



Hewlett Packard
Enterprise

KeyView

Software Version: 11.5

Filter SDK Java Programming Guide

Document Release Date: October 2017

Software Release Date: October 2017

Legal notices

Warranty

The only warranties for Hewlett Packard Enterprise Development LP products and services are set forth in the express warranty statements accompanying such products and services. Nothing herein should be construed as constituting an additional warranty. HPE shall not be liable for technical or editorial errors or omissions contained herein.

The information contained herein is subject to change without notice.

Restricted rights legend

Confidential computer software. Valid license from HPE required for possession, use or copying. Consistent with FAR 12.211 and 12.212, Commercial Computer Software, Computer Software Documentation, and Technical Data for Commercial Items are licensed to the U.S. Government under vendor's standard commercial license.

Copyright notice

© Copyright 2016-2017 Hewlett Packard Enterprise Development LP

Trademark notices

Adobe™ is a trademark of Adobe Systems Incorporated.

Microsoft® and Windows® are U.S. registered trademarks of Microsoft Corporation.

UNIX® is a registered trademark of The Open Group.

This product includes an interface of the 'zlib' general purpose compression library, which is Copyright © 1995-2002 Jean-loup Gailly and Mark Adler.

Documentation updates

The title page of this document contains the following identifying information:

- Software Version number, which indicates the software version.
- Document Release Date, which changes each time the document is updated.
- Software Release Date, which indicates the release date of this version of the software.

To check for recent software updates, go to <https://downloads.autonomy.com/productDownloads.jsp>.

To verify that you are using the most recent edition of a document, go to <https://softwaresupport.hpe.com/group/softwaresupport/search-result?doctype=online help>.

This site requires that you register for an HPE Passport and sign in. To register for an HPE Passport ID, go to <https://hpp12.passport.hpe.com/hppcf/login.do>.

You will also receive updated or new editions if you subscribe to the appropriate product support service. Contact your HPE sales representative for details.

Support

Visit the HPE Software Support Online web site at <https://softwaresupport.hpe.com>.

This web site provides contact information and details about the products, services, and support that HPE Software offers.

HPE Software online support provides customer self-solve capabilities. It provides a fast and efficient way to access interactive technical support tools needed to manage your business. As a valued support customer, you can benefit by using the support web site to:

- Search for knowledge documents of interest
- Submit and track support cases and enhancement requests
- Access product documentation
- Manage support contracts
- Look up HPE support contacts
- Review information about available services
- Enter into discussions with other software customers
- Research and register for software training

Most of the support areas require that you register as an HPE Passport user and sign in. Many also require a support contract.

To register for an HPE Passport ID, go to <https://hpp12.passport.hpe.com/hppcf/login.do>.

To find more information about access levels, go to <https://softwaresupport.hpe.com/web/softwaresupport/access-levels>.

To check for recent software updates, go to <https://downloads.autonomy.com/productDownloads.jsp>.

Contents

Part I: Overview of Filter SDK	11
Chapter 1: Introducing Filter SDK	12
Overview	12
Features	12
Platforms, Compilers, and Dependencies	13
Supported Platforms	13
Supported Compilers	13
Software Dependencies	14
Windows Installation	15
UNIX Installation	16
Package Contents	17
License Information	17
Enable Advanced Document Readers	18
Update License Information	18
Directory Structure	19
Chapter 2: Getting Started	21
Architectural Overview	21
Enhance Performance	23
File Caching	23
Filtering	24
Subfile Extraction	24
Use the Java Implementation of the API	24
Input/Output Operations	25
Filter in File or Stream Mode	25
Multithreaded Filtering	26
Before Running Your Application	27
The Filter Process Model	27
Filter API	27
File Extraction API	28
Persist the Child Process	28
In the API	28
In the formats.ini File	28
Run Filter In Process	29
In the API	29
In the formats.ini File	29
Run File Extraction Functions Out of Process	29
Restart the File Extraction Server	29
Out-of-Process Logging	30
Enable Out-of-Process Logging	30
Set the Verbosity Level	30
Enable Windows Minidump	31

Keep Log Files	31
Run File Detection In or Out of Process	31
Specify the Process Type In the formats.ini File	32
Specify the Process Type In the API	32
Part II: Use Filter SDK	33
Chapter 3: Use the File Extraction API	34
Introduction	34
Extract Subfiles	35
Extract Images	36
Recreate a File Hierarchy	36
Create a Root Node	36
Example	37
Extract Mail Metadata	38
Default Metadata Set	38
Extract the Default Metadata Set	38
Microsoft Outlook (MSG) Metadata	39
Extract MSG-Specific Metadata	40
Microsoft Outlook Express (EML) and Mailbox (MBX) Metadata	40
Extract EML- or MBX-Specific Metadata	41
Lotus Notes Database (NSF) Metadata	41
Extract NSF-Specific Metadata	41
Microsoft Personal Folders File (PST) Metadata	41
MAPI Properties	42
Extract PST-Specific Metadata	42
Exclude Metadata from the Extracted Text File	43
Extract Subfiles from Outlook Files	43
Extract Subfiles from Outlook Express Files	43
Extract Subfiles from Mailbox Files	43
Extract Subfiles from Outlook Personal Folders Files	44
Use the Native or MAPI-based Reader	44
Use the Native PST Reader (pstnsr)	45
Use the MAPI Reader (pstsar)	45
System Requirements	46
MAPI Attachment Methods	46
Open Secured PST Files	47
Detect PST Files While the Outlook Client is Running	47
Extract Subfiles from Lotus Domino XML Language Files	47
Extract .DXL Files to HTML	48
Extract Subfiles from Lotus Notes Database Files	48
System Requirements	49
Installation and Configuration	49
Windows	49
Solaris	50
AIX 5.x	50

Linux	50
Open Secured NSF Files	51
Format Note Subfiles	51
Extract Subfiles from PDF Files	51
Improve Performance for PDFs with Many Small Images	51
Extract Embedded OLE Objects	52
Extract Subfiles from ZIP Files	52
Default File Names for Extracted Subfiles	52
Default File Name for Mail Formats	52
Default File Name for Embedded OLE Objects	53
Exclude Japanese Guide Text	54
Chapter 4: Use the Filter API	55
Generate an Error Log	55
Enable or Disable Error Logging	56
Change the Path and File Name of the Log File	56
Report Memory Errors	57
Specify a Memory Guard	57
Report the File Name in Stream Mode	57
Example	58
Specify the Maximum Size of the Log File	58
Extract Metadata	58
Extract Metadata for File Filtering	59
Extract Metadata for Stream Filtering	59
Example	59
Convert Character Sets	61
Determine the Character Set of the Output Text	61
Guidelines for Character Set Conversion	61
Set the Character Set During Filtering	62
Set the Character Set During Subfile Extraction	62
Prevent the Default Conversion of a Character Set	63
Extract Tracked Deleted Text	63
Filter PDF Files	64
Filter PDF Files to a Logical Reading Order	64
Rotated Text	66
Extract Custom Metadata from PDF Files	66
Skip Embedded Fonts	67
Control Hyphenation	68
Filter Spreadsheet Files	69
Filter Worksheet Names	69
Filter Hidden Text in Microsoft Excel Files	69
Specify Date and Time Format on UNIX Systems	69
Filter Very Large Numbers in Spreadsheet Cells to Precision Numbers	70
Extract Microsoft Excel Formulas	70
Filter HTML Files	72
Filter XML Files	72
Configure Element Extraction for XML Documents	73

Configure Headers and Footers	77
Error Messages	77
Tab Delimited Output for Embedded Tables	80
Exclude Japanese Guide Text	81
Chapter 5: Sample Programs	82
Introduction	82
ExtractFilter	83
FilterTest	84
FilterFileToFile	87
FilterStreamToStream	88
FilterFileToStream	88
FilterStreamToFile	89
FilterFileByChunk	90
FilterStreamByChunk	91
Appendixes	93
Appendix A: Supported Formats	94
Supported Formats	94
Archive Formats	96
Binary Format	98
Computer-Aided Design Formats	98
Database Formats	100
Desktop Publishing	100
Display Formats	101
Graphic Formats	101
Mail Formats	104
Multimedia Formats	107
Presentation Formats	108
Spreadsheet Formats	110
Text and Markup Formats	112
Word Processing Formats	113
Supported Formats (Detected)	118
Appendix B: Character Sets	125
Multibyte and Bidirectional Support	125
Coded Character Sets	133
Appendix C: File Formats and Extensions	138
File Format and Extension Table	138
Appendix D: Extract and Format Lotus Notes Subfiles	163
Overview	163
Customize XML Templates	163
Use Demo Templates	164
Use Old Templates	164
Disable XML Templates	164
Template Elements and Attributes	165

Conditional Elements	165
Control Elements	166
Data Elements	167
Date and Time Formats	169
Lotus Notes Date and Time Formats	169
KeyView Date and Time Formats	170
Appendix E: File Format Detection	176
Introduction	176
Extract Format Information	176
Determine Format Support	176
Example formats.ini file entries	177
Refine Detection of Text Files	177
Allow Consecutive NULL Bytes in a Text File	178
Translate Format Information	179
Distinguish Between Formats	179
Determine a Document Reader	180
Category Values in formats.ini	180
Appendix F: List of Required Files for Redistribution	199
Core Files	199
Support Files	200
Document Readers	200
Appendix G: Develop a Custom Reader	207
Introduction	207
How to Write a Custom Reader	208
Naming Conventions	208
Basic Steps	209
Token Buffer	209
Macros	210
Reader Interface	211
Function Flow	211
Example Development of fffFillBuffer()	212
Implementation 1—fpFillBuffer() Function	212
Structure of Implementation 1	213
Problems with Implementation 1	213
Implementation 2—Processing a Large Token Stream	213
Structure of Implementation 2	214
Problems with Implementation 2	215
Boundary Conditions	215
Implementation 3—Interrupting Structured Access Layer Calls	216
Structure of Implementation 3	217
Development Tips	218
Functions	219
xxxsrAutoDet()	219
xxxAllocateContext()	220
xxxFreeContext()	221

xxxInitDoc()	221
xxxFillBuffer()	222
xxxGetSummaryInfo()	223
xxxOpenStream()	224
xxxCloseStream()	224
xxxCharSet()	225
Appendix H: Password Protected Files	226
Supported Password Protected File Types	226
Open Password Protected Container Files	227
Filter Password Protected Files	227
Send documentation feedback	228

Part I: Overview of Filter SDK

This section provides an overview of the Filter SDK and describes how to use the Java implementation of the API.

Chapter 1: Introducing Filter SDK

This section describes the Filter SDK package.

• Overview	12
• Features	12
• Platforms, Compilers, and Dependencies	13
• Windows Installation	15
• UNIX Installation	16
• Package Contents	17
• License Information	17
• Directory Structure	19

Overview

KeyView Filter SDK enables you to incorporate text extraction functionality into your own applications. It extracts text and metadata from a wide variety of file formats on numerous platforms, and can automatically recognize over 300 document types. It supports both file-based and stream-based I/O operations, and provides in-process or out-of-process filtering.

Filter SDK is part of the KeyView suite of products. KeyView provides high-speed text extraction, conversion to web-ready HTML and well-formed XML, and high-fidelity document viewing.

Features

- Document readers are threadsafe. The benefit of a threadsafe technology is that you can successfully extract text from hundreds of documents simultaneously. Documents are not queued for sequential filtering, but are actually filtered at the same time.
- Filter supports popular word processing, spreadsheet, and presentation formats. Body text, endnotes, footnotes, and additional items such as document metadata are all included as part of the filtering process.
- Sample programs are provided to demonstrate the functionality of the APIs.
- You can extract files embedded within files, such as email attachments or embedded OLE objects, by using the File Extraction API.
- You can configure memory management. If using the C API, you can provide your own memory allocator to the document readers.
- Filter allows for redirected input and output. You can provide an input stream that is not restricted to file system access.
- Filter automatically recognizes the file type being filtered and uses the appropriate filter. Your application does not need to rely on file name extensions to determine file types.
- You can filter documents to specific character encodings, such as Unicode or UTF-8.

- You can use Filter SDK in conjunction with other KeyView technologies, such as the Index, Highlight, and Annotate APIs.
- You can write custom document readers for formats not directly supported by KeyView.

Platforms, Compilers, and Dependencies

This section lists the supported platforms, supported compilers, and software dependencies for the KeyView software.

Supported Platforms

- CentOS 7
- FreeBSD 8.1 x86
- IBM AIX L6.1 PowerPC 32-bit and 64-bit
- IBM AIX L7.1 PowerPC 32-bit and 64-bit
- Mac OS X Mountain Lion 10.8 or higher on 32- and 64-bit Apple-Intel architecture
- Microsoft Windows Vista Business Edition x86 and x64. Other editions of Vista have not been tested, but are likely supported.
- Microsoft Windows 2008 Server Enterprise Edition x86 and x64
- Microsoft Windows 2008 Server R2
- Microsoft Windows 7 x86 and x64
- Microsoft Windows 8 x86 and x64
- Oracle Solaris 10 SPARC
- Oracle Solaris 10 x86 and x64
- Red Hat Enterprise Linux 5.0 x86 and x64
- Red Hat Enterprise Linux 6.0 x86 and x64
- SuSE Linux Enterprise Server 10, 10.1, 11, x86 and x64

Supported Compilers

Platform	Architecture	Compiler Name	Compiler Version
Microsoft Windows	x86	cl	Microsoft 32-bit C/C++ Optimizing Compiler Version 16.00.30319.01 for x86
	x64	cl	Microsoft C/C++ Optimizing Compiler Version 16.00.30319.01 for x64
Sun Solaris	x86 64-bit	Sun Studio 12	Sun C 5.9 SunOS_i386 Patch 124868-01 2007/07/12
	SPARC 64-bit	Sun Studio	Sun C 5.8 Patch 121015-06 2007/10/03

Platform	Architecture	Compiler Name	Compiler Version
		11	
Linux	x86	gcc / g++	3.4.3 (Redhat 4), 4.1.0 (SuSE Linux 10)
	x64	gcc / g++	4.1.0 (Redhat 4), 4.1.0 (SuSE Linux 10)
IBM AIX	Power	xlC_r / cc_r	IBM XL C/C++ Enterprise Edition V8.0
Mac OSX	Apple-Intel 32-bit and 64-bit	LLVM	Apple LLVM 5.1 (clang-503.0.40) (based on LLVM 3.4svn)
FreeBSD	BSD x86	gcc / g++	4.2.1 [FreeBSD] 20070719

Supported Compilers for Java and .NET Components

Component	Compiler
Java components	Java 1.5
.NET components	Microsoft Visual J# 2005 Compiler 8.00.50727.42

Software Dependencies

Some KeyView components require specific third-party software:

- Java Runtime Environment (JRE) or Java Software Developer Kit (JDK) version 1.5 is required for Java API and graphics conversion in Export SDK.
- Outlook 2002 client or later versions is required when processing Microsoft Outlook Personal Folders (PST) files using the MAPI-based reader (`pstsr`). The native PST reader (`pstnsr`) does not require an Outlook client.

NOTE:

If you are using 32-bit KeyView, you must install 32-bit Outlook. If you are using 64-bit KeyView, you must install 64-bit Outlook.

If the bit editions do not match, an error message from Microsoft Office Outlook is displayed:

Either there is a no default mail client or the current mail client cannot fulfill the messaging request. Please run Microsoft Outlook and set it as the default mail client.

Additionally, KeyView displays the following return code:

Error 32: KVErrror_PSTAccessFailed.

- Lotus Notes or Lotus Domino is required for Lotus Notes database (NSF) file processing. The minimum requirement is 6.5.1, but version 8.5 is recommended.

- Microsoft .NET Framework SDK version 2.0, Microsoft .NET Framework version 2.0 Redistributable Package is required if you are programming in a .NET environment.
- Microsoft Visual C++ 2013 and Microsoft Visual C++ 2010 Redistributables (Windows only).

Windows Installation

To install the SDK on Windows, use the following procedure.

To install the SDK

1. Run the installation program, `KeyViewProductNameSDK_VersionNumber_OS.exe`, where *ProductName* is the name of the product, *VersionNumber* is the product version number, and *OS* is the operating system.

For example:

`KeyViewFilterSDK_11.5_Windows_X86_64.exe`


The installation wizard opens.

2. Read the instructions and click **Next**.

The License Agreement page opens.

3. Read the agreement. If you agree to the terms, click **I accept the agreement**, and then click **Next**.

The Installation Directory page opens.

4. Select the directory in which to install the SDK. To specify a directory other than the default, click  , and then specify another directory. After choosing where to install the SDK, click **Next**.

The License Key page opens.

5. Type the company name and license key that were provided when you purchased KeyView, and then click **Next**.
 - The company name is case sensitive.
 - The license key is a string that contains 31 characters.

NOTE:

The installation program validates the company name and license key and generates the file `install\OS\bin\kv.lic` (where *install* is your chosen installation folder and *OS* is the name of the operating system platform). The license information is validated when the KeyView API is used. If you do not enter a license key at this step, or if you enter invalid information, the KeyView SDK is installed, but the API does not function. When you obtain a valid license key, you can either re-install the KeyView SDK, or manually update the license key file (`kv.lic`) with the new information. For more information, see [License Information, on page 17](#).

The Pre-Installation Summary dialog box opens.

6. Review the settings, and then click **Next**.

The SDK is installed.

7. Click **Finish**.

UNIX Installation

To install the SDK, use one of the following procedures.

To install the SDK from the graphical interface

- Run the installation program and follow the on-screen instructions.

To install the SDK from the console

1. Run the installation program from the console as follows:

```
./KeyViewFilterSDK_VersionNumber_Platform.exe --mode text
```

where:

VersionNumber is the product version.

Platform is the name of the platform.

