indexer:

- The PDF indexer supports DBCS languages. However, IBM does not provide any DBCS fonts. You can purchase DBCS fonts from Adobe. The PDF indexer

supports all DBCS fonts, except encrypted Japanese fonts. See “Specifying the location of Adobe fonts” on page 29 if you plan to use DBCS font files.

- When loading data using the PDF indexer, the locale must be set appropriately for the code page of the documents. For example, if the code page of the documents is 954, set the locale environment variable to ja_JP or some other locale that correctly identifies upper and lower case characters in code page 954.
- Data values that you specify on TRIGGER and FIELD parameters must be encoded in the same code page as the document. For example, if the characters in the document are encoded in code page 1252, any data values that you specify on TRIGGER and FIELD parameters must be encoded in code page 1252. Examples of data values that you might specify include TRIGGER string values and FIELD default and constant values.
- Previous versions of the PDF indexer required that you change the output from the graphical indexer from character to EBCDIC hexadecimal to index with a DBCS TRIGGER value. The DBCS TRIGGER parameter will now accept the character data from the graphical indexer. You are no longer required to change the TRIGGER value to EBCDIC hexadecimal. Existing applications that are currently coded using the EBCDIC hexadecimal format will still work without modification.

There still might be situations where a DBCS hexadecimal TRIGGER value is required. If you need to use a hexadecimal DBCS TRIGGER value, you first mark the DBCS TRIGGER value by using the graphical indexer. Once the indexing parameters are created, you must then edit the indexing parameters in the Application definition and replace the marked-up TRIGGER value with its EBCDIC hexadecimal equivalent. The TRIGGER value may be entered with or without spaces between the DBCS characters. Also, the shift out/shift in characters are not longer required between the DBCS characters.

An example of the TRIGGER indexing parameter specified with EBCDIC hexadecimal follows:

```
TRIGGER1=UL(5.40,1.92),LR(6.00,2.07),*,X'0E438D438E439443A60F'
```

For more information about NLS in Content Manager OnDemand, see the *IBM Content Manager Content Manager OnDemand for i: Common Server Planning and Installation Guide*.

Parameter reference

This parameter reference assumes that you will use the Start Monitor for OnDemand (STRMONOND) command, Add Report to OnDemand (ADDRPTOND) command, or ARSLOAD API to process your input files. When you use any one of these three methods to process input files, the PDF indexer ignores any values that you may provide for the INDEXDD, INPUTDD, MSGDD, OUTPUTDD, and PARMDD parameters. If you run the ARSPDOCI API from the command prompt or call it from a user-defined program, then you must provide values for the INPUTDD, OUTPUTDD, and PARMDD parameters and verify that the default values for the INDEXDD and MSGDD parameters are correct.

If you must include spaces in a value for an option of a PDF Indexer parameter, enclose the entire value in quotation marks.

BOOKMARKS

Indicates whether or not to copy the bookmarks from the original document to the new documents. The default is YES, which indicates that all the bookmarks from the original document are copied to each new document created by the PDF

Indexer. Some of these bookmarks might no longer be valid. If the original document contains many bookmarks, the size of the new documents can be reduced by not copying the bookmarks.

Required?

No

Default Value

YES

Syntax

BOOKMARKS=[YES | NO]

Options and values

YES Bookmarks are copied to each new document (default).

NO Bookmarks are not copied to new documents.

COORDINATES

Identifies the metrics used for x,y coordinates in the FIELD and TRIGGER parameters.

Required?

No

Default Value

IN

Syntax

COORDINATES=*metric*

Options and values

The *metric* can be:

- IN
The coordinate metrics are specified in inches (the default).
- CM
The coordinate metrics are specified in centimeters.
- MM
The coordinate metrics are specified in millimeters.

FIELD

Identifies the location of index data and can provide default and constant index values. You must define at least one field. You can define up to 128 fields. You can define two types of fields: a *trigger field*, which is based on the location of a trigger string value and a *constant field*, which provides the actual index value that is stored in the database.

Required?

Yes

Default Value

<none>

Trigger field syntax

FIELD n =ul(x,y),lr(x,y),page[, (TRIGGER= n ,BASE={0 | TRIGGER}, MASK='field_mask',DEFAULT='value')]

Options and values:

- *n*
The field parameter identifier. When adding a field parameter, use the next available number, beginning with 1 (one).
- **ul**(*x,y*)
The coordinates for the upper left corner of the field string box. The field string box is the smallest rectangle that completely encloses the field string value (one or more words on the page). The PDF indexer must find the field string value inside the field string box. The supported range of values is 0 (zero) to 45, page width and length, in inches.
- **lr**(*x,y*)
The coordinates for the lower right corner of the field string box. The field string box is the smallest rectangle that completely encloses the field string value (one or more words on the page). The PDF indexer must find the field string value inside the field string box. The supported range of values is 0 (zero) to 45, page width and length, in inches.
- *page*
The sheet number where the PDF indexer begins searching for the field, relative to a trigger or 0 (zero) for the same page as the trigger. If you specify BASE=0, the *page* value can be -16 to 16. If you specify BASE=TRIGGER, the *page* value must be 0 (zero), which is relative to the sheet number where the trigger string value is located.
- **TRIGGER**=*n*
Identifies the trigger parameter used to locate the field. This is an optional keyword, but the default is TRIGGER1. Replace *n* with the number of a defined TRIGGER parameter.
- **BASE**={0 | TRIGGER}
Determines whether the PDF indexer uses the upper left coordinates of the trigger string box to locate the field. Choose from 0 (zero) or TRIGGER. If BASE=0, the PDF indexer adds zero to the field string box coordinates. If BASE=TRIGGER, the PDF indexer adds the upper left coordinates of the location of the trigger string box to the coordinates provided for the field string box. This is an optional keyword, but the default is BASE=0.

You should use BASE=0 if the field data always starts in a specific area on the page. You should use BASE=TRIGGER if the field is not always located in the same area on every page, but is always located a specific distance from a trigger. This capability is useful when the number of lines on a page varies, causing the location of field values to change. For example, given the following parameters:

```
TRIGGER2=ul(4,4),lr(5,8),1,'Total'  
FIELD2=ul(1,0),lr(2,1),0,(TRIGGER=2,BASE=TRIGGER)
```

The trigger string value can be found in a one by four inch rectangle. The PDF indexer always locates the field in a one inch box, one inch to the right of the location of the trigger string value. If the PDF indexer finds the trigger string value in location ul(4,4),lr(5,5), it attempts to find the field in location ul(5,4),lr(6,5). If the PDF indexer finds the trigger string value in location ul(4,6),lr(5,7), it attempts to find the field in location ul(5,6),lr(6,7).

Note: A field that is based on the location of a trigger (BASE=TRIGGER) can be defined at any location on the page that contains the trigger. Previously, a field that was based on the location of a trigger had to be defined to the right and below the upper left point of the trigger. With this change, the *x* or *y* values can be negative, so long as the resulting absolute field coordinates of the field string

rectangle are still in the range of $0 \leq x \leq 45$ and $0 \leq y \leq 45$. The $ul(x,y)$ and $lr(x,y)$ coordinates of the FIELD parameter are relative offsets from the $ul(x,y)$ coordinates of the trigger. For example, suppose the field string rectangle is located at $ul(1,1)$, $lr(2,2)$ which is an absolute location on the page. If the trigger string rectangle is located at $ul(5,5)$, $lr(7,7)$, then the field coordinates would be $ul(-4,-4)$, $lr(-3,-3)$.

- **MASK='field_mask'**

The pattern of symbols that the PDF indexer matches to data located in the field. When you define a field that includes a mask, an INDEX parameter based on the field cannot reference any other fields. Valid mask symbols can include:

@ Matches alphabetic characters. For example:

```
MASK='@@@@@@@@@@@@@@'
```

Causes the PDF indexer to match a 15-character alphabetic field, such as a name.

Matches numeric characters. For example:

```
MASK='#####'
```

Causes the PDF indexer to match a 10-character numeric field, such as an account number.

^ Matches any non-blank character.

% Matches the blank character and numeric characters.

= Matches any character.

Note: The string that you specify for the mask can contain any character. For example, given the following definitions:

```
TRIGGER2=*,25,'ACCOUNT'
FIELD2=0,38,11,(TRIGGER=2,BASE=0,MASK='0000-###-#')
```

The PDF indexer selects the field only if the data in the field columns contains an eleven-character string comprised of any letter, three zeros, a dash character, any four numbers, a dash character, and any number.

- **DEFAULT='value'**

Defines the default index value when there are no words within the coordinates provided for the field string box. You might specify the default value in hexadecimal. A field that is based on a floating trigger must contain a default value.

For example, assume that an application program generates statements that contain an audit field. The contents of the field can be PASSED or FAILED. However, if a statement has not been audited, the application program does not generate a value. In that case, there are no words within the field string box. To store a default value in the database for unaudited records, define the field as follows:

```
FIELD3=ul(8,1),lr(8.5,1.25),1,(DEFAULT='NOT AUDITED')
```

The PDF indexer assigns the index associated with FIELD3 the value NOT AUDITED, if the field string box is blank.