2. Read the welcome message and instructions and press `Enter`.
The first page of the license agreement is displayed.
3. Read the license information, pressing `Enter` to continue through the text. After you finish reading the text, and if you accept the agreement, type `y` and press `Enter`.
You are asked to choose an installation folder.
4. Type an absolute path or press `Enter` to accept the default location.
You are asked for license information.
5. At the **Company Name** prompt, type the company name that was provided when you purchased KeyView, and then press `Enter`. The company name is case sensitive.
6. At the **License Key** prompt, type the license key that was provided when you purchased KeyView, and then press `Enter`. The license key is a string that contains 31 characters.

NOTE:

The installation program generates the file *install\OS\bin\kv.lic* (where *install* is your chosen installation folder and *OS* is the name of the operating system platform). The license information is validated when the KeyView API is used. If you do not enter a license key at this step, or if you enter invalid information, the KeyView SDK is installed but the API does not function. When you obtain a valid license key, you can either re-install the KeyView SDK, or manually update the license key file (*kv.lic*) with the new information. For more information, see [License Information, on the next page](#).

The Pre-Installation summary is displayed.

7. If you are satisfied with the information displayed in the summary, press `Enter`.
The SDK is installed.

Package Contents

The Filter SDK installation contains:

- All the libraries and executables necessary for extracting text from a wide variety of formats.
- The include files that define the functions and structures used by the application to establish an interface with Filter:

<code>adapi.h</code>	<code>kvfilter.h</code>
<code>adinfo.h</code>	<code>kvioobj.h</code>
<code>kvcfsr.h</code>	<code>kvtoken.h</code>
<code>kvxtract.h</code>	<code>kvtypes.h</code>
<code>kvfilt.h</code>	<code>kvxtract.h</code>
<code>kvfilt2.h</code>	<code>kwautdef.h</code>

- The Java API implemented in the package `com.verity.api.filter` contained in the file `KeyView.jar`.
- The .NET API implemented in the namespace `Autonomy.API.Filter` in the library `FilterDotNet.dll`.
- The C++ SDK, which can be found in the `cppapi` folder.
- Sample programs that demonstrate File Extraction and Filter functionality using the APIs.
- The files necessary to create a custom document reader, and the source for a sample document reader for UTF-8. See [Develop a Custom Reader, on page 207](#).

License Information

During installation, the installation program validates the organization name and license key that you enter, and generates the `install/OS/bin/kv.lic` file, where `install` is the directory in which you installed KeyView, and `OS` is the operating system. This file is opened and validated when the KeyView API is used.

The `kv.lic` file contains the organization name and the 31-digit license key you specified during installation. The contents of a `kv.lic` file looks similar to the following:

```
Company Name
XXXXXXXX-XXXXXXX-XXXXXXX-XXXXXXX
```

The license key controls whether the following are enabled:

- the full version of the KeyView SDK
- the trial version of the KeyView SDK
- language detection and advanced document readers—The following components are considered advanced features, and are licensed separately:

- Microsoft Outlook Personal Folders (PST) reader (`pstsr` and `pstnsr`)
- Lotus Notes database (NSF) reader (`nsfsr`)
- Mailbox (MBX) reader (`mbxsr`)
- Character set detection library (`kvlangdetect`)

If you change the license key at any time, you must update the licensing information in the `kv.lic` file. See [Update License Information](#).

Enable Advanced Document Readers

To enable advanced readers in one of the KeyView SDKs, you must obtain an appropriate license key from HPE and update the installed license key with the new information as described in [Update License Information](#).

If you are enabling the MBX reader in an existing installation of Filter, in addition to updating the license key, change the parameter `208=eml` to `208=mbx` in the `formats.ini` file.

Update License Information

If you currently have an evaluation version of KeyView and have purchased a full version of the SDK, or you are adding a document reader (for example, the PST reader), you must update the license information that was installed with the original version of the KeyView SDK.

If you installed a full version of KeyView, but did not enter licensing information at the time of installation, you must also update the license information.

To update the information, do one of the following:

- Manually update the license information that is stored in the text file named `kv.lic`.
- Re-install the product and enter the new license information when prompted.

To update the KeyView license information

1. Open the license key file, `kv.lic`, in a text editor. The file is in the `install\OS\bin` directory, where `install` is the directory in which you installed KeyView, and `OS` is the operating system. The file contains the following text:

```
COMPANY NAME  
XXXXXXX-XXXXXXX-XXXXXXX-XXXXXXX
```

2. Replace the text `COMPANY NAME` with the company name that appears at the top of the License Key Sheet provided by HPE. Enter the text exactly as it appears in the document.
3. Replace the characters `XXXXXX-XXXXXXX-XXXXXXX-XXXXXXX` with the appropriate license key from the License Key Sheet provided by HPE. The license key is listed in the **Key** column in the **Standalone Products** table. The key is a string that contains 31 characters, for example, `2TQD22D-2M6FV66-2KPF23S-2GEM5AB`. Enter the characters exactly as they appear in the document, including the dashes, but do not include a leading or trailing space.
4. The finished `kv.lic` file looks similar to the following:

```
Autonomy
```

24QD22D-2M6FV66-2KPF23S-2G8M59B

5. Save the `kv.lic` file.

Directory Structure

The following table describes the directories created during the Filter SDK installation. The variable *install* is the path name of the Filter installation directory (for example, `/usr/autonomy/KeyviewFilterSDK` on UNIX, or `C:\Program Files\Autonomy\KeyviewFilterSDK` on Windows).

The variable *OS* is the operating system for which the SDK is installed. For example, the `bin` directory on a standard 32-bit Windows installation would be located at `C:\Program Files\Autonomy\KeyviewFilterSDK\WINDOWS\bin`.

Installed directory structure

Directory	Description
<i>install</i> \OS\bin	Contains the libraries, the format detection file <code>formats.ini</code> , the license key file <code>kv.lic</code> , and other supporting files.
<i>install</i> \OS\lib	(Solaris installations only) Contains the redistributable <code>libstlport.so.1</code> library, which is required to run KeyView on Solaris platforms.
<i>install</i> \dotnetapi	Contains the source files for the .NET API.
<i>install</i> \dotnetapi\dotnethelp	Contains the help for the .NET API.
<i>install</i> \dotnetapi\sample	Contains the sample programs for the .NET API.
<i>install</i> \cppapi	Contains the source files for the C++ API.
<i>install</i> \cppapi\sample	Contains the sample programs for the C++ API.
<i>install</i> \guide	Contains the KeyView Filter SDK programming guides in PDF and HTML format.
<i>install</i> \include	Contains the header files required for Filter.
<i>install</i> \javaapi\javadoc	Contains the Javadoc for the Java API.
<i>install</i> \javaapi\sample	Contains the source files and sample programs for the Java API.
<i>install</i> \rel_notes	Contains the <i>KeyView Filter SDK Release Notes</i> in PDF format.
<i>install</i> \samples\filter	Contains a sample program demonstrating the Filter interface for the C API.

Installed directory structure, continued

Directory	Description
<i>install</i> \samples\filterca	Contains a C sample program demonstrating extraction of a content access stream.
<i>install</i> \samples\pdfini	Contains the initialization file used to extract custom metadata from PDF documents.
<i>install</i> \samples\tstxtract	Contains a C sample program demonstrating the File Extraction interface.
<i>install</i> \samples\utf8sr	Contains the source for the sample document reader for UTF-8 files. You can use this to create your own custom document readers.
<i>install</i> \samples\utf8sr\bin	Contains the C program <code>filtertest</code> . You can use this program to test your custom document readers. See Develop a Custom Reader, on page 207 .

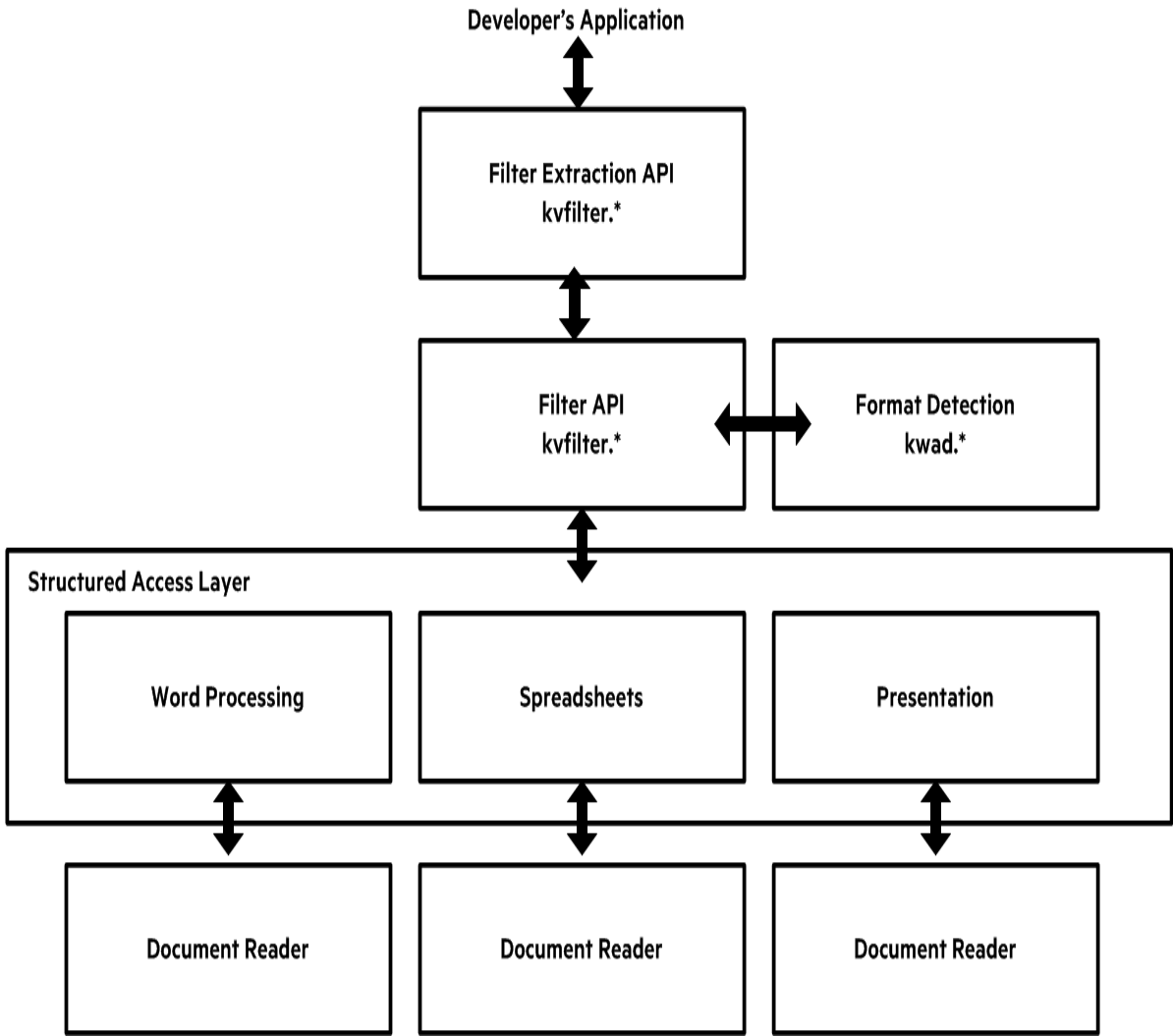
Chapter 2: Getting Started

This section provides an overview of Filter SDK, and describes how to use the Java implementation of the API.

- [Architectural Overview](#) 21
- [Enhance Performance](#) 23
- [Filtering](#) 24
- [Subfile Extraction](#) 24
- [Use the Java Implementation of the API](#) 24
- [The Filter Process Model](#) 27
- [Run File Detection In or Out of Process](#) 31

Architectural Overview

The general architecture of the KeyView Filter technology is the same across all supported platforms and is illustrated in the following diagram.



Each component is described in the following table.

Architectural Components

Component	Description
Developer's Application	The developer's application interfaces directly with the Filter API through either a C-language, Java or .NET implementation.
File Extraction API	The File Extraction API opens a file and extracts the file's subfiles so they are exposed for filtering. See Use the File Extraction API, on page 34 .
Filter API	The Filter API exposes the filtering functionality and controls all other modules during the filtering process. See Use the Filter API, on page 55 .
Format Detection	This module determines the file type of the input stream, allowing the Filter API to return that information to the developer's application, or to load the appropriate structured access layer for further processing. See File Format Detection, on page

Architectural Components, continued

Component	Description
	176 for more information format detection.
Structured Access Layer	There are three modules that reside in the structured access layer—one each for word processing, spreadsheet, and presentation formats. The file detection result determines which structured access layer module is used during the filtering process. That module loads the appropriate document reader and proceeds with text extraction or metadata retrieval.
Document Readers	Each document reader reads a specific file format and sends a text stream of the document to the structured access layer. Each filter is loaded as required by the structured access layer. See Document Readers, on page 200 for a complete list of document readers.

Enhance Performance

KeyView is designed for optimal performance out of the box. However, there are some parameters that you can adjust to improve system performance according to your needs.

File Caching

To reduce the frequency of I/O operations, and consequently improve performance, the KeyView readers load file data into memory. The readers then read the data from the cache rather than the physical disk. You can configure the amount of memory used for file caching through the `formats.ini` file. Generally, when you increase the memory, performance will improve.

By default, KeyView uses a maximum of 1 MB of memory for each thread. If the file data is larger than 1MB, up to 1MB of data is cached and the data beyond 1 MB is read from disk. The minimum amount of memory that can be used for file caching is 64 KB.

To determine a reasonable value, divide the maximum amount of memory you want KeyView to use for file caching by the total number of threads. For example, if you want KeyView to use a maximum of 50MB of memory and have 10 threads, set the value to 5 MB.

To modify the memory allocated for file caching, change the value for the following parameter in the [DiskCache] section of the `formats.ini` file:

```
DiskCacheSize=1024
```

The value is in kilobytes. If this parameter is not set or is set to 0 (zero), the minimum value of 64 KB is used.

The `formats.ini` file is in the directory `install\OS\bin`, where `install` is the pathname of the Filter installation directory and `OS` is the name of the operating system.

Filtering

Filter SDK enables you to *filter* many different types of documents. Filtering is the process of extracting the text from a document without the application-specific markup. However, the filtering process can also include the following:

- Subfile extraction—exposes all subfiles for filtering. See [Use the File Extraction API, on page 34](#).
- File format extraction—detects a file's format, and reports the information to the API, which in turn reports the information to the developer's application. See [File Format Detection, on page 176](#).
- Metadata extraction—extracts selected metadata (document properties) from a file. See [Extract Metadata, on page 58](#).
- Character set conversion—controls the character set of both the input and the output text. See [Convert Character Sets, on page 61](#).

Subfile Extraction

To filter a file, you must first determine whether the file contains any subfiles (attachments, embedded OLE objects, and so on). A file that contains subfiles is called a *container* file. Archive files (such as ZIP), mail messages with attachments (such as Microsoft Outlook Express), mail stores (such as Microsoft Outlook Personal Folders), and compound documents with embedded OLE objects (such as a Microsoft Word document with an embedded Excel chart) are examples of container files.

If the file is a container file, the container must be opened and its subfiles extracted using the File Extraction interface. The extraction process is done repeatedly until all subfiles are extracted and exposed for filtering. Once a subfile is extracted, you can use the Filter API to filter the file.

If a file is not a container, you should pass it directly to the Filter API for filtering without extraction.

The ExtractFilter sample program demonstrates this logic for extracting and filtering files. See [Use the File Extraction API, on page 34](#) for more information.

Use the Java Implementation of the API

The Java version of the Filter API provides an interface to the core functionality of the C API. It contains one primary class (*Filter*) that wraps the filter functionality of the C API. It is implemented in the package `com.verity.api.filter` contained in the file `KeyView.jar`. The JAR file is in the directory `install\javaapi`, where `install` is the path name of the Filter installation directory.

For more information on the Java API, see the Javadoc in the directory `install\javaapi\javadoc`, and [Sample Programs, on page 82](#).

Input/Output Operations

In the Filter Java API, input and output can be either a physical file accessed through a file path, or a Java stream. Depending on the method signature you use, you can create the following filtering processes:

- filter an input file to output file
- filter an input file to an output stream
- filter an input stream to an output stream
- filter an input stream to an output file
- filter an input file and return one chunk of data at a time
- filter an input stream and return one chunk of data at a time

Many methods in the Java API have method signatures supporting one or more of these filtering processes. When you select a method, ensure you use the correct signature for the desired input and output type.

The input source can be set by calling the `setInputSource` method, or when using the `doFilter`, `canFilter`, `canFilterEx`, `getDocFormatInfo`, or `getSummaryInfo` methods. The latter methods take the input source as one of their parameters.

NOTE:

When the input source is from a Java stream, Filter creates an internal buffer from the stream. If the input is a large file, HPE recommends that you use a file as the input source.

Filter in File or Stream Mode

To filter files using the methods in the `Filter` class

1. Instantiate a `Filter` object using either the default constructor or the constructor that sets the output character set and filter flags:
 - a. Use the default constructor `Filter()`. For example:

```
m_objFilter = new Filter();
```
 - b. Use the constructor `Filter(java.lang.String outputCharSet, long filterFlags)`. For example:

```
m_objFilter = new Filter(outputCharSet, Filter.FILTERFLAG_OOPLOGON);
```

The Filter flags provide instructions on how to process a file or stream. For example, they specify whether an error log is generated during filtering (`FILTERFLAG_OOPLOGON`) or whether headers and footers are extracted from the document (`FILTERFLAG_HEADERFOOTERTAGS`).

NOTE:

Filter runs out of process by default. See [The Filter Process Model, on page 27](#) for more information.

2. Set the location of the Filter libraries by calling the `setFilterDirectory(java.lang.String`

directory) method. These libraries are normally stored in the directory `install\OS\bin`, where `install` is the path name of the Filter installation directory and `OS` is the name of the operating system. For example:

```
m_objFilter.setFilterDirectory(m_filterDirectory);
```

3. Set the input source as either a file or input stream by calling the `setInputSource` method.

```
m_objFilter.setInputSource(m_extractDir + filename);
```

4. Filter the file or stream by calling either the `filterTo` or `doFilterChunk` method. The `filterTo` method extracts the data to a file or a stream. The `doFilterChunk` method extracts one chunk of data from a file or a stream. It must be called repeatedly until the entire buffer is filtered.

If filtering in file mode, use the following code:

```
{  
    m_objFilter.filterTo(m_extractDir + filename + m_extension);  
}
```

If filtering in stream mode, use the following code:

```
{  
    outf = new File(m_extractDir + filename + m_extension);  
    fos = new FileOutputStream(outf);  
    m_objFilter.filterTo(fos);  
    fos.close();  
}
```

5. Terminate the filtering session and free allocated system resources by calling the `shutdownFilter()` method.

```
m_objFilter.shutdownFilter();
```

Multithreaded Filtering

To ensure multithreaded filter processes are thread-safe, you must create a unique Filter context for every thread by instantiating a `Filter` object. In addition, threads must not share context objects, and the same context object must be used for all API calls in the same thread. Creating a context object for every thread does not affect performance because the context object uses minimal resources.

For example, your Java code should have the following logic in a thread:

```
m_objFilter = new Filter();  
m_objFilter.setFilterDirectory(m_filterDirectory);  
m_objFilter.setInputSource(infile);  
m_objFilter.getDocFormatInfo();  
  
if (objFilter.canFilter() == true)  
  
    m_objFilter.filterTo(outfile);  
  
m_objFilter.shutdownFilter();
```

Before Running Your Application

Before running your application you must set the library path using one of the following methods:

- On Windows, add the location of `KeyViewFilter.dll` to the `PATH` environment variable.
- On Solaris, Linux, and HP-UX IA-64, add the location of `libKeyViewFilter.so` to the `LD_LIBRARY_PATH` environment variable.
- On HP-UX PA-RISC, add the location of `libKeyViewFilter.sl` to the `SHLIB_PATH` environment variable.
- On AIX, add the location of `libKeyViewFilter.a` to the `LIBPATH` environment variable.
- You can also specify the library path as a system property as follows:

```
java -Djava.library.path=filter_bin_directory ...
```

The Filter Process Model

By default, Filter runs independently from the calling application process. This is called *out-of-process* filtering. Out-of-process filtering protects the stability of the calling application in the rare case when a malformed document causes Filter to fail. You can configure Filter to run in the same process as the calling application. This is called *in-process* filtering. However, it is strongly recommended you run Filter out of process whenever possible.

With the exception of Solaris and AIX, the creation of child processes on UNIX adheres to Portable Operating System Interface (POSIX) standards. Solaris and AIX use thread semantics. If required, a version of `kvfilter` with POSIX thread semantics is available for Solaris and AIX. For Solaris, the file is `kvfilter_posix.so`. For AIX, the file is `kvfilter_nsl.a`. These files must be renamed `kvfilter.so` or `kvfilter.a` to be used by Filter.

To monitor and debug filtering operations during out-of-process filtering, you can generate an error log at run time. See [Generate an Error Log, on page 55](#).

The following methods run in process or out of process:

Filter API

`canFilter`

`canFilterEx`

`doFilter`

`doFilterChunk`

`getSummaryInfo`

`GetDocFormatInfo`

File Extraction API

<code>extCloseDocument</code>	<code>extGetMainFileInfo</code>
<code>extGetSubFileInfo</code>	<code>extOpenDocument</code>
<code>extGetSubFileMetadata</code>	<code>getSummaryInfo</code>
<code>extExtractSubFile</code>	<code>KVGetExtractInterface()</code>

Other Filter API methods always run in process.

Persist the Child Process

By default, in out-of-process filtering, the parent process maintains a persistent connection with the child server after each file is filtered. When the connection is preserved in this way, subsequent filtering requests are processed more quickly because the server is already prepared to receive data.

You can restart the server at regular intervals by using a method or a configuration setting.

In the API

To force KeyView to restart, call the `refreshFilterKV00P()` method.

```
public void refreshFilterKV00P();
```

In the formats.ini File

To control whether Filter persists the server, use the `kvoopRefresh` parameter in the `[FilterSDK_Config]` section of the `formats.ini` file:

`kvoopRefresh= 0` When this is set to 0 (zero), the connection to the server is persisted for as long as the parent process is running or until the server fails. This is the default.

`kvoopRefresh= n` When this is set to `n`, the connection is persisted for `n` filter requests. After the `n`th request, the server is shutdown and restarted before processing the next request.

For example, if `kvooprefresh=5`, the connection to the server is persisted for 5 filter requests. For the 6th request, the server is shutdown and restarted.

To control whether the parent process attempts to filter a file after the file has caused the server to fail, use the `kvoopRetry` parameter in the `[FilterSDK_Config]` section of the `formats.ini` file:

`kvoopRetry= 0` When this is set to 0 and the server fails, the parent process does not resend the file to a new server.

`kvoopRetry= n` When this is set to `n` (a positive number) and the server fails, the parent process resends the file to a new server `n` times. By default, the `kvoopRetry` is set to 1, and the file is resent to a server once.

The `formats.ini` file is in the directory `install\OS\bin`, where `install` is the path name of the Filter installation directory and `OS` is the name of the operating system.

NOTE:

The `kvoopRefresh` and `kvoopRetry` parameters do not apply when running the File Extraction functions out of process. See [Run File Extraction Functions Out of Process](#), below.

Run Filter In Process

By default, Filter runs out of process. However, you can enable in-process filtering through the API or in the `formats.ini` file. If the type of process is not specified in the `formats.ini` or in the API, then Filter is run out of process. If the type of process is specified in the `formats.ini` *and* in the API, the setting in the API takes precedence.

In the API

To run Filter in process, instantiate the Filter object using the constructor `Filter(java.lang.String outputCharSet, long filterFlags)`, and set the `filterFlags` argument to `FILTERFLAG_INPROCESS`.

```
objFilter = new Filter(outputCharSet, Filter.FILTERFLAG_INPROCESS);
```

In the formats.ini File

To run Filter in process, set the following parameter in the `[FilterSDK_Config]` section of the `formats.ini` file to 1:

```
default_inprocess=1
```

By default this is set to 0 (zero), which enables out-of-process filtering.

The `formats.ini` file is in the directory `install\OS\bin`, where `install` is the path name of the Filter installation directory and `OS` is the name of the operating system.

Run File Extraction Functions Out of Process

The out-of-process setting specified when you create the Filter object or in the `formats.ini` is automatically propagated to the File Extraction API. When you extract subfiles from container files and pass the files for filtering out of process, Filter generates a server called `kvoop.exe` for filtering and a duplicate server also called `kvoop.exe` for file extraction. These servers are independent, so if the filtering service stops responding, the file extraction service can continue extracting files uninterrupted.

Restart the File Extraction Server

If the file extraction server fails on a file and throws the exception `KVError_InvalidOopDriverSignature` or `KVError_InvalidOopServiceSignature`, you must restart the server by recreating the Filter object, and process the source file again.

Out-of-Process Logging

Logging is available for out-of-process filtering. The kvoop server can now create a log file that captures information on the files being processed, storing one entry per process. The generated log file is called `xxxx_kvoop.log`, where `xxxx` is a unique number identifying the process.

In the rare case when the kvoop server fails, you can use the log files to determine which file caused the failure. After processing is complete and the system shuts down, the logs are automatically deleted. To keep the log files after processing is successfully completed, see [Keep Log Files, on the next page](#).

NOTE:

Out-of-process logging is not supported on AIX.

Enable Out-of-Process Logging

To enable out-of-process logging, set the `KVOOP_LOGS_DIR` environment variable to the directory in which you want the log files to be stored. By default, logging is not enabled.

On UNIX, the variable is set as follows:

```
setenv KVOOP_LOGS_DIR /tmp
```

On Windows, the variable is set as follows:

```
set KVOOP_LOGS_DIR=c:\tmp
```

The following log file is created in the directory:

process_id_kvoop.log

where *process_id* is a numeric value representing the logged process. New messages are appended to the file, and truncation is disabled by default.

If KeyView terminates unexpectedly and Windows minidump is enabled, a *process_id_crash_info.txt* file is generated (see [Enable Windows Minidump, on the next page](#)). If logging was not been enabled at the time of termination, this file contains instructions on how to enable logging.

Set the Verbosity Level

You can control how much information is written to the file by setting the `KVOOP_LOG_VERBOSITY` environment variable. For example:

```
set KVOOP_LOG_VERBOSITY=1
```

The variable can be set to the following:

- 1 Include only error messages.
- 2 Include errors and warnings.

- 3 Include errors, warnings, and general information. This is the default.
- 4 Include all possible information. This setting is useful for debugging purposes.

Enable Windows Minidump

KeyView can use the Windows minidump feature to provide additional logging information, which can be useful for debugging purposes.

The Windows minidump is disabled by default. To enable the Windows minidump, set `KV00P_DUMP_ENABLE=1`. If an unexpected termination occurs after the minidump is enabled, three files are generated:

- *process_id_crash_info.txt*. This file contains KV00P state and runtime information at the time of termination. If logging was not enabled at the time of termination, this file contains instructions on how to enable logging.
- *process_id_process_list.txt*. This file contains information from the DLLs that were loaded at the time of the termination.
- *process_id_report.dmp*. The Windows dump file, which contains further information about the termination. You can open it with either a Windows debugger or `autnhelper.exe` (you must copy this file to the same directory).

You can control the amount of information presented in the Windows dump file by creating the following files in the directory:

```
dumper.NORMAL  
dumper.WITHDATASEG  
dumper.WITHFULLMEMORY  
dumper.WITHHANDLEDATA
```

Keep Log Files

After processing is complete and the system is shut down, the log files are automatically deleted from the directory. To keep the log files after a successful run, set the `KV00P_KEEP_LOGS` environment variable.

On UNIX, set the variable as follows:

```
setenv KV00P_KEEP_LOGS 1
```

On Windows, set the variable as follows:

```
set KV00P_KEEP_LOGS=1
```

Run File Detection In or Out of Process

By default, detection runs in out-of-process mode. However, you can enable in-process detection through the API or in the `formats.ini` file. If the type of process is not specified in the `formats.ini` or in the API, Filter runs in out-of-process mode. If the type of process is specified in the `formats.ini` and in the API, the setting in the API takes precedence.

Specify the Process Type In the formats.ini File

Add the `default_detect_inprocess` flag to a `[FilterSDK_Config]` section in the `formats.ini` file to control the default behavior for detection. Set `default_detect_inprocess` to `0` for out-of-process detection, and `1` for in-process detection. For example:

```
[FilterSDK Config]
default_detect_inprocess=0
```

If this flag is not specified, the file detection behavior is determined by the `default_inprocess` flag for filtering. For example, if you set `default_inprocess` to `1`, filtering and file detection runs in in-process mode by default; if you set `default_inprocess` to `0`, filtering and file detection runs in out-of-process mode by default.

If both the `default_inprocess` and `default_detect_inprocess` flags are set, then `default_inprocess` controls the default filtering behavior and `default_detect_inprocess` controls the default file detection behavior.

Specify the Process Type In the API

To run detection in in-process mode, instantiate the `Filter` object by using the constructor `Filter(java.lang.String outputCharSet, long filterFlags)`, and set the `filterFlags` argument to `FILTERFLAG_DETECTINPROCESS`. To run detection in out-of-process mode, set `FILTERFLAG_DETECTOUTOFPROCESS`.

```
objFilter = new Filter(outputCharSet, Filter.FILTERFLAG_DETECTINPROCESS);
```


Part II: Use Filter SDK

This section explains how to perform some basic tasks by using the File Extraction and Filter APIs, and describes the sample programs.

Chapter 3: Use the File Extraction API

This section describes how to extract subfiles from a container file using the File Extraction API.

- [Introduction](#)34
- [Extract Subfiles](#)35
- [Extract Images](#)36
- [Recreate a File Hierarchy](#)36
- [Extract Mail Metadata](#)38
- [Extract Subfiles from Outlook Files](#)43
- [Extract Subfiles from Outlook Express Files](#)43
- [Extract Subfiles from Mailbox Files](#)43
- [Extract Subfiles from Outlook Personal Folders Files](#)44
- [Extract Subfiles from Lotus Domino XML Language Files](#)47
- [Extract Subfiles from Lotus Notes Database Files](#)48
- [Extract Subfiles from PDF Files](#)51
- [Extract Embedded OLE Objects](#)52
- [Extract Subfiles from ZIP Files](#)52
- [Default File Names for Extracted Subfiles](#)52
- [Exclude Japanese Guide Text](#)54

Introduction

To filter a file, you must first determine whether the file contains any subfiles (attachments, embedded OLE objects, and so on). A file that contains subfiles is called a *container* file. A container file has a main file (parent) and subfiles (children) embedded in the main file.

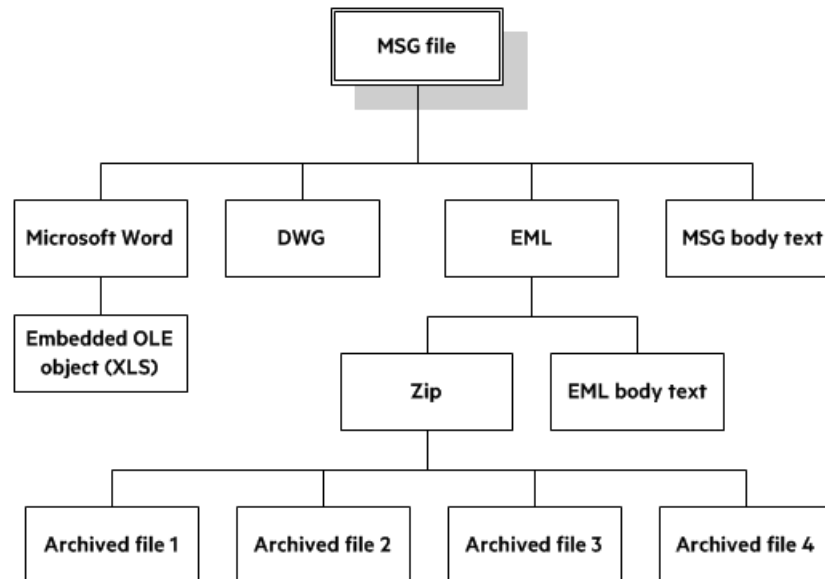
The following are examples of container files:

- Archive files such as ZIP, TAR, and RAR.
- Mail messages such as Outlook (MSG) and Outlook Express (EML).
- Mail stores such as Microsoft Outlook Personal Folders (PST), Mailbox (MBX), and Lotus Notes database (NSF).
- PDF files that contain file attachments.
- Compound documents with embedded OLE objects such as a Microsoft Word document with an embedded Excel chart.

NOTE: [Supported Formats](#), on page 94 indicates which formats are treated as container files and are supported by the File Extraction API.

The subfiles might also be container files, creating a file hierarchy of multiple levels. For example, an MSG file (the root parent) might contain three attachments:

- a Microsoft Word document that contains an embedded Microsoft Excel spreadsheet.
- an AutoCAD drawing file (DWG).
- an EML file with an attached Zip file, which in turn contains four archived files.



NOTE: The parent MSG file contains four first-level children. The body text of a message file, although not a standalone file in the container, is considered a child of the parent file.

Extract Subfiles

To filter all files in a container file, the container must be opened and its subfiles extracted to either a file or a stream using the *File Extraction API*. The extraction process is done repeatedly until all subfiles are extracted and exposed for filtering. Once a subfile is extracted, you can call Filter API methods to filter the data.

If you require a container file, including subfiles, to be filtered to a single file, you must extract all files from the container, filter the files, and then append each filtered output to its parent.

To extract subfiles, follow this general procedure

1. Open the source file by calling the `extOpenDocument` method. This call defines the parameters necessary to open a file for extraction.
2. Determine whether the main file is a container file (contains subfiles) by calling the `extGetMainFileInfo()` method.
3. If the call to `extGetMainFileInfo()` determined the source file is a container file, proceed to [Determine whether the subfile is itself a container \(contains subfiles\) by calling the `extGetSubFileInfo` method., below](#); otherwise, filter the file.
4. Determine whether the subfile is itself a container (contains subfiles) by calling the

`extGetSubFileInfo` method.

5. Extract the subfile by calling the `extExtractSubFile` method.
6. If the call to `extGetSubFileInfo` determined the subfile is a container file, repeat [Open the source file by calling the `extOpenDocument` method](#). This call defines the parameters necessary to open a file for extraction. , on the previous page through [Extract the subfile by calling the `extExtractSubFile` method.](#), above until all subfiles are extracted and the lowest level of subfiles is reached; otherwise, filter the file.

Extract Images

You can use the File Extraction API to extract images within the file by specifying the following in the `formats.ini` file:

```
[Options]
ExtractImages=TRUE
```

If you set this option, images within the file behave in the same way as any other subfile. Extracted images have the name `image[X].[Y]`, where `[X]` is an integer, and `[Y]` is the extension. The format of the image is the same as the format in which it is stored in the document.

This option can also be enabled by passing `KVFLT_EXTRACTIMAGES` to the `fpFilterConfig` function.

Recreate a File Hierarchy

When a container file is extracted, any relationships between the subfiles in the container are not maintained. However, the File Extraction interface provides information that enables you to recreate the hierarchy. The hierarchy can be used to create a directory structure in a file system, or to categorize documents according to their relationship to each other. For example, if you use `KeyView` to generate text for a search engine, the hierarchical information enables your users to search for a document based on the document's parent or sibling. In addition, when the document is returned to the user, the parent and sibling documents can be returned as recommendations.