Examples:

The following field parameter causes the PDF indexer to locate the field at the coordinates provided for the field string box. The field is based on TRIGGER1 and

located on the same page as TRIGGER1. Specify BASE=0 because the field string box always appears in a specific location on the page.

```
TRIGGER1=u1(0,0),lr(.75,.25),*, 'Page 0001'  
FIELD1=u1(1,1),lr(3.25,1.25),0, (TRIGGER=1,BASE=0)
```

Hexadecimal default value:

```
TRIGGER1 = u1(4.5,1.25),lr(5.75,1.5), *, 'ACCOUNT'  
FIELD1 = u1(6.6,1.25),lr(7.2,1.25),0, (default=x'30313233')  
INDEX1 = 'Account',FIELD1, (TYPE=GROUP)
```

Field based on a floating trigger:

```
TRIGGER1=UL(5.75,0.71),LR(7.93,1.06),*, 'Bill Summary'  
TRIGGER2=UL(1.82,7.56),LR(3.40,7.85),*, 'Account Number', (TYPE=FLOAT)  
FIELD1=UL(1.90,7.74),LR(3.24,8.04),0, (TRIGGER=2,BASE=0,DEFAULT='N/A')  
FIELD2=UL(5.79,0.13),LR(8.25,0.34),0, (TRIGGER=1,BASE=0)  
INDEX1='acctnum',FIELD1, (TYPE=GROUP)  
INDEX2='name',FIELD2, (TYPE=GROUP)
```

Constant field syntax

FIELD*n*='constant'

Options and values:

- *n*

The field parameter identifier. When adding a field parameter, use the next available number, beginning with 1 (one).

- 'constant'

The literal (constant) string value of the field. This is the index value stored in the database. The constant value can be 1 to 250 bytes in length. The PDF indexer does not validate the type or content of the constant. You might specify the constant value in hexadecimal. See an example in “Examples.”

Examples:

The following field parameter causes the PDF indexer to store the same text string in each INDEX1 value it creates.

```
FIELD1='000000000'  
INDEX1='acct',FIELD1
```

The following field parameters cause the PDF indexer to concatenate a constant value with the index value extracted from the data. The PDF indexer concatenates the constant value specified in the FIELD1 parameter to each index value located using the FIELD2 parameter. The concatenated string value is stored in the database. In this example, the account number field in the data is 14 bytes in length. However, the account number in the database is 19 bytes in length. Use a constant field to concatenate a constant five byte prefix (0000-) to all account numbers extracted from the data.

```
FIELD1='0000-'  
FIELD2=u1(2,2),lr(2.5,2.25),0, (TRIGGER=1,BASE=0)  
INDEX1='acct_num',FIELD1,FIELD2
```

Hexadecimal constant field:

```
FIELD1 = X'4D524830303252'  
FIELD2 = u1(6.6,1.25),lr(7.2,1.25),0, (default=x'30313233')  
INDEX1 = 'Account',FIELD1,FIELD2, (TYPE=GROUP)
```

You can combine a hexadecimal value and a value that is extracted from the document in an index:


```
FIELD1 = X'4D524830303252'  
FIELD2 = ul(6.0,1.4),lr(7.2,1.75),0  
INDEX1 = 'Account',FIELD1,FIELD2,(TYPE=GROUP)
```

Related parameters

INDEX parameter on page “INDEX” on page 37.

TRIGGER parameter on page “TRIGGER” on page 43.

FONTLIB

Identifies the directory or directories in which fonts are stored. Specify any valid path. The PDF indexer searches for fonts in the order that the paths are listed. If a font is referenced in an input file but not embedded in the file, the PDF indexer attempts to locate the font in the directory or directories listed on the FONTLIB parameter.

If the font is not found, the indexing might fail. If the customer purchases additional fonts and installs them on the system, the additional fonts can be found at indexing time if they are referenced in an input file and are present in one of the directories specified on the FONTLIB parameter.

Required?

No

Default Value

/QIBM/ProdData/Content Manager OnDemand/Adobe/fonts

Syntax

FONTLIB=*pathlist*

Options and values

The *pathlist* is a colon-separated string of one or more valid path names. For example:

```
FONTLIB=/QIBM/ProdData/Content Manager OnDemand/fontlib:/QIBM/ProdData/Content Manager OnDemand/V9/f
```

The PDF Indexer searches the paths in the order in which they are specified. Delimit path names with the colon (:) character.

A maximum of 6 paths can be specified.

HEX

Indicates whether the hexadecimal strings in the TRIGGER and FIELD parameters are in UTF-8 or EBCDIC. The default is EBCDIC.

Required?

No

Default Value

NO

- On i, the default is EBCDIC.
- On all other platforms, the default is UTF-8.

Syntax

HEX=*code page indicator*

Options and values

The *code page indicator* can be:

U The hexadecimal strings are specified in UTF-8.

INDEX

Identifies the index name and the field or fields on which the index is based. You must specify at least one index parameter. You can specify up to 128 index parameters. When you create index parameters, IBM recommends that you name the index the same as the application group database field name.

Required?

Yes

Default Value

<none>

Syntax

```
INDEXn='name',FIELDnn[,...FIELDnn]
```

Options and values

- *n*

The index parameter identifier. When adding an index parameter, use the next available number, beginning with 1 (one).

- 'name'

Determines the index name associated with the actual index value. For example, assume INDEX1 is to contain account numbers. The string *acct_num* would be a meaningful index name. The index value of INDEX1 would be an actual account number, for example, 000123456789.

The index name is a string from 1 to 250 bytes in length.

If the name is not the same as an application group database field name, the load fails. The name used in the INDEX parameter must match the Load ID Name value provided on the Load Information tab for the application group database field name.

- **FIELDnn**

The name of the field parameter or parameters that the PDF indexer uses to locate the index. You can specify a maximum of 128 field parameters. When the index value is constructed, the total length of all the concatenated index values cannot exceed 2000 characters.

If the index contains a field which is based on a floating trigger, it must be the only field in the index.

Examples

The following index parameter causes the PDF indexer to create group-level indexes for date index values (the PDF indexer supports group-level indexes only). When the index value changes, the PDF indexer closes the current group and begins a new group.

```
TRIGGER1=UL(5.75,0.71),LR(7.93,1.06),*, 'Bill Summary'  
FIELD1=UL(5.79,0.13),LR(8.25,0.34),0,(TRIGGER=1,BASE=0)  
INDEX1='report_date',FIELD1
```

The following index parameters cause the PDF indexer to create group-level indexes for customer name and account number index values. The PDF indexer closes the current group and begins a new group when either the customer name or the account number index value changes.

```
TRIGGER1=UL(5.75,0.71),LR(7.93,1.06),*, 'Bill Summary'  
FIELD1=UL(5.79,0.13),LR(8.25,0.34),0, (TRIGGER=1,BASE=0)  
FIELD2=UL(1.90,7.74),LR(3.24,8.04),0, (TRIGGER=1,BASE=0)  
INDEX1='name',FIELD1  
INDEX2='acct_num',FIELD2
```

The following index parameters cause the PDF Indexer to create group-level indexes for customer name and balance index values. The PDF Indexer closes the current group and begins a new group only when the customer name index value changes.

```
TRIGGER1=UL(5.75,0.71),LR(7.93,1.06),*, 'Bill Summary'  
TRIGGER2=UL(3.13,3.27),LR(5.59,4.32),*, 'Total Balance', (TYPE=FLOAT)  
FIELD1=UL(5.79,0.13),LR(8.25,0.34),0, (TRIGGER=1,BASE=0)  
FIELD2=UL(1.90,7.74),LR(3.24,8.04),0, (TRIGGER=2,BASE=TRIGGER,DEFAULT='N/A'))  
INDEX1='name',FIELD1  
INDEX2='balance',FIELD2
```

Related parameters

FIELD parameter on page “FIELD” on page 32.

INDEXDD

Determines the name or the full path name of the index object file. The PDF indexer writes indexing information to the index object file. If you specify the file name without a path, the PDF indexer puts the index object file in the current directory. If you do not specify the INDEXDD parameter, the PDF indexer writes indexing information to the file INDEX.

Required?

No

When you process input files with the Start Monitor for OnDemand (STRMONOND) command, Add Report to OnDemand (ADDRPTOND) command, or ARSLOAD API, the PDF indexer ignores any value that you might supply for the INDEXDD parameter. If you process input files with the ARSPDOCI API, verify the value of the INDEXDD parameter.

Default Value

INDEX

Syntax

INDEXDD=*filename*

Options and values

The *filename* is a valid filename or full path name.

INDEXMODE

Determines whether the PDF Indexer uses metadata indexes instead of triggers, fields, and indexes. If not specified, the PDF Indexer uses the trigger, field, and index parameters to perform the indexing.

Required?

No

Default value

<none>

If INDEXMODE is specified along with TRIGGER, FIELD, or INDEX parameters, they are ignored.

Syntax

`INDEXMODE=mode`

Options and values

The *mode* can be: METADATA - Use metadata indexes

Example

The following parameters cause the IBM Content Manager OnDemand PDF Indexer to extract metadata indexes and create a resource file. No other parameters are required.

```
RESTYPE=ALL
INDEXMODE=METADATA
```

INDEXSTARTBY

Determines the page number by which the PDF indexer must locate the first group (document) within the input file. The first group is identified when all of the triggers and fields are found. Fields based on floating triggers are ignored when determining the INDEXSTARTBY page. For example, with the following parameters:

```
TRIGGER1=u1(4.72,1.28),lr(5.36,1.45),*, 'ACCOUNT'
TRIGGER2=u1(6.11,1.43),lr(6.79,1.59),1, 'SUMMARY'
INDEX1='Account',FIELD1,FIELD2
FIELD1=u1(6.11,1.29).lr(6.63,1.45),2
FIELD2=u1(6.69,1.29),lr(7.04,1.45),2
INDEX2='Total',FIELD3
FIELD3=u1(6.11,1.43),lr(6.79,1.59),2
INDEXSTARTBY=3
```

The word ACCOUNT must be found on a page in the location described by TRIGGER1. The word SUMMARY must be found on a page following the page on which ACCOUNT was found, in the location specified by TRIGGER2. In addition, there must be one or more words found for fields FIELD1, FIELD2, and FIELD3 in the locations specified by FIELD1, FIELD2, and FIELD3 which are located on a page that is two pages after the page on which TRIGGER1 was found.