The information needed to recreate a file's hierarchy is provided in the call to `extGetSubFileInfo`. Call this method to retrieve an object of the `ExtSubFileInfo` class, then use the `getParentIndex()` and `getChildArray()` methods in this object to retrieve information about the subfile's parent and children. Since you can only retrieve the first-level children in a subfile, you must call `extGetSubFileInfo` repeatedly until information for the leaf-node children is extracted.

Create a Root Node

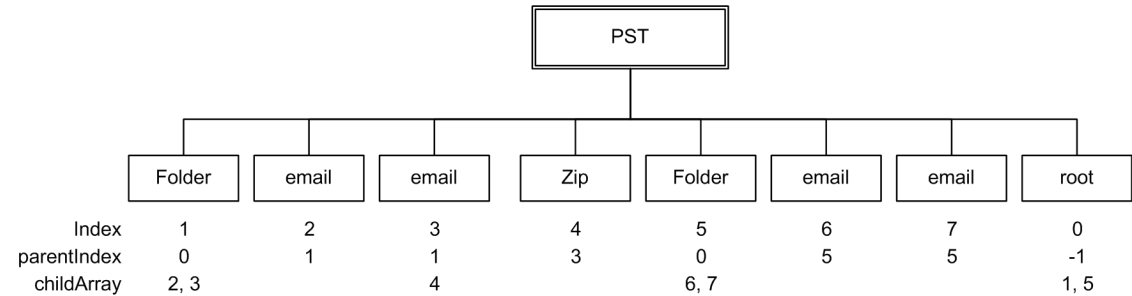
Because of their structure, some container files do not contain a subfile or folder which acts as a root directory on which the hierarchy can be based. For example, subfiles in a Zip archive can be extracted, but none of the subfiles represent the root of the hierarchy. In this case, an artificial *root node* must be created at the top of the file hierarchy as a point of reference for each child, and ultimately to recreate the relationships. This artificial root node is an internal object, and is extracted to disk as a directory called `root`. Its index number is 0.

To create a root node, call the `setCreateNode` method in the `ExtOpenDocConfig` object, and pass `ExtOpenDocConfig` to the `extOpenDocument` method. When a root node is created, the value returned from the `getNumSubFiles` method in the `ExtMainFileInfo` object includes the root node. For example, when you call `extGetMainFileInfo` on a Microsoft Word document with three embedded OLE objects and the root node is disabled, the number of subfiles is 3. If you create a root node, the number of subfiles is 4.

Example

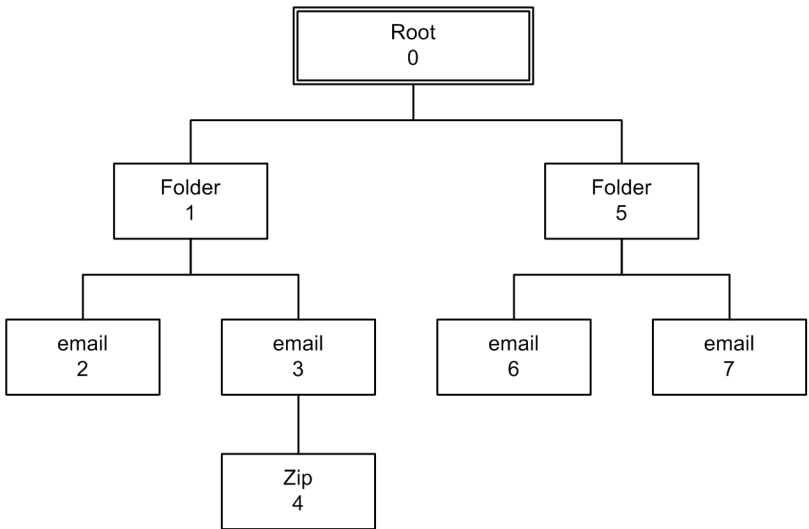
For example, you might extract a PST file that contains seven subfiles with a root node enabled. The call to `extGetMainFileInfo()` returns the number of subfiles as 8 (seven subfiles and one root node). [Extracted PST File, below](#) shows the structure and the available hierarchy information after the subfiles are extracted:

Extracted PST File



The `parentIndex` specifies the index number of a subfile's parent. The `childArray` specifies an array of a subfile's children. With this information, you can recreate the hierarchy shown in [Recreated File Hierarchy, below](#):

Recreated File Hierarchy



Extract Mail Metadata

You can extract metadata such as subject, sender, and recipient from MSG, EML, MBX, PST, and NSF files by calling the `extGetSubFileMetadata()` method. You can extract a predefined set of metadata fields, individual fields, or both, that are unique to a file format.

Default Metadata Set

KeyView internally defines a set of common mail metadata fields that can be extracted as a group from mail formats. This default metadata set is listed in [Default Mail Metadata List, below](#). When you retrieve all metadata for a file—that is, pass `NULL` for the array of metadata—the complete set of default metadata, not all available metadata in the file, is returned.

Default Mail Metadata List

Field Name (string to specify)	Description
From	The display name and email address of the sender.
Sent	The time the message was sent.
To	The display names and email addresses of the recipients.
Cc	The display names and email addresses of recipients who receive copies of the email.
Bcc	The display names and email addresses of recipients who received blind copies of the email.
Subject	The text in the subject line of the message.
Priority	The priority applied to the message.

Because mail formats use different terms for the same fields, the format's reader maps the default field name to the appropriate format-specific name. For example, when retrieving the default metadata set, the NSF field *Importance* is mapped to the name *Priority* and is returned.

You can also extract the default field names individually by passing the field name (such as *From*, *To*, and *Subject*); however, in this case, the string is not mapped to the format-specific name. For example, if you pass *Priority* in the call, you will retrieve the contents of the *Priority* field from an MBX file, but will not retrieve the contents of the *Importance* field from an NSF file.

NOTE: You cannot pass the field names listed in [Default Mail Metadata List, above](#) individually for PST files. However, you can pass either the MAPI tag number or one of the constants in the Filter class as integers. See [Microsoft Personal Folders File \(PST\) Metadata, on page 41](#).

Extract the Default Metadata Set

To extract the default metadata set, call the `extGetSubFileMetadata(long docContextID, int nSubFileIndex, ExtSubFileMetaConfig config)` method.

For example:

```
ExtSubFileMetaConfig metaConfig = new ExtSubFileMetaConfig();  
ExtSubFileMetadata subfilemeta = null;  
subfilemeta = m_objFilter.extGetSubFileMetadata(extContextID, index, metaConfig);
```

Microsoft Outlook (MSG) Metadata

In addition to the default metadata set, the metadata fields listed in [MSG-Specific Metadata List](#), below can be extracted for MSG files. The field name must be passed to `metaNameArray` in the call to the `extGetSubFileMetadata()` method.

MSG-Specific Metadata List

Field Name (string to specify)	Description
AttachFileName	An attachment's long file name and extension, excluding path.
ConversationTopic	The topic of the first message in a conversation thread. A conversation thread is a series of messages and replies. This is the first message's subject with any prefix removed.
CreationTime	The time the message or attachment was created. This value is displayed in the Sent field in the message's Properties dialog in Outlook.
InternetMessageID	The identifier for messages that come in over the Internet. This is the MAPI property <code>PR_INTERNET_MESSAGE_ID</code> . This property is not in the MAPI headers or MAPI documentation.
LastModificationTime	The time the message or attachment was last modified. This value is displayed in the Modified field in the message's Properties dialog in Outlook.
Location	The physical location of the event specified in the Outlook calendar entry.
MessageID	The message transfer system (MTS) identifier for the message transfer agent (MTA). This value is displayed on the Message ID tab in the message's Properties dialog in Outlook.
Received	The date and time a message was delivered. This value is displayed in the Received field in the message's Properties dialog in Outlook.
Sender	The name and email address of the message sender. This value is a concatenation of two MAPI properties in the following format:

MSG-Specific Metadata List, continued

Field Name (string to specify)	Description
	"PR_SENDER_NAME" <PR_SENDER_EMAIL_ADDRESS> The Sender value might be the same as or different than the default metadata From value (see Default Metadata Set, on page 38), depending on which MAPI properties exist in the MSG file.
Sensitivity	The value indicating the message sender's opinion of the sensitivity of a message, such as Personal, Private, or Confidential. This value is displayed in the Sensitivity field in the message's Properties dialog in Outlook.
TransportMsgHeaders	Contains transport-specific message envelope information. This value corresponds to the MAPI property PR_TRANSPORT_MESSAGE_HEADERS.
StartDate	Contains an appointment start date. This value corresponds to the PR_START_DATE MAPI property.
EndDate	Contains an appointment end date. This value corresponds to the PR_END_DATE MAPI property.

Extract MSG-Specific Metadata

To extract specific metadata fields from an MSG file, use the method `extGetSubFileMetadata(long docContextID, int nSubFileIndex, java.lang.String[] metaNameArray, ExtSubFileMetaConfig config)` and pass the field name defined in [MSG-Specific Metadata List, on the previous page](#) to `metaNameArray` (the string is not case sensitive).

For example, the following code extracts the contents of the `ConversationTopic` and `MessageID` fields:

```
ExtSubFileMetaConfig metaConfig = new ExtSubFileMetaConfig();  
ExtSubFileMetadata subfilemeta = null;  
String[] metaNameArray = {"conversationtopic", "MessageID"};  
subfilemeta = m_objFilter.extGetSubFileMetadata(extContextID, index, metaNameArray,  
metaConfig);
```

Microsoft Outlook Express (EML) and Mailbox (MBX) Metadata

In addition to the default metadata set, you can extract any metadata field that exists in the header of an EML or MBX file by passing the field's name. If the name is a valid field in the file, the contents of the field are returned. For example, to retrieve the name of the last mail server that received the message before it was delivered, you can pass the string "Received".

Extract EML- or MBX-Specific Metadata

To extract specific metadata fields from an EML or MBX file, use the method `extGetSubFileMetadata` (`long docContextID`, `int nSubFileIndex`, `java.lang.String[] metaNameArray`, `ExtSubFileMetaConfig config`) and pass the metadata name to `metaNameArray` (the string is not case sensitive).

For example, the following code extracts the contents of the Received and Mime-version fields:

```
ExtSubFileMetaConfig metaConfig = new ExtSubFileMetaConfig();  
ExtSubFileMetadata subfilemeta = null;  
String[] metaNameArray = {"Received", "Mime-version"};  
subfilemeta = m_objFilter.extGetSubFileMetadata(extContextID, index, metaNameArray,  
metaConfig);
```

Lotus Notes Database (NSF) Metadata

In addition to the default metadata set, you can extract any Lotus field name that exists in an NSF file by passing the field's name. (You can extract fields from mail NSF files and non-mail NSF files.) If the name is a valid field in the file, the field is returned. For example, to retrieve the date a document in an NSF file was last accessed, you would pass the string `"$LastAccessedDB"`.

NOTE: A complete list of NSF fields are provided in the Lotus Notes file `stdnames.h`. This header file is available in the Lotus API Toolkit.

Extract NSF-Specific Metadata

To extract specific metadata fields from an NSF file, use the method `extGetSubFileMetadata` (`long docContextID`, `int nSubFileIndex`, `java.lang.String[] metaNameArray`, `ExtSubFileMetaConfig config`) and pass the metadata name to `metaNameArray` (the string is not case sensitive).

For example, the following code extracts the contents of the Description and Categories fields:

```
ExtSubFileMetaConfig metaConfig = new ExtSubFileMetaConfig();  
ExtSubFileMetadata subfilemeta = null;  
String[] metaNameArray = {"description", "Categories"};  
subfilemeta = m_objFilter.extGetSubFileMetadata(extContextID, index, metaNameArray,  
metaConfig);
```

Microsoft Personal Folders File (PST) Metadata

In addition to the default metadata set, you can extract Messaging Application Programming Interface (MAPI) properties from a PST file. These properties describe elements (subject, sender, recipient, and so on) of Outlook items within the PST file. Since the properties are stored in the PST file itself, they

can be retrieved before the contents of the PST are extracted. This enables you to determine whether an Outlook item should be extracted based on a subfile's attributes. MAPI properties are also stored for Outlook attachments that are not mail messages (such as an attached Microsoft Word document or Lotus 1-2-3 file).

MAPI Properties

Each MAPI property is identified by a property tag, which is a constant that contains the property type and a unique identifier. For example, the property that indicates whether a message has attachments has the following components:

Property	PR_HASATTACH
Identifier	0x0E1B
Property type	PT_BOOLEAN (000B)
Property tag	0x0E1B000B

The Microsoft MAPI documentation on the Microsoft Developer Network website lists all available MAPI properties, their tags, and types.

You can retrieve any MAPI property that is of one of the MAPI property types listed below:

PT_I2	PT_DOUBLE	PT_STRING8
PT_I4	PT_FLOAT	PT_TSTRING
PT_BINARY	PT_LONG	PT_SYSTIME
PT_BOOLEAN	PT_SHORT	PT_UNICODE

NOTE: Properties with a PT_TSTRING type have the property type recompiled to either a Unicode string (PT_UNICODE) or to an ANSI string (PT_STRING8) depending on the operating system's character set. To retrieve the Unicode property, pass in the Unicode version of the tag. For example, the property tag for PR_SUBJECT is either 0x0037001E for an ANSI string, or 0x0037001F for a Unicode string.

Extract PST-Specific Metadata

In the call to extract subfile metadata, you can pass either the MAPI tag number (such as 0x0070001e) or one of the constants in the Filter class (such as KVPR_SUBJECT). These constants are a subset of MAPI properties and use a KeyView naming convention. For example, the property PR_CONVERSATION_TOPIC is defined as KVPR_CONVERSATION_TOPIC. If the property you want to retrieve is not defined as a constant in the Filter class, you must pass the MAPI tag number.

To extract specific MAPI properties from a PST file, use the method `extGetSubFileMetadata(long docContextID, int nSubFileIndex, int[] metaNameArray, ExtSubFileMetaConfig config)` and pass the tag number or constant to `metaNameArray`.

For example, the following code extracts the MAPI properties PR_SUBJECT and PR_ALTERNATE_RECIPIENT:

```
ExtSubFileMetaConfig metaConfig = new ExtSubFileMetaConfig();  
ExtSubFileMetadata subfilemeta = null;
```

```
int[] metaNameArray = {Filter.KVPR_SUBJECT, 0x3A010102};  
  
subfilemeta = m_objFilter.extGetSubFileMetadata(extContextID, index, metaNameArray,  
metaConfig);
```

Exclude Metadata from the Extracted Text File

When a mail message is extracted, the message text and header information (To, From, Sent, and so on) is also extracted. You can prevent the header information from appearing in the text file.

To exclude the header information, call the `setExcludeMailHeader()` method of the `ExtSubFileExtractConfig` object, and pass `ExtSubFileExtractConfig` to the `extExtractSubFile` method. For example:

```
m_excludeMailHeader = true;  
  
extconfig = new ExtSubFileExtractConfig();  
  
extconfig.setExcludeMailHeader(m_excludeMailHeader);  
  
extinfo = m_objFilter.extExtractSubFile(extContextID, i, extconfig);
```

Extract Subfiles from Outlook Files

When you extract an Outlook file (MSG) to disk, the message text and header information (To, From, Sent, and so on) is extracted to a text file. (If you do not want the header information to appear in the text file, see [Exclude Metadata from the Extracted Text File, above](#).) If the Outlook file contains a non-mail attachment, the attachment is extracted in its native format to a subdirectory. If the Outlook file contains a mail attachment, the attachment's message text and any attachments are extracted to a subdirectory.

Extract Subfiles from Outlook Express Files

When you extract an Outlook Express (EML) file to disk, the message text and header information (To, From, Sent, and so on) is extracted to a text file. (If you do not want the header information to appear in the text file, see [Exclude Metadata from the Extracted Text File, above](#).) If the Outlook Express file contains a non-mail attachment, the attachment is extracted in its native format to the same directory as the message text file. If the Outlook Express file contains a mail attachment, the complete attachment (including message text and attachments), the message text file, and any non-mail attachments are extracted to the same directory as the main message.

NOTE: When the MBX reader (`mbxsr`) is enabled, it is used to filter MBX and EML files. If the MBX reader is not enabled, the EML reader (`emlsr`) is used.

Extract Subfiles from Mailbox Files

A Mailbox (MBX) file is a collection of individual emails compiled with RFC 822 and RFC 2045 - 2049 (MIME), and divided by message separators. There are many mail applications that export to an MBX

format, such as Eudora Email and Mozilla Thunderbird.

When an MBX file is extracted to disk, the message text and header information (To, From, Sent, and so on) from each mail file are extracted to text files. (If you do not want the header information to appear in the text file, see [Exclude Metadata from the Extracted Text File, on the previous page.](#))

In Eudora MBX files, attachments are inserted as a link and are stored externally from the message. These attachments are not extracted, but the path to the attachment is returned in the call to the `extGetSubFileInfo` method. You can write code to retrieve the attachment based on the returned path.

For MBX files from other clients, KeyView extracts attachments when they are embedded in the message.

NOTE: The Mailbox (MBX) reader is an advanced feature and is sold and licensed separately. To enable this reader in a KeyView SDK, you must obtain the appropriate license key from HPE.

Extract Subfiles from Outlook Personal Folders Files

KeyView can extract Outlook items such as messages, appointments, contacts, tasks, notes, and journal entries from a PST file. When a PST file is extracted to disk, the body text and header information (To, From, Sent, and so on) from each Outlook item is extracted to a text file. (If you do not want the header information to appear in the text file, see [Exclude Metadata from the Extracted Text File, on the previous page.](#))

You can also extract messages from PST files as MSG files, including all their attachments, using the `setSaveAsMSG()` method in the `ExtSubFileExtractConfig` class.