In the example, the first group in the file must start on either page one, page two, or page three. If TRIGGER1 is found on page one, then TRIGGER2 must be found on page two and FIELD1, FIELD2, and FIELD3 must be found on page three.

The PDF indexer stops processing if it does not locate the first group by the specified page number. This parameter is optional, but the default is that the PDF indexer must locate the first group on the first page of the input file. This parameter is helpful if the input file contains header pages. For example, if the input file contains two header pages, you can specify a page number one greater than the number of header pages (INDEXSTARTBY=3) so that the PDF indexer will stop processing only if it does not locate the first group by the third page in the input data.

Note: When you use INDEXSTARTBY to skip header pages, the PDF indexer does not copy non-indexed pages to the output file or store them in Content Manager OnDemand. For example, if you specify INDEXSTARTBY=3 and the first group is found on page three, then pages one and two are not copied to the output file or

stored in Content Manager OnDemand. If you specify INDEXSTARTBY=3 and the first group is found on page two, then page one is not copied to the output file or stored in Content Manager OnDemand.

However, if a field based on a floating trigger is collected from a page which occurs before the first group, the PDF Indexer includes the page where the field was found, and the pages between where the field was found and the start of the first group as part of the first group. The field will be part of the first group in the Search Results.

Required?

No

Default Value

1

Syntax

INDEXSTARTBY=*value*

Options and values

The *value* is the page number by which the PDF indexer must locate the first group (document) in the input file.

INPUTDD

Identifies the name or the full path name of the PDF input file that the PDF indexer will process.

Required?

No

When you process input files with the Start Monitor for OnDemand (STRMONOND) command, Add Report to OnDemand (ADDRPTOND) command, or ARSLOAD API, the PDF indexer ignores any value that you might supply for the INPUTDD parameter. If you process input files with the ARSPDOCI API, you must specify a value for the INPUTDD parameter.

Default Value

<none>

Syntax

INPUTDD=*name*

Options and values

The *name* is the file name or full path name of the input file. If you specify the file name without a path, the PDF indexer searches the current directory for the specified file.

MSGDD

Determines the name or the full path name of the file where the PDF indexer writes error messages. If you do not specify the MSGDD parameter, the PDF indexer writes messages to the display (interactive) or the job log (batch).

Required?

No

When you process input files with the Start Monitor for OnDemand (STRMONOND) command, Add Report to OnDemand (ADDRPTOND) command, or ARSLOAD API, the PDF indexer ignores any value that you

might supply for the MSGDD parameter. If you process input files with the ARSPDOCI API, verify the value of the MSGDD parameter.

Default Value

the display (interactive) or the job log (batch), which are sometimes referred to as stderr (standard error)

Syntax

MSGDD=*name*

Options and values

The *name* is the file name or full path name where the PDF indexer writes error messages. If you specify the file name without a path, the PDF indexer places the error file in the current directory.

OUTPUTDD

Identifies the name or the full path name of the output file.

Required?

No

When you process input files with the Start Monitor for OnDemand (STRMONOND) command, Add Report to OnDemand (ADDRPTOND) command, or ARSLOAD API, the PDF indexer ignores any value that you might supply for the OUTPUTDD parameter. If you process input files with the ARSPDOCI API, you must specify a value for the OUTPUTDD parameter.

Default Value

<none>

Syntax

OUTPUTDD=*name*

Options and values

The *name* is the file name or full path name of the output file. If you specify the file name without a path, the PDF indexer puts the output file in the current directory.

PARMDD

Identifies the name or the full path name of the file that contains the indexing parameters used to process the input data.

Required?

No

When you process input files with the Start Monitor for OnDemand (STRMONOND) command, Add Report to OnDemand (ADDRPTOND) command, or ARSLOAD API, the PDF indexer ignores any value that you might supply for the PARMDD parameter. If you process input files with the ARSPDOCI API, you must specify a value for the PARMDD parameter.

Default Value

<none>

Syntax

PARMDD=*name*

Options and values

The *name* is the file name or full path name of the file that contains the indexing parameters. If you specify the file name without a path, the PDF indexer searches for the file in the current directory.

REMOVERES

Indicates whether or not to remove unused resources before the indexer collects resources and creates the indexes. The input file is examined and a new copy is saved in the Content Manager OnDemand temporary directory. This new copy is then used for processing and the original input file is not changed. You can change the location of the temporary directory by specifying the PDF parameter TEMPDIR. Ensure that the temporary directory has enough space to hold the file. If a file contains many unused resources, you can greatly reduce the size of the resource file and speed up the indexing process by using this parameter. If a file does not contain any unused resources, then do not specify this parameter. You can use this parameter without resource collection.

Required?

No

Default Value

NO

Syntax

REMOVERES=*value*

Options and values

The *value* indicates whether or not to remove unused resources before the indexer collects resources and creates the indexes. The *value* can be one of the following:

- YES The unused resources are removed before the indexer collects resources (if requested) and creates the indexes.
- NO The unused resources are not removed before the indexer collects resources (if requested) and creates the indexes.

RESOBJDD

Specifies the name or the full path name of the resource object file. The PDF indexer collects resources to the resource object file. If you specify the file name without a path, the PDF indexer puts the resource object file in the current directory. Use the RESOBJDD parameter in conjunction with the RESTYPE parameter for the PDF indexer to collect resources.

Required?

No

When you process input files with the Start Monitor for OnDemand (STRMONOND) command, Add Report to OnDemand (ADDRPTOND) command, or ARSLOAD API, the PDF indexer ignores any value that you might supply for the RESOBJDD parameter. If you process input files with the ARSPDOCI API, you must specify a value for the RESOBJDD parameter.

Default Value

<none>

Syntax

RESOBJDD=*filename*

Options and values

The *filename* is a valid file name or full path name. If the PDF file does not contain resources, no RESOBJDD file is produced.

RESTYPE

Determines the types of PDF print resources that the PDF indexer should collect and include in the resource group file.

Required?

No

Default Value

NONE

Syntax

RESTYPE={ NONE | ALL | [FONT] [,IMAGE] }

Options and values

NONE

No resource file is created.

ALL All fonts and images are collected in the resource file.

FONT Fonts are collected in the resource file.

IMAGE

Images are collected in the resource file.

TEMPDIR

Determines the name of the directory that the PDF indexer uses for temporary work space.

Required?

No

Default Value

/arstmp

Syntax

TEMPDIR=*directory*

Options and values

The *directory* is a valid directory name.

TRACEDD

For more information, see “Trace facility” on page 49.

TRIGGER

Identifies locations and string values required to uniquely identify the beginning of a group and the locations and string values of fields used to define indexes. You must define at least one trigger, and can define up to 16 triggers.

Required?

Yes

Default Value

<none>

Syntax

TRIGGER n =ul(x,y),lr(x,y),page,'value'[, (TYPE = {GROUP | FLOAT})]

Options and values

- n

The trigger parameter identifier. When adding a trigger parameter, use the next available number, beginning with 1 (one).
- ul(x,y)

The coordinates for the upper left corner of the trigger string box. The trigger string box is the smallest rectangle that completely encloses the trigger string value (one or more words on the page). The PDF indexer must find the trigger string value inside the trigger string box. The supported range of values is 0 (zero) to 45, page width and length, in inches.
- lr(x,y)

The coordinates for the lower right corner of the trigger string box. The trigger string box is the smallest rectangle that completely encloses the trigger string value (one or more words on the page). The PDF indexer must find the trigger string value inside the trigger string box. The supported range of values are 0 (zero) to 45, page width and length, in inches.
- page

The page number in the input file on which the trigger string value must be located.

 - For TRIGGER1, the *page* value must be an asterisk (*), to specify that the trigger string value can be located on any page in the input file. The PDF indexer begins searching on the first page in the input file. The PDF indexer continues searching until the trigger string value is located, the INDEXSTARTBY value is reached, or the last page of the input file is searched, whichever occurs first. If the PDF indexer reaches the INDEXSTARTBY value or the last page and the trigger string value is not found, then an error occurs and indexing stops.
 - For all other triggers, the *page* value can be 0 (zero) to 16, relative to TRIGGER1. For example, the page value 0 (zero) means that the trigger is located on the same page as TRIGGER1; the value 1 (one) means that the trigger is located on the page after the page that contains TRIGGER1; and so forth. For TRIGGER2 through TRIGGER16, the trigger string value can be a maximum of 16 pages from TRIGGER1.
- 'value'

The actual string value the PDF indexer uses to match the input data. The string value is case sensitive. The value is one or more words that can be found on a page. If the trigger is represented by a double byte or Unicode font in the document, enter the trigger string in hexadecimal. You can use hexadecimal and non-hexadecimal triggers together. See “Examples” on page 45 for a hexadecimal trigger.
- TYPE

The default trigger type is GROUP.