If an Outlook item contains a non-mail attachment, the attachment is extracted in its native format to a subdirectory. If an Outlook item contains an Outlook attachment, the attached item's body text and any attachments are extracted to a subdirectory.

NOTE: The Microsoft Outlook Personal Folders (PST) reader is an advanced feature and is sold and licensed separately. To enable this reader in a KeyView SDK, you must obtain the appropriate license key from HPE.

Use the Native or MAPI-based Reader

KeyView accesses PST files in one of two ways:

- indirectly using the Microsoft Messaging Application Programming Interface (MAPI) reader named `pstsr`.
- directly using the native PST reader named `pstnsr`.

On UNIX platforms, the native reader is always used to process PST files because the MAPI-based reader runs only on Windows x86 and x64. On Windows, you can specify either reader; however, the MAPI-based reader is used by default.

The differences between the two readers are summarized in the following table:

Feature/Requirement	Native Reader (pstnsr)	MAPI-based Reader (pstsrr)
All platforms supported	Yes	Windows x86 and x64 only
Outlook client required	No	Yes
MAPI properties supported	Yes All properties defined in mapitags.h. Object properties are not supported.	Yes Extracts properties defined in mapitags.h. Object properties are not supported.
Password protection supported	Yes	Yes (using the setPassword method)
Compressible encryption supported	Yes	Yes
High encryption supported	No	Yes

To use the MAPI-based reader for PST files, change the PST entry in the `formats.ini` file as follows:

```
297=pst
```

To use the native reader for PST files, change the PST entry in the `formats.ini` file as follows:

```
297=pstn
```

NOTE: You must make sure that the PST that you are extracting is not open in the Outlook client, and that the Outlook process is not running.

Use the Native PST Reader (pstnsr)

The native PST reader accesses PST files directly without relying on the Microsoft interface to the PST format. It runs on both Windows and UNIX, and does not require an Outlook client on the system processing the PST files. However, the native reader does not support password-protected PST files that use high encryption.

Use the MAPI Reader (pstsrr)

The `pstsrr` reader accesses PST files indirectly by using Microsoft's Messaging Application Programming Interface (MAPI). MAPI is a standard Windows message interface that enables different mail programs and other mail-aware applications (such as word processors and spreadsheets) to exchange messages and attachments with each other. MAPI allows KeyView to open a PST file, traverse the folders and Outlook items, and extract the items inside the PST file.

NOTE: When extracting subfiles from PST files, information on the distribution list used in an email is extracted to a file called `emailname.dist`. This applies to the MAPI reader (`pstsrr`) only.

System Requirements

MAPI is supported on Windows platforms only, so you can filter PST files on Windows only. Because MAPI relies on functionality in Microsoft Outlook, a Microsoft Outlook client must be installed on the same machine as the application filtering PST files, and must be the default email application. KeyView supports the following PST formats and Outlook clients:

- Outlook 97 or higher PST files
- Outlook 2002 or Outlook 2003 clients

NOTE: The Outlook client must be the same version as, or newer than, the version of Outlook that generated the PST file.

NOTE: The bit edition of Microsoft Outlook must match that of the KeyView software. For example, if 32-bit KeyView is used, 32-bit Outlook must be installed. If 64-bit KeyView is used, 64-bit Outlook must be installed.

If the bit editions do not match, an error message from Microsoft Office Outlook is displayed:

Either there is a no default mail client or the current mail client cannot fulfill the messaging request. Please run Microsoft Outlook and set it as the default mail client.

Additionally, KeyView displays the following return code:

Error 32: KVErrror_PSTAccessFailed.

MAPI Attachment Methods

The way in which you can access the contents of a PST message attachment is determined by the MAPI *attachment method* applied to the attachment. For example, if the attachment is an embedded OLE object, it uses the ATTACH_OLE attachment method. KeyView can access message attachments that use the following attachment methods:

ATTACH_BY_VALUE

ATTACH_EMBEDDED_MSG

ATTACH_OLE

ATTACH_BY_REFERENCE

ATTACH_BY_REF_ONLY

ATTACH_BY_REF_RESOLVE

Attachments using the ATTACH_BY_VALUE, ATTACH_EMBEDDED_MSG, or ATTACH_OLE attachment methods are extracted automatically when the PST file is extracted. An "attach by reference" method means that the attachment is not in Outlook, but Outlook contains an absolute path to the attachment. Before you can extract these types of attachments, you must retrieve the path to access the attachment.

To extract "attach by reference" attachments

1. Determine whether the attachment uses an ATTACH_BY_REFERENCE, ATTACH_BY_REF_ONLY, or ATTACH_BY_REF_RESOLVE method by retrieving the MAPI property PR_ATTACH_METHOD.
2. If the attachment uses one of the "attach by reference" methods, get the fully qualified path to the attachment by retrieving the MAPI properties PR_ATTACH_LONG_PATHNAME or PR_ATTACH_PATHNAME.
3. You can then either copy the files from their original location to the path where the PST file is extracted, or use the Filter API methods to filter the attachment.

Open Secured PST Files

KeyView enables you to specify credentials (user name and password), which are used to open a secured PST file for extraction. See [Password Protected Files, on page 226](#) for more information.

Detect PST Files While the Outlook Client is Running

If you are running an Outlook client while running the File Extraction API, the KeyView format detection module (kwad) might not be able to open the PST file to determine the file's format because Outlook has the file locked. In this case, you can do one of the following:

- Close Outlook when using the Extraction API
- Detect PST files by extension only and bypass the format detection module. To enable this option, add the following lines to the `formats.ini` file.

```
[container_flags]
```

```
detectPSTbyExtension=1
```

NOTE: The `detectPSTbyExtension` option only applies when you are using the MAPI reader (pstr).

NOTE: If you use this option, you must make sure in your code that valid PST files are passed to KeyView because the format detection module will not be available to verify the file type and pass the file to the appropriate reader.

Extract Subfiles from Lotus Domino XML Language Files

When you extract a Lotus Domino XML Language (.DXL) file, the message text and header information (*To*, *From*, *Sent*, and so on) is extracted to a text file.

NOTE: To prevent header information from being extracted, see [Exclude Metadata from the Extracted Text File, on page 43](#).

You can make sure that dates and times extracted from Lotus Domino .DXL files are displayed in a uniform format.

To extract custom date/time formats

- In the `formats.ini` file, set the `DateTimeFormat` option in the `[dxlsr]` section. For example:

```
[dxlsr]
DateTimeFormat=%m/%d/%Y %I:%M:%S %p
```

In this example, dates and times are extracted in the following format:

`02/11/2003 11:36:09 AM`

The format arguments are the same as those for the `strftime()` function. See <http://msdn.microsoft.com/en-us/library/fe06s4ak%28VS.71%29.aspx> for more information.

Extract .DXL Files to HTML

You can use the file extraction API to process .DXL files with an XSLT engine. The XSLT engine then transforms the extracted .DXL to .mail HTML files.

To extract .DXL files to HTML

- Set the following options in the `formats.ini` file:

```
[nsfsr]
ExportDXL=1
ExportDXL_PureXML=1
[dxlsr]
LNDDParser=2
```

Extract Subfiles from Lotus Notes Database Files

A Lotus Notes database is a single file that contains multiple documents called *notes*. Notes include design notes (such as forms, views, folders, navigators, outlines, pages, framesets, agents, and resources), data document notes, profile document notes, access control list notes, and collection (index) notes. KeyView can extract text items, attachments, and OLE objects from *data document notes* only. Data document notes include emails, journal entries, discussion threads, documents (Microsoft Office and Lotus SmartSuite), and so on.

All components of a note are prefixed by field names such as "SendTo:", "Subject:", and "Body:". When a note is extracted, the field names are not included in the extracted output; only the field values are extracted.

When a mail message in an NSF file is extracted to disk, the body text and header information (such as the values from the `SendTo`, `From`, and `DeliveredDate` fields) in each message is extracted to a text file. (If you do not want the header information to appear in the message text file, see [Exclude Metadata from the Extracted Text File, on page 43.](#))

NOTE: The Lotus Notes Database (NSF) reader is an advanced feature and is sold and licensed separately. To enable this reader in a KeyView SDK, you must obtain the appropriate license key from HPE.

System Requirements

The NSF format is proprietary. Therefore, KeyView accesses NSF files indirectly by using the Lotus Notes API. Because the NSF reader relies on functionality in Lotus Notes, a Lotus Notes client or Lotus Domino server must be installed and configured on the same machine on which the application filtering NSF files is installed. On UNIX and Linux, the Lotus Domino server is required. On Windows, the Lotus Notes client or Lotus Domino server is required.

KeyView supports the following Lotus Notes clients and Domino servers:

- Lotus Notes 6.5.1
- Lotus Domino 6.5.1

KeyView supports NSF files on the same platforms supported by Lotus Notes and Lotus Domino:

- Windows XP x86 (Service Pack 1 and 2)
- Windows 2000 x86 (Service Pack 2)
- Solaris 8.0 and 9.0 (built on Solaris 8.0)
- Red Hat Enterprise Linux AS 3.0 (x86)
- SuSE Linux Enterprise Server 8 and 9 (x86)
- IBM AIX 5.1, 5L version 5.2

Installation and Configuration

Before KeyView can filter NSF files, you must set up the Lotus Notes client or Lotus Domino server. Full configuration is not required. The following steps outline the minimal setup for NSF filtering.

Windows

1. Install the Lotus Notes client or Lotus Domino server. You do not need to configure the client or server.
2. Make sure that the `notes.ini` file is in the proper location.
 - If Lotus Notes is installed, the file should appear in the `install\lotus\notes` directory, where `install` is the installation directory.
 - If only Lotus Domino is installed, the file should appear in the `install\lotus\domino` directory, where `install` is the installation directory.

If the file does not exist, create an ASCII file named `notes.ini`, and add the following text:

```
[Notes]
```

3. Add the KeyView `bin` directory and the `install\lotus\notes` or `install\lotus\domino` directory to the `PATH` environment variable (the KeyView `bin` directory must be first in the path). HPE recommends that you add the KeyView `bin` directory because the Lotus Notes or Domino server installation might contain older KeyView OEM libraries.

Solaris

1. Install Lotus Domino server. You do not need to configure the server.
2. Make sure that the file `notes.ini` is in the `install/lotus/notes/latest/sunspa` directory, where `install` is the directory where Lotus Notes is installed. If the file does not exist, create an ASCII file named `notes.ini`, and add the following text:

```
[Notes]
```

3. Add the `install/lotus/notes/latest/sunspa` directory to the PATH environment variable:

```
setenv PATH install/lotus/notes/latest/sunspa:$PATH
```

4. Add the `install/lotus/notes/latest/sunspa` and the KeyView bin directory to the LD_LIBRARY_PATH environment variable:

```
setenv LD_LIBRARY_PATH keyview_bin:install/lotus/notes/latest/sunspa:$LD_LIBRARY_PATH
```

where `keyview_bin` is the location of the KeyView bin directory. HPE recommends that you add the KeyView bin directory because the Lotus Notes installation might contain older KeyView OEM libraries.

AIX 5.x

1. Install the `bos.iocp.rte` file set if it is not already installed, and reboot the machine. See the Lotus Domino server documentation for more information.
2. Install Lotus Domino server. You do not need to configure the server.
3. Make sure that the file `notes.ini` is in the `install/lotus/notes/latest/ibmpow` directory, where `install` is the directory where Lotus Notes is installed. If the file does not exist, create an ASCII file named `notes.ini`, and add the following text:

```
[Notes]
```

4. Add the `install/lotus/notes/latest/ibmpow` directory to the PATH environment variable:

```
setenv PATH install/lotus/notes/latest/ibmpow:$PATH
```

5. Add the `install/lotus/notes/latest/ibmpow` and the KeyView bin directory to the LIBPATH environment variable:

```
setenv LIBPATH keyview_bin:install/lotus/notes/latest/ibmpow:$LIBPATH
```

where `keyview_bin` is the location of the KeyView bin directory. HPE recommends that you add the KeyView bin directory because the Lotus Notes installation might contain older KeyView OEM libraries.

Linux

1. Install Lotus Domino server. You do not need to configure the server.
2. Make sure the file `notes.ini` is in the `install/lotus/notes/latest/linux` directory, where `install` is the directory where Lotus Notes is installed. If the file does not exist, create an ASCII

file named `notes.ini`, and add the following text:

```
[Notes]
```

3. Add the `install/lotus/notes/latest/linux` directory to the PATH environment variable:

```
setenv PATH install/lotus/notes/latest/linux:$PATH
```

4. Add the `install/lotus/notes/latest/linux` and the KeyView bin directory to the LD_LIBRARY_PATH environment variable:

```
setenv LD_LIBRARY_PATH keyview_bin:install/lotus/notes/latest/linux:$LD_LIBRARY_PATH
```

where `keyview_bin` is the location of the KeyView bin directory. HPE recommends that you add the KeyView bin directory because the Lotus Notes installation might contain older KeyView OEM libraries.

Open Secured NSF Files

KeyView enables you to specify a user ID file and password to use to open a secured NSF file for extraction. See [Password Protected Files, on page 226](#) for more information.

Format Note Subfiles

The KeyView NSF reader uses XML templates to format note subfiles. You can customize the templates to approximate the look and feel of the original notes as closely as possible. For more information, see [Extract and Format Lotus Notes Subfiles, on page 163](#).

Extract Subfiles from PDF Files

KeyView can extract document-level and page-level attachments from a PDF document. Document-level attachments are added by using the **Attach A File** tool, and can include links to or from the parent document or to other file attachments. Page-level attachments are added as comments by using various tools. Page-level or comment attachments display the File Attachment icon or the Speaker icon on the page where they are located.

When a PDF file is extracted to disk, the PDF file is extracted to a directory and the PDF's attachments are saved in their native format to the same directory as the original PDF file.

Improve Performance for PDFs with Many Small Images

To improve performance when processing PDF files that contain many small images, you can choose to ignore images unless they exceed a minimum width and/or height. If an image is smaller than the minimum width or height, KeyView does not extract the image.

For example, to ignore images that are less than 16 pixels wide or less than 16 pixels in height, add the following to the `[pdf_flags]` section of the `formats.ini` file:

```
[pdf_flags]  
process_images_with_min_width=16  
process_images_with_min_height=16
```

Extract Embedded OLE Objects

The File Extraction API can extract embedded OLE objects from the following types of documents:

- Lotus Notes (DXL)
- Microsoft Excel
- Microsoft Word
- Microsoft PowerPoint
- Microsoft Outlook
- Microsoft Visio
- Microsoft Project
- OASIS Open Document
- Rich Text Format (RTF)

When an embedded OLE object is extracted from its parent file, the location of the embedded file in the original document is not available. The parent and child are extracted as separate files.

Extract Subfiles from ZIP Files

You can extract ZIP files that are not password-protected by using the general method (see [Extract Subfiles, on page 35](#)). However, some ZIP files use password protection, in which case you must use a different method to enter the required credentials. See [Password Protected Files, on page 226](#) for more information.

Default File Names for Extracted Subfiles

When a file name is not specified in the call to `extExtractSubFile`, in some cases, a default file name is applied to the extracted subfile.

Default File Name for Mail Formats

To avoid naming conflicts and problems with long file names, KeyView applies its own names to the extracted mail folders and mail items when a name is not supplied in the call to `extExtractSubFile`. A non-mail attachment retains its original file name and extension.

When the contents of a mail store or the message body of a mail message are extracted, the extracted file names might include the following:

- The first valid eight characters of the original folder name or "Subject" line of the mail message. If the "Subject" line is empty, the characters `kvext` are used, where `ext` is the format's extension. For example, the characters would be `"kvmsg"` for MSG, and `"kvnsf"` for NSF.

The following special characters are considered invalid and are ignored:

any non-printing character with a value less than 0x1F

angle brackets (< >)	double quotation mark (")
asterisk (*)	forward slash (/)
back slash (\)	pipe ()
colon (:)	question mark (?)

For notes, the file name is derived from the first 24 characters of the note text. For contact entries, the file name is derived from the full name of the contact.

- The characters `_kvn`, where `n` is an integer incremented from 0 for each extracted item.
- One of the following extensions:

Type	File Extension
email message	.mail
calendar appointment	.cal
contact entry	.cont
task entry	.task
note	.note
journal entry	.jrn1
distribution list	.dist
posting note	.post

If the type cannot be determined for an MSG or PST file, the file is given a `.mail` extension.

If the type cannot be determined for an NSF file, the file is given a `.tmp` extension.

For example, an MSG mail message with the subject line "RE: Product roadmap" that contains the Microsoft Excel attachment `release_schedule.xls` is extracted as

```
RE produ_kv0.mail
```

```
release_schedule.xls
```

If an extracted message contains an embedded OLE object or any attachment that does not have a name, the object or attachment is extracted as `_kv#.tmp`.

Default File Name for Embedded OLE Objects

KeyView can apply a default name to an extracted embedded OLE object when a name is not supplied in the call to `extExtractSubFile`. When an embedded OLE object is extracted, the extracted file name might include the following:

- The first valid eight characters of the main file. The following special characters are considered invalid and are ignored:

any non-printing character with a value less than 0x1F

angle brackets (< >) double quotation mark (")

asterisk (*) forward slash (/)

back slash (\) pipe (|)

colon (:) question mark (?)

- The characters `_kvn`, where `n` is an integer incremented from 0 for each extracted object.
- If KeyView can determine the embedded OLE is a Microsoft Office document, the original extension is used. If the file type cannot be determined, the file is given a `.tmp` extension.

For example, let us say a Microsoft Word document (`sales_quarterly.doc`) contains two embedded OLE objects: a Microsoft Excel file called `west_region.xls`, and a bitmap created in the Word document. The embedded objects would be extracted as

`sales_qu_kv0.xls`

`sales_qu_kv1.tmp`

Exclude Japanese Guide Text

This option prevents output of Japanese phonetic guide text when Microsoft Excel (`.xlsx`) files are processed.

To prevent output of Japanese phonetic guide text

- Set `NoPhoneticGuides` to `TRUE` in the `formats.ini` file:

```
[Options]
NoPhoneticGuides=TRUE
```

You can also enable this option programatically when filtering by passing `KVFLT_NOPHONETICGUIDES` to `fpFilterConfig`.

Chapter 4: Use the Filter API

This section describes how to perform some basic filtering tasks by using the Filter API.

• Generate an Error Log	55
• Extract Metadata	58
• Convert Character Sets	61
• Extract Tracked Deleted Text	63
• Filter PDF Files	64
• Filter Spreadsheet Files	69
• Filter HTML Files	72
• Filter XML Files	72
• Configure Headers and Footers	77
• Error Messages	77
• Tab Delimited Output for Embedded Tables	80
• Exclude Japanese Guide Text	81

Generate an Error Log

You can monitor and debug filtering operations by enabling a detailed error log. This allows you to see errors that are generated at run time and to track problem files in stream or file mode.

NOTE:

Error logs are not generated when in-process filtering is enabled.

The error log might include the following information:

- Generated error messages.
- Time stamp.
- Path and file name of the file in which the error occurred.
- Length of the file in which the error occurred. If the name of the original file or the name of the temporary file are not obtained in stream mode, the file length is reported.

The following is a sample log file:

```
-KVOOPE 12 # Time: 11:14:32 # File Len = 68140
-KVOOPE 13 # Time: 11:23:05 # H:\files\WP\Word97\fnldmsa.doc
-KVOOPE 5 # Time: 12:15:54 # H:\files\SS\XL2000\corporate.xls
-KVOOPE 5 # Time: 12:45:19 # H:\files\WP\WPerf5\wp501.doc
-KVOOPE 12 # Time: 14:25:33 # H:\files\PG\PPoint95\95.ppt
-KVOOPE 26 # Time: 16:26:04 # File Len = 19117568
-KVOOPE 10 # Time: 20:27:40 # File Len = 19117568
```

You can specify the information that is written to the log file using either the API or environment variables. To configure a log file for a single filtering session, use environment variables. To configure a log file for all

filtering sessions, use the API. Configuring the log file using the API overrides the same settings in the environment variables. You can also specify additional settings in the `formats.ini` file

You can configure the following features of the log file:

- Enable or disable logging. See [Enable or Disable Error Logging, below](#).
- Change the default path and file name of the log file. See [Change the Path and File Name of the Log File, below](#).
- Include memory errors in the log file. See [Report Memory Errors, on the next page](#).
- Specify a memory guard that is used to generate memory overwrite errors in the log. See [Specify a Memory Guard, on the next page](#).
- Include the input file name in the log file when filtering a stream. See [Report the File Name in Stream Mode, on the next page](#).
- Specify the maximum size of the log file. See [Specify the Maximum Size of the Log File, on page 58](#).

Enable or Disable Error Logging

You can enable or disable error logging using either the API or environment variables. By default, a file called `kvoop.log` is created in the system temporary directory; however, you can change the path and file name of this file (see [Change the Path and File Name of the Log File, below](#)).

Use the API

To enable or disable logging in the API, instantiate the `Filter` object using the constructor `Filter(java.lang.String outputCharSet, long filterFlags)`, and set the `filterFlags` argument to either `FILTERFLAG_OOPLOGON` or `FILTERFLAG_OOPLOGOFF`. For example:

```
objFilter = new Filter(outputCharSet, Filter.FILTERFLAG_OOPLOGON);
```

Use Environment Variables

To enable logging, add the environment variable `KVOOPLOGON`, and set the variable value to 1. To disable logging, do not set the environment variable `KVOOPLOGON`.

Change the Path and File Name of the Log File

You can change the default path and file name of the log file. The default is `C:\temp\kvoop.log` on Windows and `/tmp/kvoop.log` on UNIX.

To change the path and file name of the log file, add the following to the `formats.ini` file:

```
[kvooplog]
KvoopLogName=filepath
```

The `formats.ini` file is in the directory `install\OS\bin`, where `install` is the path name of the Filter installation directory and `OS` is the name of the operating system.

Report Memory Errors

You can report memory leaks and memory overwrites in the log file by enabling the memory trace system, either by using the API or environment variables. If the memory trace system is enabled, the error messages for memory leaks and memory overwrites (KError_MemoryLeak and KError_MemoryOverwrite, respectively) are reported in the log file when they are generated. The error messages are listed in [Error Messages, on page 77](#).

NOTE:

To report memory overwrites, you must also set a memory guard. See [Specify a Memory Guard, below](#).

Use the API

To enable or disable the memory trace system in the API, instantiate the Filter object using the constructor `Filter(java.lang.String outputCharSet, long filterFlags)`, and set the `filterFlags` argument to either `FILTERFLAG_OOPMEMTRACEON` or `FILTERFLAG_OOPMEMTRACEOFF`. For example:

```
objFilter = new Filter(outputCharSet, Filter.FILTERFLAG_OOPMEMTRACE);
```

Use Environment Variables

To enable the memory trace system, add the `KVOOPMT` environment variable, and set its value to 1. To disable the memory trace system, do not set the `KVOOPMT` environment variable.

Specify a Memory Guard

To report memory overwrites in the log file, you must set a memory guard that protects against memory overwrites. Normally, this is set in the range of 100-200 bytes. For example, if a memory guard of 100 is set and 20 bytes of memory are specified, a total of 120 bytes of memory are allocated. The additional memory is used to monitor and identify memory overwrites.

To configure the memory guard, add the following section to the `formats.ini` file:

```
[Kvooplog]
mg=100
```

Report the File Name in Stream Mode

When you run Filter in file mode, the file name is always reported in the log file. To report the file name in stream mode, you must extract it through the API.

To add the input file name to the log

1. Create an instance of `ConfigOption` with the following arguments:
 - a. Set the `OptionType` to `CFG_SET00PSRCFILE`.
 - b. Set the `OptionValue` to 0.
 - c. Set `OptionData` to the `input_filename`.
2. Call the `setConfigOption` method, and pass in the `ConfigOption` instance.

Example

```
if((filterFlags & Filter.FILTERFLAG_OOPLOGON) == Filter.FILTERFLAG_OOPLOGON)
{
    ConfigOption config = new ConfigOption(Filter.CFG_SET00PSRCFILE, 0, inFile);
    objFilter.setConfigOption(config);
}
```

Specify the Maximum Size of the Log File

You can specify the maximum size of the log file. When this size is reached and new entries are logged, either the first entry in the file is overwritten or the new entries are not reported.

To configure the maximum log size and whether old entries are overwritten, add the following section to the `formats.ini` file:

```
[Kvooplog]
LogFileSize=10
OverWriteLog=1
```

Option	Description
LogFileSize	This option specifies the maximum size of the log file in KB. The minimum is 1 K. If a size is not specified, the default 2 MB is used.
OverWriteLog	This option determines whether the log file is overwritten when the maximum log file size (<code>LogFileSize</code>) is reached. If you set this option to 1, the first entry in the log file is overwritten. If you set this option to 0, new entries are not reported in the log file.

Extract Metadata

When a file format supports metadata, KeyView can extract and process that information. Metadata includes document information fields such as title, author, creation date, and file size. Depending on the file's format, metadata is referred to in a number of ways: for example, "summary information," "OLE summary information," "file information," and "document properties."

The metadata in mail formats (MSG and EML) and mail stores (PST, NSF, and MBX) is extracted differently than other formats. For information on extracting metadata from these formats, see [Extract Mail Metadata, on page 38](#).

NOTE:

KeyView can only extract metadata from a document if metadata is defined in the document, and if the document reader can extract metadata for the file format. The section [Supported Formats, on page 94](#) lists the file formats for which metadata can be extracted. KeyView does not generate metadata automatically from the document contents.

The sample program `FilterTest` demonstrates how to extract metadata. See [Sample Programs, on page 82](#).

Extract Metadata for File Filtering

To extract metadata for file filtering

1. Optionally, set the input source using the `setInputSource(java.lang.String inFile)` method of the `Filter` object.
2. If the input source was set in step 1, call the `getSummaryInfo()` method of the `Filter` object to retrieve an object of the `SummaryInfo` class. Otherwise, call the `getSummaryInfo(java.lang.String inFile)` method.
3. Use the methods of the `SummaryInfo` object to retrieve the metadata information.

Extract Metadata for Stream Filtering

To extract metadata for stream filtering

1. Optionally, set the input source using the `setInputSource(java.io.InputStream input)` method of the `Filter` object.
2. If the input source was set in step 1, call the `getSummaryInfo()` method of the `Filter` object to retrieve an object of the `SummaryInfo` class. Otherwise, call the `getSummaryInfo(java.io.InputStream in)` method.
3. Use the methods of the `SummaryInfo` object to retrieve the metadata information.

Example

Below is an example of a call to `getSummaryInfo()`:

```
SummaryInfo[] sinfo = objFilter.getSummaryInfo();
if(sinfo != null)
{
    System.out.println("\nSummary info has been extracted.");
    fos_sum = new FileOutputStream(summaryOutFile);
    DataOutputStream dos_sum = new DataOutputStream(fos_sum);
    for(int i=0; i<sinfo.length; i++)
    {
        if(sinfo[i].getElementName() != null)
        {
            dos_sum.writeBytes("Element name: " + sinfo[i].getElementName() + "\n");
            dos_sum.writeBytes("Element type: " + sinfo[i].getSumInfoType() + "\n");
        }
    }
}
```

```

        if(sinfo[i].getIsValid() == true)
        {
            if(sinfo[i].isDateTimeType())
            {
                dos_sum.writeBytes("Date/time: ");
                dos_sum.writeBytes(sinfo[i].getDateTime());
            }
            else
            {
                byte[] data = sinfo[i].getData();
                if(data != null)
                {
                    dos_sum.writeBytes("Element data: ");
                    dos_sum.write(data);
                }
            }
        }
        dos_sum.writeBytes("\n\n");
    }
}
dos_sum.close();
fos_sum.close;
}
sinfo=null;

```

The `SummaryInfo` class stores the metadata extraction results. After calling the `Filter.getSummaryInfo()` method, call the get methods provided by each instance of this class to extract metadata:

`getElementName` Gets the name of the metadata element.

`()`

`getSumInfoType` Specifies the data type of the metadata element. The possible types are:

`()`

- `KV_String`—The value in the metadata field is a string.
- `KV_Int4`—The value in the metadata field is an integer.
- `KV_DateTime`—The value in the metadata field is a date and time. This type is a 64-bit value representing the number of 100-nanosecond intervals since January 1, 1601 (Windows FILETIME EPOCH). You might need to convert this value into another format.

You can also use the `isDateTimeType()` method to determine whether a metadata element is of date/time type, and then use the `getDateTime()` method to obtain the date/time in the form of a string.

- `KV_ClipBoard`—Currently not supported.
- `KV_Bool`—The value in the metadata field is a boolean.
- `KV_Unicode`—The value in the metadata field is a Unicode string.
- `KV_IEEE8`—The value in the metadata field is an IEEE 8-byte integer.
- `KV_Other`—The value in the metadata field is user-defined.

<code>getIsValid()</code>	Specifies whether the data value is present in the document. <code>true</code> specifies that the value is valid. For example, if the “Title” element was not populated in the document, <code>getIsValid</code> would return <code>false</code> .
<code>isDateTimeType()</code>	Determines whether the metadata element is of date/time type.
<code>getDateTime()</code>	Gets the date and time in the form of a string. If the metadata element is of <code>KV_DateTime</code> type, call this method to get the date and time in the form of a string, for example “Wed Jun 30 21:49:08 1993” or “135 Minutes”.
<code>getData()</code>	Gets the content of the element. If type is <code>KV_Int4</code> or <code>KV_Bool</code> , data contains the actual value. Otherwise, data is a pointer to the actual value. <code>KV_DateTime</code> and <code>KV_IEEE8</code> point to an 8-byte value. <code>KV_String</code> and <code>KV_Unicode</code> point to the beginning of the string that contains the text. <code>KV_Unicode</code> is replaced with <code>KV_String</code> when the UNICODE value has been character mapped to the desired output character set.

Convert Character Sets

Filter can convert the character set of a source document to an arbitrary character set specified in the API, or to the character set of the operating system on which the output text is viewed. For this conversion to occur, a source character set *must* be identified. The source character set can either be determined by the document reader, or can be set in the API. The section [Supported Formats, on page 94](#) lists file formats for which character set information can be determined by the document reader. The character sets are defined as constants in the Filter class.

Determine the Character Set of the Output Text

To determine the output character set of a filtered document, Filter considers the following:

- Whether the document reader can determine the character set of the file format. If the document reader cannot determine the character set information for the document type, set the source character set in the API.
- Whether the *source* character set is specified in the API.
- Whether the *target* character set is specified in the API.

Guidelines for Character Set Conversion

Below are some rules for the determination of character set mapping:

- If the source is not determined by the document reader or configured in the API, then the character set of the output text is always unknown, regardless of the target character set configuration. The document cannot be converted to a target character set or the operating system’s code page unless the source character set is known.

- If the target character set is *not* specified in the API, and the source character set is identified, then the operating system's code page is used for the output text.
- If the source character set is identified, and the target character set is specified in the API, then the target character set specified in the API is used for the output text.
- For documents that contain multiple character sets, HPE recommends that the target character set be forced to UNICODE or UTF-8.

The following table illustrates how Filter determines the character set of the output text.

Determining the Output Character Set—Example

Source charset read by Filter	Source charset specified in API	Target charset specified in API	Output charset
No	No	No	no conversion
No	KVCS_936	No	OS code page
No	No	UNICODE	no conversion
No	KVCS_936	UNICODE	UNICODE
Yes	No	No	OS code page
Yes	KVCS_936	No	OS code page
Yes	No	UNICODE	UNICODE
Yes	KVCS_936	UNICODE	UNICODE

Set the Character Set During Filtering

You can convert the character set of a file at the time the file is filtered.

To specify the source character set, use the `setSourceCharSet(java.lang.String charset)` method. For example:

```
objFilter.setSourceCharSet(sourceCharSet);
```

To specify the target character set, instantiate the Filter object using the constructor `Filter(java.lang.String outputCharSet, long filterFlags)`. For example:

```
objFilter = new Filter(outputCharSet, filterFlags);
```

Set the Character Set During Subfile Extraction

You can convert the character set of a subfile at the time the subfile is extracted from the container and before it is filtered. This is most often used to set the character set of a mail message's body text. See [Filter PDF Files, on page 64](#) for more information.

To specify the source and target character set of a subfile

1. Use the methods of the `ExtSubFileExtractConfig` object to set the source and target character set.
2. Call the `extExtractSubFile` method of the `Filter` object and pass in the `ExtSubFileExtractConfig` object. For example:

```
extconfig = new ExtSubFileExtractConfig();
extconfig.setSourceCharset(m_sourceCharset);
extconfig.setTargetCharset(m_outputCharset);
extinfo = m_objFilter.extExtractSubFile(extContextID, i, extconfig);
```

Prevent the Default Conversion of a Character Set

You can prevent the default conversion of text to the operating system code page, and specify that `Filter` retain the original character encoding of the document when it is available. Any document identified as containing more than one character encoding is converted to the first encoding encountered in the file.

To prevent the default conversion, instantiate the `Filter` object using the constructor `Filter` (`java.lang.String outputCharset, long filterFlags`), and set the `filterFlags` argument to `FILTERFLAG_NODEFAULTCHARSETCONVERT`. For example:

```
objFilter = new Filter(outputCharset, Filter.FILTERFLAG_NODEFAULTCHARSETCONVERT);
```

This setting overrides the source or target character set specified in the API.

Extract Tracked Deleted Text

The revision tracking feature in applications—such as Microsoft Word's **Track Changes**—marks changes to a document (typically, strikethrough for deleted text and underline for inserted text) and tracks each change by reviewer name and date. If revision tracking was enabled when text was deleted from a source document, you can configure `Filter` to extract the deleted text. `Filter` does not extract the reviewer name and revision date. Deleted text is excluded from the filtered output by default.

To extract deleted text from a document and include it in the filtered output, call the `includeRevisionMark` method. For example:

```
if(inclRevisionMark == true)
{
    objFilter.includeRevisionMark();
}
```

To reset the flag and exclude deleted text from the filtered output, call the `excludeRevisionMark` method. For example:

```
if(inclRevisionMark == false)
{
    objFilter.excludeRevisionMark();
}
```

Filter PDF Files

Filter has special configuration options that allow greater control over the conversion of Adobe Acrobat PDF files.

Filter PDF Files to a Logical Reading Order

The PDF format is primarily designed for presentation and printing of brochures, magazines, forms, reports, and other materials with complex visual designs. Most PDF files do not contain the *logical structure* of the original document—the correct reading order, for example, and the presence and meaning of significant elements such as headers, footers, columns, tables, and so on.

KeyView can filter a PDF file either by using the file’s internal unstructured paragraph flow, or by applying a structure to the paragraphs to reproduce the logical reading order of the visual page. Logical reading order enables KeyView to output PDF files that contain languages that read from right-to-left (such as Hebrew and Arabic) in the correct reading direction.

NOTE:
The algorithm used to reproduce the reading order of a PDF page is based on common page layouts. The paragraph flow generated for PDFs with unique or complex page designs might not emulate the original reading order exactly.

For example, page design elements such as drop caps, callouts that cross column boundaries, and significant changes in font size might disrupt the logical flow of the output text.

By default, KeyView produces an *unstructured* text stream for PDF files. This means that PDF paragraphs are extracted in the order in which they are stored in the file, not the order in which they appear on the visual page. For example, a three-column article could be output with the headers and title at the end of the output file, and the second column extracted before the first column. Although this output does not represent a logical reading order, it accurately reflects the internal structure of the PDF.