 - TRIGGER1 must be a group trigger. Valid trigger types are the following:
 - GROUP Triggers that identify the beginning of a group. Group triggers, and the fields based on them, define the extent of a group.
 - FLOAT Triggers that identify field data that might not occur in the same location on each page, the same page in each group, or in each group. The PDF Indexer searches within the trigger string box for every occurrence of the trigger. When it is found, any fields based on it will be collected. If a

field is not found, the default value defined for the field will be used.
Fields based on floating triggers cannot define the extent of a group.

Remember:

1. Trigger1 must be a group trigger.
 2. Fields based on floating trigger must contain a default value.
 3. Fields based on floating triggers cannot be combined with any other field in an index.
 4. At least one index must contain a field (or fields) based on a group trigger
- For Float Triggers, the page value must be an asterisk (*), to specify that the trigger string value can be located on any page in the input file. The PDF Indexer begins searching on the first page in the input file. If the trigger is not found, this situation is not considered an error.

Examples

TRIGGER1:

The following TRIGGER1 parameter causes the PDF indexer to search the specified location on every page of the input data for the specified string. You must define TRIGGER1 and the page value for TRIGGER1 must be an asterisk.

```
TRIGGER1=ul(0,0),lr(.75,.25),*, 'Page 0001'
```

Group triggers:

The following trigger parameter causes the PDF indexer to attempt to match the string value Account Number within the coordinates provided for the trigger string box. Specifying a page number of zero (0) for TRIGGER2 means that it can be found on the same page as TRIGGER1.

```
TRIGGER2=ul(1,2.25),lr(2,2.5),0, 'Account Number'
```

The following trigger parameter causes the PDF indexer to attempt to match the string value Total within the coordinates provided for the trigger string box. In this example, a one by four inch trigger string box is defined, because the vertical position of the trigger on the page may vary. For example, assume that the page contains account numbers and balances with a total for all of the accounts listed. There can be one or more accounts listed. The location of the total varies, depending on the number of accounts listed. The field parameter is based on the trigger so that the PDF indexer can locate the field regardless of the actual location of the trigger string value. The field is a one inch box that always begins one inch to the right of the trigger. After locating the trigger string value, the PDF indexer adds the upper left coordinates of the trigger string box to the coordinates provided for the field. Specifying a page number of one (1) for TRIGGER2 means that it can be found on the page following TRIGGER1.

```
TRIGGER2=ul(4,4),lr(5,8),1, 'Total'  
FIELD2=ul(1,0),lr(2,1),0, (TRIGGER=2, BASE=TRIGGER)
```

Float trigger

The following trigger parameter causes the PDF Indexer to attempt to match the string value Total Balance within the coordinates provided for the trigger string box. The field is on the same page as the trigger.

```
TRIGGER2=UL(0.57,0.71),LR(0.89,2.40),*, 'Total Balance', (TYPE=FLOAT)  
FIELD2=UL(1.06,1.77),LR(3.29,2.06),0, (TRIGGER=2, BASE=0, DEFAULT='N/A')  
INDEX2='Balance', FIELD2, (TYPE=GROUP)
```

Hexadecimal trigger:

The following example shows how to code a trigger that represents two side-by-side UTF-8 characters in a document. In this example, each UTF-8 character consists of three bytes. Do not code the index name in hexadecimal.

```
TRIGGER1=UL(1.54,5.40),LR(1.79,5.53),*,X'E6AC8AE79B8A'  
FIELD1=UL(2.29,3.86),LR(3.34,4.04),0,(TRIGGER=1,BASE=0)  
INDEX1='emp_name',FIELD1,(TYPE=GROUP)
```

In this example, hexadecimal and non-hexadecimal triggers are used together:

```
TRIGGER1=UL(6.49,1.72),LR(6.89,1.93),*,X'E8BD8920E7A7BB'  
TRIGGER2=UL(7.02,2.34),LR(7.53,2.60),0,'Page 1'
```

Related parameters

The FIELD parameter on page “FIELD” on page 32.

Message reference

The PDF indexer creates a message list at the end of each indexing run. A return code of 0 (zero) means that processing completed without any errors.

The PDF indexer detects a number of error conditions that can be logically grouped into several categories:

- **Informational**

When the PDF indexer processes a file, it issues informational messages that allow the user to determine if the correct processing parameters have been specified. These messages can assist in providing an audit trail.

- **Warning**

The PDF indexer issues a warning message and a return code of 4 (four) when the fidelity of the document may be in question.

- **Error**

The PDF indexer issues an error message and return code of 8 (eight) or 16 (sixteen) and terminates processing the current input file. Most error conditions detected by the PDF indexer fall into this category. The exact method of termination may vary. For certain severe errors, the PDF indexer may fail with a segment fault. This is generally the case when some system service fails. In other cases, the PDF indexer terminates with the appropriate error messages written either to standard error or to a file. When the PDF indexer is invoked by the ARSLOAD program, error messages are automatically written to the system log. If you run the ARSPDOCI command, you can specify the name or the full path name of the file to hold the processing messages by using the **MSGDD** parameter.

- **Internal Error**

The PDF indexer issues an error message and return code of 16 (sixteen) and terminates processing the current input file.

See *IBM DB2® Content Manager Content Manager OnDemand: Messages and Codes*, SC27-1379 for a list of the messages that may be generated by the PDF indexer, along with explanations of the messages and actions that you can take to respond to the messages. The messages that are generated by the PDF indexer are listed in the Common Server section of the messages publication.

ARSPDOCI reference

Purpose

Generate index data for a PDF file.

The ARSPDOCI program uses the identified locations of text strings on a page of a PDF document to produce a text index file as well as a byte offset indexed PDF document. You can use the ARSPDUMP program to list the locations of text strings in a document. See “ARSPDUMP reference” on page 48 for more information.

Syntax

Note: The following syntax should be used only when you run the ARSPDOCI program from the command line or call it from a user-defined program.

```
▶ ARSPDOCI [COORDINATES=metric] FIELD n [= spec]
▶ [FONTLIB=pathList] INDEX n [= spec] [INDEXDD=fileName]
▶ [INDEXSTARTBY=pageNumber] INPUTDD=fileName [MSGDD=fileName]
▶ OUTPUTDD=fileName [PARMDD=fileName] [TEMPDIR=fileSystem]
▶ TRIGGER n [= spec]
```

Description

The ARSPDOCI program can be used to index a PDF file. The ARSLOAD program automatically calls the ARSPDOCI program if the input data type is PDF and the indexer is PDF. If you need to index a PDF file and you do not want to use the ARSLOAD program to process the file, then you can run the ARSPDOCI program from the command line or call it from a program.

The ARSPDOCI program requires two input files: a PDF document and a parameter file.

If the customer purchases additional fonts and installs them on the system, the additional fonts can be embedded at indexing time if they are referenced in an input file and the location is specified on the FONTLIB parameter. See “FONTLIB” on page 36 for more information.

Parameters

Refer to “Parameter reference” on page 31 for details about the parameters that you can specify when you run the ARSPDOCI program from the command line or a user-defined program.

IFS location

`/usr/bin/arspdoci`

The executable program.

ARSPDUMP reference

Purpose

Print the locations of text strings on a page.

The ARSPDUMP program lists the locations of text strings on a page in a PDF file. The output of the ARSPDUMP program contains a list of the text strings on the page and the coordinates for each string. You can use the information that is generated by the ARSPDUMP program to create the parameter file that is used by the ARSPDOCI program to index PDF files. See “ARSPDOCI reference” on page 47 for more information.

Syntax

```
▶▶ ARSPDUMP -f inputFile [-F fontFile] [-h] [-o outputFile]
▶ -p sheetNumber [-t tempDir]
```

Description

The ARSPDUMP program can be used to identify the locations of text strings on a page in a PDF file.

The output of the ARSPDUMP program contains a list of the text strings on the page and the coordinates for each string.

If a font is referenced in a PDF file, but not embedded, then the ARSPDUMP program attempts to find the font using information provided with the **-F** parameter. If the ARSPDUMP program does not find the font, then it uses a substitute Adobe Type 1 font.

Parameters

-f *inputFile*

The file name or full path name of the PDF file to process.

-F *fontDir*

Identifies directories in which fonts are stored. Specify any valid path. Use the colon (:) character to separate path names. The ARSPDUMP program searches the paths in the order in which they are specified. If you do not specify this flag and name a font directory, then the ARSPDUMP program attempts to locate fonts in the /QIBM/ProdData/Content Manager OnDemand/Adobe/fonts directory.

-h Lists the parameters and their descriptions for the ARSPDUMP program.

-o *outputFile*

The file name or full path name of the file into which the ARSPDUMP program writes output messages. If you do not specify this flag and name a file, then the ARSPDUMP program writes output to the display (interactive) or the job log (batch).

-p *sheetNumber*

The number of the page in the PDF file that you want the ARSPDUMP

program to process. This is the page that contains the text strings that you want to use to define triggers and fields. The sheet number is the order of the page as it appears in the file, beginning with the number 1 (one), for the first page in the file. Contrast with page identifier, which is user-defined information that identifies each page (for example, iv, 5, and 17-3).

-t tempDir

Identifies the directory that the ARSPDUMP program uses for temporary work space. Specify any valid directory name. If you do not specify this flag and name a directory, then the ARSPDUMP program uses the /arstmp directory for temporary work space.