You can configure KeyView to produce a *structured* text stream that flows in a specified direction. This means that PDF paragraphs are extracted in the order (logical reading order) and direction (left-to-right or right-to-left) in which they appear on the page.

The following paragraph direction options are available:

Paragraph Direction Option	Description
Left-to-right	Paragraphs flow logically and read from left to right. You should specify this option when most of your documents are in a language that uses a left-to-right reading order, such as English or German.
Right-to-left	Paragraphs flow logically and read from right to left. You should specify this option when most of your documents are in a language that uses a right-to-left reading order, such as Hebrew or Arabic.
Dynamic	Paragraphs flow logically. The PDF filter determines the paragraph direction for each PDF page, and then sets the direction

Paragraph Direction Option	Description
	accordingly. Filter uses this option when a paragraph direction is not specified.

NOTE:

Filtering might be slower when logical reading order is enabled. For optimal speed, use an unstructured paragraph flow.

The paragraph direction options control the direction of paragraphs on a page; they do not control the text direction in a paragraph. For example, a PDF file might contain English paragraphs in three columns that read from left to right, but 80% of the second paragraph might contain Hebrew characters. If the left-to-right logical reading order is enabled, the paragraphs are ordered logically in the output—title paragraph, then paragraph 1, 2, 3, and so on—and flow from the top left of the first column to the bottom right of the third column. However, the *text* direction of the second paragraph is determined independently of the page by the PDF filter, and is output from right to left.

NOTE:

Extraction of metadata is not affected by the paragraph direction setting. The characters and words in metadata fields are extracted in the correct reading direction regardless of whether logical reading order is enabled.

Enable Logical Reading Order

You can enable logical reading order by using either the API or the `formats.ini` file. Setting the paragraph direction in the API overrides the setting in the `formats.ini` file.

Use the Java API

To enable PDF logical reading order in the API, use the `setPDFLogicalOrder(int orderFlag)` method, and set the `orderFlag` argument to one of the following flags:

Flag	Description
PDF_LOGICAL_ORDER_LTR	Logical reading order and left-to-right paragraph direction
PDF_LOGICAL_ORDER_RTL	Logical reading order and right-to-left paragraph direction
PDF_LOGICAL_ORDER_AUTO	Logical reading order. The PDF reader determines the paragraph direction for each PDF page, and then sets the direction accordingly. Filter uses this option when a paragraph direction is not specified.
PDF_LOGICAL_ORDER_RAW	Unstructured paragraph flow. This is the default behavior. If logical reading order is enabled, and you want to return to an unstructured paragraph flow, set this flag.

For example:

```
objFilter.setPDFLogicalOrder(Filter.PDF_LOGICAL_ORDER_RTL);
```

The `FilterTest` sample program demonstrates this method. See [FilterTest](#), on page 84.

Use the `formats.ini` File

The `formats.ini` file is in the directory `install\OS\bin`, where `install` is the path name of the Filter installation directory and `OS` is the name of the operating system.

To enable logical reading order by using the `formats.ini` file

1. Change the PDF reader entry in the `[Formats]` section of the `formats.ini` file as follows:

```
[Formats]
200=1pdf
```

2. Optionally, add the following section to the end of the `formats.ini` file:

```
[pdf_flags]
pdf_direction=paragraph_direction
```

where `paragraph_direction` is one of the following:

Flag	Description
LPDF_LTR	Left-to-right paragraph direction
LPDF_RTL	Right-to-left paragraph direction
LPDF_AUTO	The PDF filter determines the paragraph direction for each PDF page, and then sets the direction accordingly. Filter uses this option when a paragraph direction is not specified.
LPDF_RAW	Unstructured paragraph flow. This is the default behavior. If logical reading order is enabled, and you want to return to an unstructured paragraph flow, set this flag.

Rotated Text

When a PDF that contains rotated text is filtered, the rotated text is extracted after the text at the end of the PDF page on which the rotated text appears. If the PDF is filtered with logical order enabled, and the amount of rotated text on a page surpasses a predefined threshold, the page is automatically output as an unstructured text stream. You cannot configure this threshold.

Extract Custom Metadata from PDF Files

To extract custom metadata from your PDF files, add the custom metadata names to the `pdfsr.ini` file provided, and copy the modified file to the `bin` directory. You can then extract metadata as you normally would.

The `pdfsr.ini` is in the directory `samples\pdfini`, and has the following structure:

```

<META>
<TOTAL>total_item_number</TOTAL>,
/metadata_tag_name datatype,
</META>

```

Parameter	Description
total item number	The total number of metadata tags that are listed.
metadata_tag_name	The metadata tag name used in the PDF files.
datatype	<p>The data type of the metadata element. The possible types are:</p> <ul style="list-style-type: none"> • KV_String • KV_Int4 • KV_DateTime • KV_ClipBoard • KV_Bool • KV_Unicode • KV_IEEE8 • KV_Other

For example:

```

<META>
<TOTAL>4</TOTAL>
/part_number      INT4
/volume           INT4
/purchase_date    DATETIME
/customer         STRING
</META>

```

Skip Embedded Fonts

Text in PDF files sometimes contain embedded fonts. If you experience difficulties filtering embedded fonts, there are options in the API, the `formats.ini` file, and the `FilterTest` sample program that you can set to skip this type of text.

NOTE:

If you choose to skip embedded fonts, none of the content that contains embedded fonts is included in the output.

Use the formats.ini File

To skip embedded fonts using the formats.ini file

- Set the following parameters:

```
[pdf_flags]
skipembeddedfont=TRUE
embedded_font_threshold=threshold
```

where `threshold` is a value between 0 and 100. A threshold of 100 skips all embedded font text; a threshold of 0 retains all embedded font text. Set `skipembeddedfont` to `TRUE` to enable the `embedded_font_threshold` parameter.

The default value of `embedded_font_threshold` is 100. If you set `skipembeddedfont` to `true` and do not specify the `embedded_font_threshold` parameter, Filter skips all embedded text.

When you use `formats.ini` to skip embedded fonts, you can also specify an *embedded font threshold*, which is an arbitrary percentage probability that the glyph in the embedded text maps to a character value in the output character set (ASCII, UTF-8, and so on).

For example, if you specify a threshold of 75, embedded text glyphs that have a 75% or greater probability of correctly matching the character in the output character set are included in the output; glyphs that have a probability of less than 75% of matching the output character set are omitted from the output.

Use the Java API

To skip embedded fonts using the Java API, set the `setSkipEmbeddedFont(boolean)` method to `true`. For example:

```
objFilter.setSkipEmbeddedFont(true);
```

The `FilterTest` sample program demonstrates this method. See [FilterTest](#), on page 84.

Control Hyphenation

There are two types of hyphens in a PDF document:

- A *soft hyphen* is added to a word by a word processor to divide the word across two lines. This is a discretionary hyphen and is used to ensure proper text flow in justified text.
- A *hard hyphen* is intentionally added to a word regardless of the word's position in the text flow. It is required by the rules of grammar or word usage. For example, compound words (such as *three-week vacation* and *self-confident*) contain hard hyphens.

By default, KeyView skips the source document's soft hyphens in the Filter output to provide more searchable text content. However, if you want to maintain the document layout, you can keep soft hyphens in the Filter output. To keep soft hyphens, you must enable the soft hyphen flag in `formats.ini` or in the API.

Use the formats.ini File

To keep soft hyphens by using the `formats.ini` file, set the following parameter:

```
[pdf_flags]
keepsofthyphen=TRUE
```

Use the Java API

To keep soft hyphens using the Java API, set the `setKeepSoftHyphen(boolean)` method to `true`. For example:

```
objFilter.setKeepSoftHyphen(true);
```

The `FilterTest` sample program demonstrates this method. See [FilterTest](#), on page 84.

Filter Spreadsheet Files

Filter has special configuration options that allow greater control over the conversion of spreadsheet files.

Filter Worksheet Names

Normally, Filter does not extract worksheet names from a spreadsheet because it is assumed the text should not be exposed. You can change this default behavior, and extract worksheet names by adding the following lines to the `formats.ini` file:

```
[Options]
getsheetnames=1
```

Filter Hidden Text in Microsoft Excel Files

Normally, Filter does not filter hidden text from a Microsoft Excel spreadsheet because it is assumed the text should not be exposed. You can change this default behavior, and extract text from hidden rows, columns, and sheets from Excel spreadsheets by adding the following lines to the `formats.ini` file:

```
[Options]
gethiddeninfo=1
```

Specify Date and Time Format on UNIX Systems

In Microsoft Excel you can choose to format dates and times according to the system locale. On Windows, KeyView uses the system locale settings to determine how these dates and times should be formatted. In other operating systems, KeyView uses the U.S. short date format (*mm/dd/yyyy*). You can change this by specifying the formats you wish to use in the `formats.ini` file.

To specify the system date and time format on UNIX systems

- In the `formats.ini` file, specify the following options:
 - `SysDateTime`. The format to use when a cell is formatted using the system format including both the date and the time.
 - `SysLongDate`. The format to use when a cell is formatted using the system long date format.
 - `SysShortDate`. The format to use when a cell is formatted using the system short date format.
 - `SysTime`. The format to use when a cell is formatted using the system time format.

NOTE:

These values cannot contain spaces.

For example, if you specify `SysDateTime=%d/%m/%Y`, dates and times are extracted in the following format:

`28/02/2008`

The format arguments are the same as those for the `strftime()` function. Refer to the following webpage for more information.

<http://linux.die.net/man/3/strftime>

Filter Very Large Numbers in Spreadsheet Cells to Precision Numbers

Numbers in Microsoft Excel files can now be extracted and written to the output without formatting. By default, numbers are extracted in the format specified by the Excel file (for example, *General*, *Currency* and *Date*). Spreadsheets might contain cells that have very large numbers in them. Excel displays the numbers in a scientific notation that rounds or truncates the numbers.

To extract numbers without formatting, add the following options in the `formats.ini` file:

```
[Options]
```

```
ignoredefnumformats=1
```

Extract Microsoft Excel Formulas

Normally, the actual value of a formula is extracted from an Excel spreadsheet; the formula from which the value is derived is not included in the output. However, KeyView enables you to include the value as well as the formula in the output. For example, if Filter is configured to extract the formula and the formula value, the output might look like this:

```
245 = SUM(B21:B26)
```

The calculated value from the cell is 245 and the formula from which the value is derived is SUM (B21:B26).

NOTE:

Depending on the complexity of the formulas, enabling formula extraction might result in slightly slower performance.

To set the extraction option for formulas, add the following lines to the `formats.ini` file:

```
[Options]
getformulastring=option
```

where `option` is one of the following:

Option	Description
0	Extract the formula value only. This is the default. If formula extraction is enabled, and you want to return to the default, set this option.
1	Extract the formula only.
2	Extract the formula and the formula value.

If a function in a formula is not supported or is invalid, and option 1 or 2 is specified, only the calculated value is extracted. See [Extract Microsoft Excel Formulas, on the previous page](#) for a list of supported functions.

When formula extraction is enabled, Filter can extract Microsoft Excel formulas that contain the functions listed in [Extract Microsoft Excel Formulas, on the previous page](#).

Supported Microsoft Excel Functions

=ABS()	=ACOS()	=AND()	=AREAS()
=ASIN()	=ATAN2()	=ATAN2()	=AVERAGE()
=CELL()	=CHAR()	=CHOOSE()	=CLEAN()
=CODE()	=COLUMN()	=COLUMNS()	=CONCATENATE()
=COS()	=COUNT()	=COUNTA()	=DATE()
=DATEVALUE()	=DAVERAGE()	=DAY()	=DCOUNT()
=DDB()	=DMAX()	=DMIN()	=DOLLAR()
=DSTDEV()	=DSUM()	=DVAR()	=EXACT()
=EXP()	=FACT()	=FALSE()	=FIND()
=FIXED()	=FV()	=GROWTH()	=HLOOKUP()
=HOUR()	=ISBLANK()	=IF()	=INDEX()
=INDIRECT()	=INT()	=IPMT()	=IRR()
=ISERR()	=ISERROR()	=ISNA()	=ISNUMBER()
=ISREF()	=ISTEXT()	=LEFT()	=LEN()
=LINEST()	=LN()	=LOG()	=LOG10()
=LOGEST()	=LOOKUP()	=LOWER()	=MATCH()

=MAX()	=MDETERM()	=MID()	=MIN()
=MINUTE()	=MINVERSE()	=MIRR()	=MMULT()
=MOD()	=MONTH()	=N()	=NA()
=NOT()	=NOW()	=NPER()	=NPV()
=OFFSET()	=OR()	=PI()	=PMT()
=PPMT()	=PRODUCT()	=PROPER()	=PV()
=RATE()	=REPLACE()	=REPT()	=RIGHT()
=ROUND()	=ROUND()	=ROW()	=ROWS()
=SEARCH()	=SECOND()	=SIGN()	=SIN()
=SLN()	=SQRT()	=STDEV()	=SUBSTITUTE()
=SUM()	=SYD()	=T()	=TAN()
=TEXT()	=TIME()	=TIMEVALUE()	=TODAY()
=TRANSPOSE()	=TREND()	=TRIM()	=TRUE()
=TYPE()	=UPPER()	=VALUE()	=VAR()
=VLOOKUP()	=WEEKDAY()	=YEAR()	

Filter HTML Files

KeyView can filter comments from HTML documents. To enable comment filtering, you must set a flag in the `formats.ini` file.

The `formats.ini` file is in the `install\OS\bin` directory, where `install` is the Filter installation directory and `OS` is the name of the operating system.

To enable filtering of comments from HTML files

1. Open the `formats.ini` file in a text editor.
2. Under `[Options]`, set the following flag.

```
GetHTMLHiddenInfo=1
```

Filter XML Files

Filter SDK enables you to extract all or selected content from source XML files. You can specify the elements and attributes extracted from a document using the API or an INI file (see [Configure Element Extraction for XML Documents, on the next page](#)). Filter detects the following XML formats:

- generic XML
- Microsoft Office 2003 XML (Word, Excel, and Visio)

- StarOffice/OpenOffice XML (text document, presentation, and spreadsheet)

See [File Format Detection, on page 176](#) for more information on format detection.

Configure Element Extraction for XML Documents

When filtering XML files, you can specify which elements and attributes are extracted according to the file's format ID or *root element*. This is useful when you want to extract only relevant text elements, such as abstracts from reports, or a list of authors from an anthology.

A root element is an element in which all other elements are contained. In the XML sample below, `book` is the root element:

```
<book>
  <title>XML Introduction</title>
  <product id="33-657" status="draft">XML Tutorial</product>
  <chapter>Introduction to XML
    <para>What is HTML</para>
    <para>What is XML</para>
  </chapter>
  <chapter>XML Syntax
    <para>Elements must have a closing tag</para>
    <para>Elements must be properly nested</para>
  </chapter>
</book>
```

For example, you could specify that when filtering files with the root element `book`, the element `title` is extracted as metadata, and only `product` elements with a `status` attribute value of `draft` are extracted. When you extract an element, the child elements within the element are also extracted. For example, if you extract the element `chapter` from the sample above, the child element `para` is also extracted.

Filter SDK defines default element extraction settings for the following XML formats:

- generic XML
- Microsoft Office 2003 XML (Word, Excel, and Visio)
- StarOffice/OpenOffice XML (text document, presentation, and spreadsheet)

These settings are defined internally and are used when filtering these file formats; however, you can modify their values.

In addition to the default extraction settings, you can also add custom settings for your own XML document types. If you do not define custom settings for your own XML document types, the settings for the generic XML are used.

Modify Element Extraction Settings

You can modify configuration settings for XML documents through either the API or the `kvxconfig.ini` file.

Use the Java API

You can use the Java API to modify the settings for the standard XML document types or add configuration settings for your own XML document types.

To modify settings

1. Declare an array of XMLConfigSet objects.
2. Create an instance of ConfigOption with the following arguments:
 - a. Set the OptionType to CFG_SETXMLCONFIGINFO.
 - b. Set the OptionValue to 0.
 - c. Set OptionData to the array object.
3. Call the setConfigOption method, and pass in the ConfigOption instance.
4. Call a filter method. For example:

```
XMLConfigSet[] XMLInfo;
ConfigOption config=new ConfigOption(Filter.CFG_SETXMLCONFIGINFO, 0, XMLInfo);
objFilter.setConfigOption(config);
```

Use an Initialization File

You can use the initialization file to modify the settings for the standard XML document types or add configuration settings for your own XML document types.

To modify settings

1. Modify the kvxconfig.ini file.
2. Use the initialization file when processing the XML file. See [Modify Element Extraction Settings in the kvxconfig.ini File](#), below.

The Java sample program FilterTest demonstrates how to use the initialization file in the filtering process. See [Sample Programs](#), on page 82.