Examples

The following example shows how to invoke the ARSPDUMP program within QSHELL to print the strings and locations of text found on page number three of sample.pdf to sample.out:

```
arspdump -f sample.pdf -o sample.out -p 3
```

See the *IBM Content Manager OnDemand for i: Common Server Administration Guide* for more information about running ARSPDUMP using QSHELL.

IFS location

/usr/bin/arspdump

The executable program.

Trace facility

Beginning with Version 5.3, an enhanced tracing capability for the PDF indexer is now available. The tracing capability provides assistance to users attempting to debug problems, such as when the system fails during the indexing and loading of PDF documents.

To trace or debug a problem with the PDF indexer, the following is required:

- The parameter file, which specifies the fields, triggers, indexes and other indexing information
- The PDF input file to process

The parameter file and PDF input file can be processed by running the PDF indexer from the command line. For example:

```
arspdoci parmdd=filen.parms inputdd=filen.pdf outputdd=filen.out indexdd=filen.ind  
tracedd=filen.trace
```

Where:

arspdoci is the name of the command-line version of the PDF indexer program

parmdd= specifies the name of the input file that contains the indexing parameters

inputdd= specifies the name of the PDF input file to process

outputdd= specifies the name of the output file that contains the indexed PDF documents created by the PDF indexer

indexdd= specifies the name of the output file that contains the index information that will be loaded into the database

tracedd= specifies the name of the output file that contains the trace information

Note: See “ARSPDOCI reference” on page 47 for more information about the parameters that may be specified when running the ARSPDOCI program.

After running the PDF indexer with the trace, the output file specified by the tracedd= parameter will contain detailed information about the processing that took place and where the PDF indexer is failing during the process. The trace information will identify whether a trigger was not found, a field was not found, the PDF data is corrupted, there was a problem extracting a PDF page from the document, or even if there is not enough memory or disk space to complete the required operations. Figure 9 on page 51 shows an example of the trace information that may be generated by the PDF indexer.

```

COORDINATES=IN
ARSPDOCI completed code get_keyword <-----
ARSPDOCI completed code get_keyword 003 ----->
TRIGGER1=UL(7.00,0.25),LR(7.70,0.57),*,'Page:'
ARSPDOCI completed code get_keyword <-----
ARSPDOCI completed code get_keyword 003 ----->
ARSPDOCI completed code parse_trigger <-----
ARSPDOCI completed code parse_quoted_parm <-----
ARSPDOCI completed code parse_quoted_parm 001 ----->
ARSPDOCI completed code parse_trigger 001 ----->
FIELD1=UL(7.00,0.48),LR(7.90,0.77),0,(TRIGGER=1,BASE=0)
ARSPDOCI completed code get_keyword <-----
ARSPDOCI completed code get_keyword 003 ----->
ARSPDOCI completed code parse_field <-----
ARSPDOCI completed code parse_subfields <-----
ARSPDOCI completed code get_keyword <-----
ARSPDOCI completed code get_keyword 003 ----->
ARSPDOCI completed code get_keyword <-----
ARSPDOCI completed code get_keyword 003 ----->
ARSPDOCI completed code parse_subfields 001 ----->
ARSPDOCI completed code parse_field 001 ----->
FIELD2=UL(6.11,1.39),LR(7.2,1.57),0,(TRIGGER=1,BASE=0)
ARSPDOCI completed code get_keyword <-----
ARSPDOCI completed code get_keyword 003 ----->
ARSPDOCI completed code parse_field <-----
ARSPDOCI completed code parse_subfields <-----
ARSPDOCI completed code get_keyword <-----

.
.
.

ARSPDOCI completed code get_keyword <-----
ARSPDOCI completed code get_keyword 003 ----->
ARSPDOCI completed code arspparm_final_sanity_check <-----
ARSPDOCI completed code arspparm_final_sanity_check 001 ----->
ARSPDOCI completed code ArspProcessOpt <-----
ARSPDOCI completed code ArspOpenIndex <-----
ARSPDOCI completed code ArspOpenIndex 001 ----->
Adobe PDF Library version -732512488.-1
Editing is : -1
Number of input pages = 130
ARSPDOCI completed code ArspProcessOpt:Calling ArspSearchDocPages()
ARSPDOCI completed code ArspSearchDocPages <-----
ARSPDOCI completed code ArspSearchDocPages: ArspCreateWordFinder()
ARSPDOCI completed code ArspSearchDocPages: PDWordFinderAcquireWordList()
ARSPDOCI completed code ArspSearchDocPages: PDDocAcquirePage()
ARSPDOCI completed code ArspSearchDocPages: ArspSearchPage()
ARSPDOCI completed code ArspSearchDocPages: PDPPageRelease()
ARSPDOCI completed code ArspSearchDocPages: PDWordFinderReleaseWordList()
Trigger(s) not found by page 1
ARSPDOCI completed code ArspSearchDocPages 004 ----->
ARSPDOCI completed code ArspProcessOpt:Calling ArspCloseIndex()
ARSPDOCI completed code ArspCloseIndex <-----
ARSPDOCI completed code ArspCloseIndex 001 ----->
ARSPDOCI completed code ArspProcessOpt:Calling PDDocClose()
ARSPDOCI completed code ArspProcessOpt 002 ----->
ARSPDOCI completed code 1
ARSPDOCI completed code ArspFreeParms ()

```

Figure 9. Trace information for the PDF indexer

Content Manager OnDemand generic indexer

Content Manager OnDemand provides the generic indexer to allow you to specify indexing information for input data that you cannot or do not want to index with the OS/400 Indexer or the PDF Indexer. For example, suppose that you want to load files into the system that were created by using a word processor. The files can be stored in the system in the same format in which they were created. The files can be retrieved from the system and viewed by using the word processor. However, because the documents do not contain PDF, SCS, SCS-extended, AFP, or LINE spooled data, you cannot index them with the other indexers that are provided with the Content Manager OnDemand product. You can specify index information about the files in the format that is used by the Generic indexer, and load the index data and files into the system. Users can then search for and retrieve the files by using the OnDemand client program.

To use the Generic indexer, you must specify all of the index data for each input file or document that you want to store in and retrieve from the system. You specify the index data in a parameter file. The parameter file contains the index fields, index values, and information about the input files or documents that you want to process. The Generic indexer retrieves the index data from the parameter file and generates the index information that is loaded into the database. OnDemand creates one index record for each input file (or document) that you specify in the parameter file. The index record contains the index values that uniquely identify a file or document in OnDemand.

The generic indexer supports group-level indexes. Group indexes are stored in the database and used to search for documents. You must specify one set of group indexes for each file or document that you want to process with the Generic indexer.

Loading data

The Content Manager OnDemand directory monitor started with the Start Monitor for OnDemand (STRMONOND) command with *DIR or *DIR2 specified for the Type parameter and the Add Report to OnDemand (ADDRPTOND) command are the two most common ways to invoke the Generic Indexer on IBM i. You can also use the ARSLOAD API.

The Generic Indexer uses the index data you provide and the input file you specify, both located in the .IND parameter file. During processing, the index information is added to the database and the input data is loaded on to the storage media defined for the particular Content Manager OnDemand application group to which the data belongs.

There are two ways to run the STRMONOND command:

- STRMONOND with TYPE(*DIR) parameter specified. The STRMONOND command runs as a monitor to periodically check a specified directory for input files to process. When running the STRMONOND command with TYPE(*DIR), the Generic indexer parameter file (.IND) is required to initiate a load process. The GROUP_FILENAME: parameter in the .IND file specifies the full path name of the actual input file to be processed.

- STRMONOND with TYPE(*DIR2) parameter specified. The STRMONOND command runs as a monitor to periodically check a specified directory for input files to process. When running the STRMONOND command with TYPE(*DIR2), a dummy file with the file type extension of .ARD is required to initiate a load process. In addition, the Generic indexer parameter file (.IND) must be located in the specified directory. The GROUP_FILENAME: parameter in the .IND file specifies the full path name of the actual input file to be processed. This is similar to running the ARSLOAD program in daemon mode.

There is one way to run the ADDRPTOND command:

- ADDRPTOND. The ADDRPTOND command is run from the command line to process a specific file. When running the ADDRPTOND command, you specify INPUT(*STMF) and provide the name of the .IND file to process in the Stream file (STMF) parameter (omitting the .IND file extension). The ADDRPTOND command adds the .IND file name extension to the name that you specify. For example, if you specify STMF(po3510), where po3510 is the name of the input file, the ADDRPTOND command looks for and processes the po3510.ind Generic indexer parameter file. The GROUP_FILENAME: parameter in the Generic indexer parameter file specifies the full path name of the actual input file to be processed. This is similar to running the ARSLOAD program in manual mode.

When the data is successfully loaded, both STRMONOND and ADDRPTOND can optionally delete the input file that is specified on the GROUP_FILENAME: parameter if the Delete processed file (DLTSPLF) or Delete input (DLTINPUT) parameters are set to *YES. For the input file to be deleted, the input file must be located in the same directory as the file that triggered the loading of the data, and the file extension must be .OUT. The system also deletes the .IND file (the Generic indexer parameter file) and the .ARD file (the dummy file that is used to initiate a load process in some cases) if the DLTSPLF or DLTINPUT parameter is set to *YES.

Example of file names for STRMONOND TYPE(*DIR):

```
po3510.IND
po3510.OUT
```

The ¹ file is the input file that triggers a load process for STRMONOND TYPE(*DIR). The po3510.IND file is the Generic indexer parameter file, and contains a GROUP_FILENAME: parameter that specifies the input po3510.OUT file to process. When the data is successfully loaded, the system deletes both files.

Example of file names for STRMONOND TYPE(*DIR2):

```
po3510.ARD
po3510.ARD.IND
po3510.ARD.OUT
```

The po3510.ARD file is the dummy file that triggers a load process for STRMONOND TYPE(*DIR2). The po3510.ARD.IND file is the Generic indexer parameter file, and contains a GROUP_FILENAME: parameter that specifies the input file to process, which is po3510.ARD.OUT. When the data is successfully loaded, the system deletes all three files.

There are two ways to run the ARSLOAD API:

Daemon mode

The ARSLOAD API runs as a daemon (monitor) to periodically check a

1. po3510.IND

specified directory for input files to process. When the ARSLOAD API is running in daemon mode, a dummy file with the file type extension of .ARD is required to initiate a load process. In addition, the Generic indexer parameter file (.IND) must be located in the specified directory. The GROUP_FILENAME: parameter in the .IND file specifies the full path name of the actual input file to be processed.

Manual mode

The ARSLOAD API is run from the qshell command line to process a specific file. When the ARSLOAD API is running in manual mode, specify only the *name* of the file to process. The ARSLOAD API adds the .IND file name extension to the name that you specify. For example, if you specify `arsload ... po3510`, where `po3510` is the name of the input file, the ARSLOAD API processes the `po3510.ind` Generic indexer parameter file. The GROUP_FILENAME: parameter in the Generic indexer parameter file specifies the full path name of the actual input file to be processed.

When the data is successfully loaded, ARSLOAD deletes the input file that is specified on the GROUP_FILENAME: parameter if the file name extension is .OUT, and for daemon mode processing, the rest of the input file name is the same as the .ARD file name. For the input file to be deleted, the input file must be located in the same directory as the file that triggered the loading of the data, and the file extension must be .OUT. The system also deletes the .IND file (the Generic indexer parameter file) and the .ARD file (the dummy file that is used to initiate a load process when the ARSLOAD program is running in daemon mode).

Example of file names in daemon processing mode:

```
po3510.ARD
po3510.ARD.IND
po3510.ARD.OUT
```

The `po3510.ARD` file is the dummy file that triggers a load process in daemon mode. The `po3510.ARD.IND` file is the Generic indexer parameter file, and contains a GROUP_FILENAME: parameter that specifies the input file to process, which is `po3510.ARD.OUT`. When the data is successfully loaded, the system deletes all three files.

If you plan to automate the data indexing and loading process on the Content Manager OnDemand server, either the input file name, specific parameters on the command used to load the data, or a monitor user exit program must identify the application group and application to load. The .IND file name extension (for STRMONOND *DIR processing) or the .ARD file name extension (for STRMONOND *DIR2 or ARSLOAD daemon processing) is required to initiate a load process. The case (uppercase or lowercase) of the extension (.ARD or .IND) is ignored. Application group and application names are case sensitive. Application group and application names might include special characters such as the blank character when using ADDRPTOND or ARSLOAD with a specific application group and application name provided. However, STRMONOND and ARSLOAD when using the MVS naming convention (-A and -G parameters) do not support archiving files that have spaces in the file name. See the *IBM Content Manager Content Manager OnDemand for i: Common Server Administration Guide* for more information about using the STRMONOND and ADDRPTOND commands and the ARSLOAD API to load data into Content Manager OnDemand.

Specifying the parameter file

The Generic indexer requires one or more input files that you want to load into the system and a parameter file that contains the indexing information for the input files. To use the Generic indexer, you must create a parameter file that contains the indexing information for the input files. This section describes the parameter file that is used by the Generic indexer.

There are three types of statements that you can specify in a parameter file:

- **Comments.** You can place a comment line anywhere in the parameter file.
- **Code page.** You must specify a code page line at the beginning of the parameter file, before you define any groups.
- **Groups.** A group represents a document that you want to index. Each group contains the application group field names and their index values, the location of the document in the input file, the number of bytes (characters) that make up the document, and the name of the input file that contains the document.

Important:

1. The parameter names in the parameter file are case sensitive and must appear in upper case. For example, `GROUP_FIELD_NAME:account` is valid, while `group_field_name:account` is not.
2. When loading data using the Generic indexer, the locale must be set appropriately for the `CODEPAGE:` parameter. For example, if `CODEPAGE:954` is specified, set the locale environment variable to `ja_JP` or some other locale that correctly identifies upper and lower case characters in code page 954.

CODEPAGE:

Specifies the code page of the input data. You must specify one and only one code page. The **CODEPAGE:** line must appear before you specify any of the groups. The **CODEPAGE:** line is required.

Important: When loading data using the Generic indexer, the locale must be set appropriately for the `CODEPAGE:` parameter. For example, if `CODEPAGE:954` is specified, set the locale environment variable to `ja_JP` or some other locale that correctly identifies upper and lower case characters in code page 954.

Syntax

`CODEPAGE:cpgid`

Options and values

The character string **CODEPAGE:** identifies the line as specifying the code page of the input data. The string `cpgid` can be any valid code page, a three to five character identifier of an IBM-registered or user-defined code page.

The **CODEPAGE:** parameter is required.

Example

The following illustrates how to specify a code page of 37 for the input data:

```
CODEPAGE:37
```

COMMENT:

Specifies a comment line. You can place comment lines anywhere in the parameter file.

Syntax

COMMENT: text on a single line

Options and values

The character string **COMMENT:** identifies the line as containing a comment. Everything after the colon character to the end of the line is ignored.

Example

The following are examples of comment lines:

```
COMMENT:  
COMMENT: this is a comment
```

GROUP_FIELD_NAME:

Specifies the name of an application group field. Each group that you specify in the parameter file must contain one **GROUP_FIELD_NAME:** line for each application group field. (The application group is where you store a file or document in Content Manager OnDemand. You specify the name of the application group to the ARSLOAD program.) Content Manager OnDemand supports up to 128 fields per application group. If the field names that you specify are different than the application group field names, then you must map the field names that you specify to the application group field names on the application Load Information page.

Specify a pair of **GROUP_FIELD_NAME:** and **GROUP_FIELD_VALUE:** lines for each application group field. For example, if the application group contains two fields, then each group that you specify in the parameter file must contain two pairs of **GROUP_FIELD_NAME:** and **GROUP_FIELD_VALUE:** lines. The following is an example of a group with two application group fields:

```
GROUP_FIELD_NAME:rdate  
GROUP_FIELD_VALUE:05/31/00  
GROUP_FIELD_NAME:studentID  
GROUP_FIELD_VALUE:0012345678
```

The group lines must appear after the **CODEPAGE:** line.

Syntax

GROUP_FIELD_NAME:applgrpFieldName

Options and values

The character string **GROUP_FIELD_NAME:** identifies the line as containing the name of an application group field. The string applgrpFieldName specifies the name of an application group field. Content Manager OnDemand ignores the case of application group field names.

Example

The following shows examples of application group field names:

```
GROUP_FIELD_NAME:rdate  
GROUP_FIELD_NAME:studentID  
GROUP_FIELD_NAME:account#
```

GROUP_FIELD_VALUE:

Specifies an index value for an application group field. Each group that you specify in the parameter file must contain one **GROUP_FIELD_VALUE:** line for each application group field. (The application group is where you store a file or document in Content Manager OnDemand. You specify the name of the application group to the ARSLOAD program.) Content Manager OnDemand

supports up to 128 fields per application group. The **GROUP_FIELD_VALUE:** line must follow the **GROUP_FIELD_NAME:** line for which you are specifying the index value.

Specify a pair of **GROUP_FIELD_NAME:** and **GROUP_FIELD_VALUE:** lines for each application group field. For example, if the application group contains two fields, then each group that you specify in the parameter file must contain two pairs of **GROUP_FIELD_NAME:** and **GROUP_FIELD_VALUE:** lines. The following is an example of a group with two application group fields:

```
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:05/31/00
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
```

The group lines must appear after the **CODEPAGE:** line.

Syntax

GROUP_FIELD_VALUE:value

Options and values

The character string **GROUP_FIELD_VALUE:** identifies the line as containing an index value for an application group field. The string value specifies the actual index value for the field.

Example

The following shows examples of index values:

```
GROUP_FIELD_VALUE:05/31/00
GROUP_FIELD_VALUE:0012345678
GROUP_FIELD_VALUE:0000-1111-2222-3333
```

GROUP_FILENAME:

The file name or full path name of the input file. If you do not specify a path, then the generic indexer searches the current directory for the specified file; however, you should always specify the full path name of the input file.

Each group that you specify in the parameter file must contain one **GROUP_FILENAME:** line. The **GROUP_FILENAME:** line must follow the **GROUP_FIELD_NAME:** and **GROUP_FIELD_VALUE:** lines that comprise a group. The following is an example of a group:

```
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:05/31/00
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
GROUP_OFFSET:0
GROUP_LENGTH:0
GROUP_FILENAME:/tmp/statements.out
```

If the **GROUP_FILENAME** line does not contain a value (blank), the Generic indexer uses the value of the **GROUP_FILENAME** line from the previous group to process the current group. In the following example, the input data for the second and third groups is retrieved from the input file that is specified for the first group:

```
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:05/31/00
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
GROUP_OFFSET:0
GROUP_LENGTH:8124
GROUP_FILENAME:/tmp/statements.out
```

```

GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:06/30/00
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
GROUP_OFFSET:8124
GROUP_LENGTH:8124
GROUP_FILENAME:
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:07/31/00
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
GROUP_OFFSET:16248
GROUP_LENGTH:8124
GROUP_FILENAME:

```

If the first **GROUP_FILENAME** line in the parameter file is blank, you must specify the name of the input file when you run the ARSLOAD program.

The group lines must appear after the **CODEPAGE:** line.

After successfully loading the data, the system deletes the input file that is specified on the **GROUP_FILENAME:** parameter if the file name extension is .OUT, and for daemon mode processing, the rest of the input file name is the same as the .ARD file name. The system also deletes the .IND file (the Generic indexer parameter file) and the .ARD file (the dummy file that is used to initiate a load process when the ARSLOAD program is running in daemon mode). See "Loading data" on page 53 for more information.

Syntax

GROUP_FILENAME:fileName

Options and values

The character string **GROUP_FILENAME:** identifies the line as containing the input file to process. The string fileName specifies the full path name of the input file. You should always specify the full path name of the input file to process. For example:

```
GROUP_FILENAME:/tmp/ondemand/inputfiles/f1b0a1600.out
```

Example

The following are valid file name lines:

```

GROUP_FILENAME:/tmp/statements
GROUP_FILENAME:D:\ARSTMP\statements
GROUP_FILENAME:/tmp/ondemand/inputfiles/f1b0a1600.out
GROUP_FILENAME:

```

GROUP_LENGTH:

Specifies the number of contiguous bytes that comprise the document to be indexed. Specify 0 (zero) to indicate the entire input file or the remainder of the input file. Each group that you specify in the parameter file must contain one **GROUP_LENGTH:** line. The **GROUP_LENGTH:** line must follow the **GROUP_FIELD_NAME:** and **GROUP_FIELD_VALUE:** lines that comprise a group. For example:

```

GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:05/31/00
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
GROUP_OFFSET:0
GROUP_LENGTH:0

```


The group lines must appear after the **CODEPAGE:** line.

Syntax

GROUP_LENGTH:value

Options and values

The character string **GROUP_LENGTH:** identifies the line as containing the byte count of the data to be indexed. The string value specifies the actual byte count. The default value is 0 (zero), for the entire (or remainder) of the file.

Example

The following illustrates how to specify length values:

```
GROUP_LENGTH:0
GROUP_LENGTH:8124
```

GROUP_OFFSET:

Specifies the starting location (byte offset) into the input file of the data to be indexed. Specify 0 (zero) for the first byte (the beginning) of the file. Each group that you specify in the parameter file must contain one **GROUP_OFFSET:** line. The **GROUP_OFFSET:** line must follow the **GROUP_FIELD_NAME:** and **GROUP_FIELD_VALUE:** lines that comprise a group. For example:

```
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:05/31/00
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
GROUP_OFFSET:0
```

The group lines must appear after the **CODEPAGE:** line.

Syntax

GROUP_OFFSET:value

Options and values

The character string **GROUP_OFFSET:** identifies the line as containing the byte offset (location) of the data to be indexed. The string value specifies the actual byte offset. Specify 0 (zero), to indicate the beginning of the file.

Example

The following illustrates offset values for three documents from the same input file. The documents are 8 KB in length.

```
GROUP_OFFSET:0
GROUP_OFFSET:8124
GROUP_OFFSET:16248
```

Parameter file examples

The following example shows how to specify indexing information for three groups or documents. Each document is indexed using two fields. The input data for each document is contained in a different input file.

```
COMMENT:
COMMENT: Generic Indexer Example 1
COMMENT: Different input file for each document
COMMENT:
COMMENT: Specify code page of the index data
CODEPAGE:37
COMMENT: Document #1
COMMENT: Index field #1
GROUP_FIELD_NAME:rdate
```



```

GROUP_FIELD_VALUE:07/13/99
COMMENT: Index field #2
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
COMMENT: document data starts at beginning of file
GROUP_OFFSET:0
COMMENT: document data goes to end of file
GROUP_LENGTH:0
GROUP_FILENAME:/arstmp/statement7.out
COMMENT: Document #2
COMMENT: Index field #1
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:08/13/99
COMMENT: Index field #2
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
GROUP_OFFSET:0
GROUP_LENGTH:0
GROUP_FILENAME:/arstmp/statement8.out
COMMENT: Document #3
COMMENT: Index field #1
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:09/13/99
COMMENT: Index field #2
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
GROUP_OFFSET:0
GROUP_LENGTH:0
GROUP_FILENAME:/arstmp/statement9.out
COMMENT:
COMMENT: End Generic Indexer Example 1

```

The following example shows how to specify indexing information for three groups (documents). Each document will be indexed using two fields. The input data for all of the documents is contained in the same input file.

```

COMMENT:
COMMENT: Generic Indexer Example 2
COMMENT: One input file contains all documents
COMMENT:
COMMENT: Specify code page of the index data
CODEPAGE:37
COMMENT: Document #1
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:07/13/99
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
COMMENT: first document starts at beginning of file (byte 0)
GROUP_OFFSET:0
COMMENT: document length 8124 bytes
GROUP_LENGTH:8124
GROUP_FILENAME:/arstmp/accounting.student information.loan.out
COMMENT: Document #2
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:08/13/99
GROUP_FIELD_NAME:studentID
GROUP_FIELD_VALUE:0012345678
COMMENT: second document starts at byte 8124
GROUP_OFFSET:8124
COMMENT: document length 8124 bytes
GROUP_LENGTH:8124
COMMENT: use prior GROUP_FILENAME:
GROUP_FILENAME:
COMMENT: Document #3
GROUP_FIELD_NAME:rdate
GROUP_FIELD_VALUE:09/13/99
GROUP_FIELD_NAME:studentID

```

```
GROUP_FIELD_VALUE:0012345678
COMMENT: third document starts at byte 16248
GROUP_OFFSET:16248
COMMENT: document length 8124 bytes
GROUP_LENGTH:8124
COMMENT: use prior GROUP_FILENAME:
GROUP_FILENAME:
COMMENT:
COMMENT: End Generic Indexer Example 2
```

Additional indexing topics

This section presents information on indexing topics not covered elsewhere in this manual, that applies to all indexers (OS/400, PDF, and Generic), unless otherwise specified.

Postprocessor program

If you require a postprocessor program to further process the index data that is extracted from your input data, you can create a custom-written program that Content Manager OnDemand calls to process all the index records immediately before loading them into the database. For the latest instructions and sample programs, go to the Content Manager OnDemand for i Support Web site at <http://www.ibm.com/software/data/ondemand/400/support.html>, and search for "postprocessor."

Index (.ind), output (.out), and resource (.res) files in IFS

You might notice files in IFS on your IBM i server that might look similar to this, for example:

```
/SP_QPRLR133_QPRTJOB_TKRUPA_067503_000003_MYSYSTEM_1040629_083851.ind
/SP_QPRLR133_QPRTJOB_TKRUPA_067503_000003_MYSYSTEM_1040629_083851.out
/SP_QPRLR133_QPRTJOB_TKRUPA_067503_000003_MYSYSTEM_1040629_083851.res
```

These are either a result of running the Add Report to OnDemand (ADDRPTOND) command with the Index Only (IDXONLY) parameter set to *YES, or from a failed archive initiated by the ADDRPTOND command, a Content Manager OnDemand monitor job, or one of the ARSxxx APIs. If a home directory exists for the user profile running the archive job, these files are located in that user's home directory. Otherwise, the files are located in the root directory, and may be a little harder to notice and maintain.

The purpose of these files is to help determine why the archive processing failed. The .ind files contain the index data captured during the processing of the file, and might help to identify the cause of the problem. If you have a large number of these files on your system, you should investigate the cause (unless you know that testing has been done with IDXONLY(*YES) specified as described above).

Delete the files if they are not needed for problem determination or testing.

Recommended order for defining triggers and fields

As a general rule, you should define triggers and fields from the top left to the bottom right of the report. This has the added benefit of making your indexer parameters easier to understand.

Defining indexes for data to be retrieved using Content Manager OnDemand Web Enablement Kit (ODWEK)

About this task

The percent sign (%) and colon (:) characters in index data will cause a failure or unpredictable results when retrieving documents using the Content Manager OnDemand Web Enablement Kit (ODWEK) interface. Care should be taken when defining index fields for use with ODWEK if the data contained in the index fields might contain percent sign (%) or colon (:) characters.

Notices

This information was developed for products and services offered in the U.S.A.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only the IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe on any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not give you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing
IBM Corporation
North Castle Drive
Armonk, NY 10504-1785
U.S.A.

For license inquiries regarding double-byte (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

Intellectual Property Licensing
Legal and Intellectual Property Law
IBM Japan Ltd.
1623-14, Shimotsuruma, Yamato-shi
Kanagawa 242-8502 Japan

The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law:INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some states do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this IBM product and use of those Web sites is at your own risk.

IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

Licensees of this program who wish to have information about it for the purpose of enabling: (i) the exchange of information between independently created programs and other programs (including this one) and (ii) the mutual use of the information which has been exchanged, should contact:

IBM Corporation
Software Interoperability Coordinator
3605 Highway 52 N
Rochester, MN 55901-7829
U.S.A.

Such information may be available, subject to appropriate terms and conditions, including in some cases, payment of a fee.

The licensed program described in this information and all licensed material available for it are provided by IBM under terms of the IBM Customer Agreement, IBM International Program License Agreement, or any equivalent agreement between us.

Any performance data contained herein was determined in a controlled environment. Therefore, results obtained in other operating environments may vary significantly. Some measurements may have been made on development-level systems and there is no guarantee that these measurements will be the same on generally available systems. Furthermore, some measurements may have been estimated through extrapolation. Actual results may vary. Users of this document should verify the applicable data for their specific environment.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental.

COPYRIGHT LICENSE:This information contains sample application programs in source language, which illustrates programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs. The sample programs are provided "AS IS", without warranty of any kind. IBM shall not be liable for any damages arising out of your use of the sample programs.

If you are viewing this information softcopy, the photographs and color illustrations may not appear.

Trademarks

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both. If these and other IBM trademarked terms are marked on their first occurrence in this information with a trademark symbol (® or ™), these symbols indicate U.S. registered or common law trademarks owned by IBM at the time this information was published. Such trademarks may also be registered or common law trademarks in other countries. A current list of IBM trademarks is available on the Web at “Copyright and trademark information” at www.ibm.com/legal/copytrade.shtml

Adobe, the Adobe logo, Acrobat, and the Acrobat logo are trademarks of Adobe Systems Incorporated, which may be registered in certain jurisdictions.

Java™ and all Java-based trademarks and logos are trademarks of Sun Microsystems, Inc. in the United States, other countries, or both.

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both.

Microsoft, Windows, and Windows NT are trademarks of Microsoft Corporation in the United States, other countries, or both.

UNIX is a registered trademark of The Open Group in the United States and other countries.

Other company, product, and service names may be trademarks or service marks of others.

Index

A

- ACRO_RES_DIR environment
 - variable 29
- Adobe font requirements
 - PDF indexer 29
- ARSPDOCI
 - COORDINATES parameter 32
 - error messages 46
 - FIELD parameter 32
 - FONTLIB parameter 36
 - INDEX parameter 37
 - INDEXDD parameter 38
 - INDEXMODE parameter 38
 - INDEXSTARTBY parameter 39
 - INPUTDD parameter 40
 - messages 46
 - MSGDD parameter 40
 - OUTPUTDD parameter 41
 - PARMDD parameter 41
 - reference 31, 47
 - TEMPDIR parameter 43
 - TRACEDD parameter 43
 - TRIGGER parameter 43
- ARSPDUMP program
 - reference 48

B

- bookmarks
 - PDF indexer 29

C

- code page
 - DBCS 30
 - generic indexer 56
 - PDF indexer 30
- CODEPAGE: parameter 56
- commands
 - ARSPDOCI 47
 - ARSPDUMP 48
- COMMENT: parameter 56
- constant field 35
- coordinate system 23
- coordinates
 - on FIELD parameter for PDF indexer 33
 - on TRIGGER parameter for PDF indexer 44
- COORDINATES parameter 32
 - flags and values 32

D

- DBCS
 - PDF indexer 30
- DBCS font files
 - environment variables 29
 - where to install 29

- DBCS fonts 29
- debugging 49
- default index value
 - FIELD parameter option 34
- document
 - generic indexer parameter 58, 59, 60

E

- environment variables
 - ACRO_RES_DIR 29
 - PSRESOURCEPATH 29
- error messages
 - ARSPDOCI program 46
 - PDF indexer 46
- examples
 - generic indexer 60

F

- FIELD parameter 32
 - constant field 35
 - default index value 34
 - flags and values 32
 - mask option 34
 - trigger field 32
- fields
 - constant field 35
 - default index value 34
 - generic indexer parameter 57
 - mask option 34
 - PDF indexer parameter 32
 - trigger field 32
- files
 - PDF indexer 30
- flags and values
 - HEX 36
 - REMOVERES 42
 - RESOBJDD 42
 - RESTYPE 43
- font files
 - where to install 29
- FONTLIB parameter 36
 - flags and values 36
- fonts
 - DBCS 29
 - NLS 29
 - PDF indexer 29, 36

G

- generic indexer
 - application group field names 57
 - code page 56
 - CODEPAGE: parameter 56
 - COMMENT: parameter 56
 - document 58, 59, 60
 - examples 60
 - field names 57
 - field values 57

- generic indexer (*continued*)
 - group indexes, defining 57
 - GROUP_FIELD_NAME:
 - parameter 57
 - GROUP_FIELD_VALUE:
 - parameter 57
 - GROUP_FILENAME: parameter 58
 - GROUP_LENGTH: parameter 59
 - GROUP_OFFSET: parameter 60
 - input file 58, 59, 60
 - national language support (NLS) 56
 - NLS 56
 - parameter file 56, 60
- graphical indexer 1
- group indexes
 - defining 37, 57
 - defining for generic indexer 57
 - GROUP_FIELD_NAME: parameter 57
 - GROUP_FIELD_VALUE: parameter 57
 - GROUP_FILENAME: parameter 58
 - GROUP_LENGTH: parameter 59
 - GROUP_OFFSET: parameter 60

H

- header pages
 - skipping 39
- HEX 36

I

- i Navigator
 - configuration v
 - Management Central v
 - system navigation v
- IFS location 47
- INDEX parameter 37
 - flags and values 37
- INDEXDD parameter 38
 - flags and values 38
- indexer parameters
 - using break=yes versus break=no 4
- indexes
 - generic indexer parameter 57
 - PDF indexer parameter 37
- indexing
 - constant field 35
 - default index value 34
 - field mask 34
 - fields for PDF indexer 32
 - group indexes 37
 - header pages 39
 - indexes 37
 - mask option 34
 - parameters 23
 - skipping header pages 39
 - trigger field 32
 - triggers 43
- INDEXMODE parameter 38
 - flags and values 38

INDEXSTARTBY parameter 39
 flags and values 39
input file
 generic indexer parameter 58, 59, 60
INPUTDD parameter 40
 flags and values 40

L

limitations
 PDF indexer 29
links
 PDF indexer 29

M

mask
 FIELD parameter option 34
messages
 ARSPDOCI program 46
 PDF indexer 46
MSGDD parameter 40
 flags and values 40

N

naming input files
 PDF indexer 30
national language support (NLS) 56
 PDF indexer 29, 30
NLS 56
 PDF indexer 29, 30

O

OUTPUTDD parameter 41
 flags and values 41

P

parameter file
 ARSPDOCI program 31
 generic indexer 60
 PDF indexer 23, 31
parameters
 ARSPDOCI program 31, 47
 ARSPDUMP program 48
 CODEPAGE: 56
 COMMENT: 56
 COORDINATES 32
 FIELD 32
 FONTLIB 36
 generic indexer 56
 GROUP_FIELD_NAME: 57
 GROUP_FIELD_VALUE: 57
 GROUP_FILENAME: 58
 GROUP_LENGTH: 59
 GROUP_OFFSET: 60
 INDEX 37
 INDEXDD 38
 indexing 48
 INDEXMODE 38
 INDEXSTARTBY 39
 INPUTDD 40
 MSGDD 40

parameters (*continued*)
 OUTPUTDD 41
 PARMDD 41
 PDF indexer 23, 31
 TEMPDIR 43
 TRACEDD 43
 TRIGGER 43
PARMDD parameter 41
 flags and values 41
PDF indexer
 Adobe font requirements 29
 Adobe PDF 47
 ARSPDOCI reference 47
 ARSPDUMP reference 48
 bookmarks 29
 code page 30
 concepts 22
 constant field 35
 coordinate system 23
 DBCS 30
 DBCS fonts 29
 default index value 34
 error messages 46
 field mask 34
 fields 32
 file naming conventions 30
 font requirements 29
 fonts 29, 36
 group indexes 37
 indexes 37
 indexing concepts 22
 limitations 29
 links 29
 mask option 34
 messages 46
 naming input files 30
 national language support (NLS) 29,
 30
 NLS 29, 30
 parameter file 23
 parameter reference 31
 printing 29
 resource collection 28
 restrictions 29
 transferring input files to 30
 trigger field 32
 triggers 43
 x, y coordinate system 23
PDF resource collection 28
printing
 PDF indexer 29
PSRESOURCEPATH environment
 variable 29

R

REMOVERES 42
report wizard 1
requirements
 Adobe font requirements 29
 fonts 29
RESOBJDD 42
restrictions
 PDF indexer 29
RESTYPE 43

S

skipping header pages 39
Syntax
 Constant field 31
 COORDINATES 31
 Field 31
 FONTLIB 31
 INDEXDD 31
 INDEXn 31
 INDEXSTARTBY 31
 INPUTDD 31
 MSGDD 31
 OUTPUTDD 31
 PARMDD 31
 TEMPDIR 31
 TRIGGER 31
system requirements
 Adobe font requirements 29
 fonts 29

T

TEMPDIR parameter 43
 flags and values 43
trace facility 49
TRACEDD parameter 43
 flags and values 43
 trace facility 49
trigger field 32
TRIGGER parameter 43
 options and values 43
triggers
 PDF indexer parameter 43
Triggers
 field syntax 31
 Group Triggers 31
 TRIGGER1 31

X

x,y coordinate system 23



Printed in USA

SC19-2793-01