Modify Element Extraction Settings in the kvxconfig.ini File

The kvxconfig.ini file contains default element extraction settings for supported XML formats. The file is in the directory install\OS\bin, where install is the path name of the Filter installation directory and OS is the name of the operating system. For example, the following entry defines extraction settings for the Microsoft Visio 2003 XML format:

```
[config3]
eKVFormat=MS_Visio_XML_Fmt
szRoot=
szInMetaElement=DocumentProperties
szExMetaElement=PreviewPicture
szInContentElement=Text
szExContentElement=
szInAttribute=
```

The following options are available:

Configuration Option	Description
eKVFormat	<p>The format ID as detected by the KeyView detection module. This determines the file type to which these extraction settings apply. See File Format Detection, on page 176 for more information on format ID values.</p> <p>If you are adding configuration settings for a custom XML document type, this is not defined.</p>
szRoot	<p>The file's root element. When the format ID is not defined, the root element is used to determine the file type to which these settings apply.</p> <p>To further qualify the element, specify its namespace. See Specify an Element's Namespace and Attribute, on the next page.</p>
szInMetaElement	<p>The elements extracted from the file as metadata. All other elements are extracted as text.</p> <p>Separate multiple entries with commas. To further qualify the element, specify its namespace, its attributes, or both. See Specify an Element's Namespace and Attribute, on the next page.</p>
szExMetaElement	<p>The child elements in the included metadata elements that are not extracted from the file as metadata. For example, the default extraction settings for the Visio XML format extract the DocumentProperties element as metadata. This element includes child elements such as Title, Subject, Author, Description, and so on. However, the child element PreviewPicture is defined in szExMetaElement because it is binary data and should not be extracted.</p> <p>You cannot exclude any metadata elements from the output for StarOffice files. All metadata is extracted regardless of this setting.</p> <p>Separate multiple entries with commas. To further qualify the element, specify its namespace, its attributes, or both. See Specify an Element's Namespace and Attribute, on the next page.</p>
szInContentElement	<p>The elements extracted from the file as content text. Enter an asterisk (*) to extract all elements including child elements.</p> <p>Separate multiple entries with commas. To further qualify the element, specify its namespace, its attributes, or both. See Specify an Element's Namespace and Attribute, on the next page.</p>
szExContentElement	<p>The child elements in the included content elements that are not extracted from the file as content text.</p>

Configuration Option	Description
	Separate multiple entries with commas. To further qualify the element, specify its namespace, its attributes, or both. See Specify an Element's Namespace and Attribute , below.
szInAttribute	<p>The attribute values extracted from the file. If attributes are not defined here, attribute values are not extracted.</p> <p>Enter the namespace (if used), element name, and attribute name in the following format:</p> <p>namespace:elementname@attributename</p> <p>For example:</p> <p>autonomy:division@name</p> <p>Separate multiple entries with commas.</p>

Specify an Element's Namespace and Attribute

To further qualify an element, you can specify that the element exist in a certain namespace and/or contain a specific attribute. To define the namespace *and* attribute of an element, enter the following:

```
ns_prefix:elemname@attribname=attribvalue
```

NOTE:

You must enclose attribute values that contain spaces in quotation marks.

For example, the entry `bg:language@id=xml` extracts a `language` element in the namespace `bg` that contains the attribute name `id` with the value of `"xml"`. This entry extracts the following element from an XML file:

```
<bg:language id="xml">XML is a simple, flexible text format derived from
SGML</bg:language>
```

but does not extract:

```
<bg:language id="sgml">SGML is a system for defining markup
languages.</bg:language>
```

or

```
<adv:language id="xml">The namespace should be a Uniform Resource Identifier
(URI).</adv:language>
```

Add Configuration Settings for Custom XML Document Types

You can define element extraction settings for custom XML document types by adding the settings to the `kvxconfig.ini` file. For example, for files that contain the root element `autxml`, you can add the following section to the end of the initialization file:

```
[config101]
eKVFormat=
```

```

szRoot=autxml
szInMetaElement=dc:title,dc:meta@title,dc:meta@name=title
szExMetaElement=

szInContentElement=aut:division@name=keyview,aut:division@name=idol,p@style="Heading 1"
szExContentElement=
szInAttribute=aut:division@name

```

The custom extraction settings must be preceded by a section heading named `[configN]`, where `N` is an integer starting at 100 and increasing by 1 for each additional file type, as in `[config100]`, `[config101]`, `[config102]`, and so on. The default extraction settings for the supported XML formats are numbered `config0` to `config99`. Currently only 0 to 6 are used.

Since a custom XML document type is not recognized by the KeyView detection module, the format ID is not defined. The file type is identified by the file's root element only.

If a custom XML document type is not defined in the `kvxconfig.ini` file or by the `setConfigOption` method, then the default extraction settings for a generic XML document are used.

Configure Headers and Footers

You can configure custom header and footer tags for word processing and spreadsheet documents by editing the `formats.ini` file.

To configure headers and footers

1. Open the `formats.ini` file.
2. In the `[Options]` section, add the following items:

```

header_start_tag=HeaderStart
header_end_tag=HeaderEnd
footer_start_tag=FooterStart
footer_end_tag=FooterEnd

```

For example:

```

header_start_tag=<myHeaderTag>
header_end_tag=</myHeaderTag>
footer_start_tag=<myFooterTag>
footer_end_tag=</myFooterTag>

```

NOTE:

You must encode custom tags in UTF-8.

Error Messages

When a KeyView exception is thrown, it might be caused by one of the following errors.

Exception	Description
KVERR_Success	Function completed successfully.
KVERR_DLLNotFound	A DLL or shared library was not found.
KVERR_OutOfCore	Memory allocation failure.
KVERR_processCancelled	Callback function returns FALSE.
KVERR_badInputStream	Invalid or corrupt input stream.
KVERR_badOutputType	Invalid output is requested.
KVERR_General	General error.
KVERR_FormatNotSupported	File format is not supported.
KVERR_PasswordProtected	File is encrypted or password-protected. KeyView only supports secure PST, NSF, and ZIP files.
KVERR_ADSNotFound	Adobe Document Server not found. This error is obsolete.
KVERR_AutoDetFail	Autodetect error.
KVERR_AutoDetNoFormat	Unable to detect file format.
KVERR_ReaderInitError	Error initializing the reader.
KVERR_NoReader	No reader available for this format.
KVERR_CreateOutputFileFailed	Unable to create output file. If the overwrite flag in <code>setOverWrite</code> is FALSE and a subfile has the same name as a file in the target path, this error is generated.
KVERR_CreateTempFileFailed	Unable to create temporary file.
KVERR_ErrorWritingToOutputFile	Error writing to output file.
KVERR_CreateProcessFailed	Error creating a child process.
KVERR_WaitForChildFailed	Wait for child process failed.
KVERR_ChildTimeOut	Child process hung/timed out.
KVERR_ArchiveFileNotFound	Attempt to extract nonexistent file.
KVERR_ArchiveFatalError	Fatal error processing an archive file.
KVError_OpenStreamFailure = KVERR_ArchiveFatalError + 1	Failed to open a stream during out-of-process filtering.
KVError_	An interface function was not found during out-of-process filtering.

Exception	Description
InterfaceFunctionNotFound	
KVError_InputFileNotFound	Could not find the input file during out-of-process filtering.
KVError_OpenOutputFileFailed	Could not open the output file during out-of-process filtering.
KVError_MemoryLeak	Memory leak occurred during out-of-process filtering.
KVError_MemoryOverwrite	Memory overwrite occurred during out-of-process filtering.
KVError_GPF	Exception occurred during out-of-process filtering.
KVError_OopCore	Memory dump was generated in a child process during out-of-process filtering.
KVError_KVoopLogFailed	Creation of out-of-process error log failed.
KVError_OverNestedFileLimit	The container file has more than the allowable number of child documents. One or more child documents were not converted. Currently, this is not used.
KVError_PSTAccessFailed	<p>The PST file could not be converted. This error might be returned when a call to <code>extOpenDocument</code> returns NULL for one of the following reasons:</p> <ul style="list-style-type: none"> • Microsoft Outlook client is not installed • Microsoft Outlook client is installed, but is not the default email client • Microsoft Outlook client is installed, but is not configured correctly • PST file is corrupt • PST file is read-only (PST files must allow read and write access) • MAPI call fails • The bit editions of Microsoft Outlook do not match the bit editions of the KeyView software. <p>For example, if 32-bit KeyView is used, 32-bit Outlook must be installed. If 64-bit KeyView is used, 64-bit Outlook must be installed.</p>
KVError_PasswordRequired	To open the file, credentials must be provided. This error might be returned when a call to <code>extOpenDocument</code> returns NULL.
KVError_InvalidArgs	The input argument or structure is invalid. This is generated by the File Extraction APIs.
KVError_OutputFileExists	A file with the same name already exists in the output directory. This error is generated when extracting a subfile from a container

Exception	Description
	file with the <code>setOverWrite</code> flag set to <code>FALSE</code> , and a file by the same name already exists in the output directory.
<code>KVError_ReaderUsageDenied</code>	<p>The current license key does not enable the document reader required to filter the file. This error might be returned when a call to <code>extOpenDocument</code> returns <code>NULL</code>.</p> <p>Some document readers are considered advanced features and are licensed separately from the KeyView SDK (for example, the PST and MBX readers). Contact your Autonomy sales representative to get an updated license key</p>
<code>KVError_OopBadConfig</code>	Information in the <code>kvxconfig.ini</code> file is incomplete and cannot be used to filter the XML file.
<code>KVError_OopBrokenPipe</code>	Data was not transferred between the parent and child processes during out-of-process filtering because either the parent or child failed.
<code>KVError_OopPipeOEF</code>	Data was not transferred between the parent and child processes during out-of-process filtering because the parent process was shutdown.
<code>KVError_IPCTimeOut</code>	Either the parent or child process is waiting for a reply or request during out-of-process filtering.
<code>KVError_InvalidOopDriverSignature</code>	A client sent a request to the File Extraction out-of-process server, but context driver does not exist on the server.
<code>KVError_InvalidOopServiceSignature</code>	<p>A client sent a request to a File Extraction out -of-process server that does not exist.</p> <p>If this error is generated on the call to <code>fpClose()</code>, it can be ignored.</p>

Tab Delimited Output for Embedded Tables

You can use KeyView to convert embedded tables in Word Processing documents (for example, Microsoft Word) to tab-delimited form, by specifying the following option in the `formats.ini` file:

```
[Options]
TabDelimitedOutput=TRUE
```

This option inserts a tab character between each cell, and a line break between each row. Tab and line break characters in the cells are replaced with spaces.

Exclude Japanese Guide Text

This option prevents output of Japanese phonetic guide text when Microsoft Excel (.xlsx) files are processed.

To prevent output of Japanese phonetic guide text

- Set NoPhoneticGuides to TRUE in the formats.ini file:

```
[Options]  
NoPhoneticGuides=TRUE
```

You can also enable this option programatically when filtering by passing KVFLT_NOPHONETICGUIDES to fpFilterConfig.

Chapter 5: Sample Programs

This section describes the sample programs provided with Filter SDK.

• Introduction	82
• ExtractFilter	83
• FilterTest	84
• FilterFileToFile	87
• FilterStreamToStream	88
• FilterFileToStream	88
• FilterStreamToFile	89
• FilterFileByChunk	90
• FilterStreamByChunk	91

Introduction

The following Java sample programs are provided:

- [ExtractFilter](#), on the next page
- [FilterTest](#), on page 84
- [FilterFileToFile](#), on page 87
- [FilterStreamToStream](#), on page 88
- [FilterFileToStream](#), on page 88
- [FilterStreamToFile](#), on page 89
- [FilterFileByChunk](#), on page 90
- [FilterStreamByChunk](#), on page 91

The source code for each program is in the directory `install\javaapi\sample`, where `install` is the path name of the Filter installation directory.

In addition to the sample programs, the following batch and C shell files that run each program are also provided:

`FilterFileToFile.bat (.csh)`

`FilterStreamToStream.bat (.csh)`

`FilterFileToStream.bat (.csh)`

`FilterStreamToFile.bat (.csh)`

`FilterFileByChunk.bat (.csh)`

`FilterStreamByChunk.bat (.csh)`

NOTE:

The `ExtractFilter` and `FilterTest` sample programs do not use batch or C shell files.

ExtractFilter

The ExtractFilter program demonstrates the File Extraction interface. The FilterTest sample program demonstrates the functionality of the Filtering interface. See [FilterTest, on the next page](#).

The ExtractFilter program demonstrates the following functionality:

- opens a document
- extracts subfiles from a document
- repeats subfile extraction until all subfiles are extracted
- enables you to specify the command-line options listed in the following table

To run ExtractFilter

1. Add the location of the javaapi\KeyView.jar file, the javaapi\sample directory, and the Filter bin directory to the CLASSPATH environment variable.
2. Type the following:

```
java -Djava.library.path=bin_directory ExtractFilter [options] bin_directory
input_file output_file
```

where,

bin_directory is the path to the Filter bin directory.

options is one or more of the options listed in the following table.

input_file is the path and file name of the source file.

output_file is the path and file name of the output file if the source file is not a container file.

Options for ExtractFilter Sample Program

Option	Description
-extonly	Extracts the subfiles from a source file but does not filter the files after extraction.
-extdir directory	Sets the directory to which the subfiles are extracted.
-ext-fbody	Extracts the formatted version of the message body (HTML or RTF) from mail files when possible.
-source-cs charset	Sets the character set of the source file. charset is a character set defined in the Filter class. See Coded Character Sets, on page 133 .
-target-cs charset	Sets the character set of the output file. charset is a character set defined in the Filter class. See Coded Character Sets, on page 133 .
-little-end	Sets the byte order for Unicode text to Little Endian.

Options for ExtractFilter Sample Program, continued

Option	Description
-is	Sets the input as a stream. The default is file.
-os	Sets the output as a stream. The default is file.
-ip	Runs file extraction in the same process as the calling application (in process). See Run Filter In Process, on page 29 .
-open-user username	Specifies the user name used to open a protected PST file.
-open-pass password	Specifies the password used to open a protected PST file.
-open-idfile idfile	Specifies the user ID file used to open a protected PST file.
-open-createroot	Creates a root directory on which a hierarchy can be based. See Create a Root Node, on page 36 .
-ext-nodir	Specifies the subfile directory structure is not created.
-ext-noheader	Excludes mail header information from extracted message body text file. See Exclude Metadata from the Extracted Text File, on page 43 .
-meta outfile	Extracts default mail metadata and writes it to a file. See Extract Mail Metadata, on page 38 .

FilterTest

The FilterTest program demonstrates most of the Filtering methods available in the Java API. It filters an input document to an output document and enables you to specify command-line options. The command-line options are listed in [Options for FilterTest Sample Program, on the next page](#).

To run FilterTest

1. Add the location of the `javaapi\KeyView.jar` file, the `javaapi\sample` directory, and the Filter bin directory to the CLASSPATH environment variable.
2. Type the following command line:

```
java -Djava.library.path=bin_directory FilterTest [options] bin_directory
input_file output_file
```

where,

bin_directory is the path to the Filter bin directory.

options is one or more of the options listed in [Options for FilterTest Sample Program, on the next page](#).

input_file is the path and file name of the source file.

output_file is the path and file name of the generated file. If a path is not specified, the file is output to the current directory.

Options for FilterTest Sample Program

Option	Description
-is	Sets the input as a stream. The default is file.
-os	Sets the output as a stream. The default is file.
-chunk	Filters an input source and returns one chunk of output data. The program calls the filter method repeatedly until the entire output buffer is processed.
-docformat filename	Extracts the file format information and writes it to a file. filename is the name of the file to which the format information is written.
-summary filename	Extracts the metadata and writes it to a file. filename is the name of the file to which the metadata is written. See Extract Metadata, on page 58 .
-getTargetCS	Extracts the character set used in the output file to the standard output.
-c charset	Sets the character set of the output file. Use the option -getTargetCS to determine whether the target character set specified is used in the output file. charset is a character set defined in the Filter class. See Coded Character Sets, on page 133 .
-cs charset	Sets the character set of the source file. charset is a character set defined in the Filter class. See Coded Character Sets, on page 133 .
-rc character	Sets a replacement character for characters that cannot be mapped. The default is a question mark (?).
-ip	Runs Filter in the same process as the calling application (in process). See Run Filter In Process, on page 29 .
-ooplog	Enables error logging. See Enable or Disable Error Logging, on page 56 . Error logs are not generated when in-process filtering is enabled.
-oopmem	Enables the memory trace system in the error logs. The memory trace system reports memory leaks and memory overwrites in the log file. See Report Memory Errors, on page 57 . Error logs are not generated when in-process filtering is enabled.
-hf	Extracts headers and footers, as well as the body text.
-hftags	Puts tags around header and footer data.
-lo	Specifies that PowerPoint PPT97 and PPTX file text data is output in a logical

Options for FilterTest Sample Program, continued

Option	Description
	reading order.
-lsbmsb	Uses LSBMSB byte order for Unicode text. LSBMSB is the "Least Significant Byte Most Significant Byte," or in other words, the byte order for Little Endian systems.
-msblsb	For Unicode text, uses MSBLSB byte order. MSBLSB is the "Most Significant Byte Least Significant Byte," or in other words, the byte order for Big Endian systems.
-bomarker	Generates the byte order marker for Unicode text.
-nodefcsconv	Prevents default conversion of document character encoding. See Prevent the Default Conversion of a Character Set, on page 63 .
-x xmlconfigfile	Filters an XML file using customized extraction settings defined in the kvxconfig.ini file. If you do not enter the full path to the INI file, the program looks for the file in the current working directory. See Filter XML Files, on page 72 .
-z tempdirectory	Specifies a temporary directory where temporary files generated by the filtering process are stored. The default is the current working directory. On Windows systems, there is a 64 K size limit to the temp directory. Once the limit is reached, you must either create a new directory or delete the contents of the existing directory; otherwise, you might receive an error message.
-ps password	Specifies a password to open a password-protected PST file. This uses the Container API which is obsolete.
-pdforder orderFlag	Specifies that PDF files are output in a logical reading order. The parameter orderFlag is one of the following: <ul style="list-style-type: none"> ltr—left-to-right paragraph direction. rtl—right-to-left paragraph direction. auto—The PDF filter determines the paragraph direction (left-to-right or right-to-left) for each PDF page, and then sets the direction accordingly. raw—Unstructured paragraph flow. See Filter PDF Files, on page 64 .
-rm	If you set this option, text that was deleted from a document with revision tracking enabled is extracted from the document and included in the filtered output. See Extract Tracked Deleted Text, on page 63 .
-embeddedfont	If you set this option, text that contains embedded fonts is not filtered from PDF documents. See Filter PDF Files, on page 64 .

FilterFileToFile

The FilterFileToFile program filters an input file to an output file using Java API methods in Filter. It demonstrates the following functions:

- filters an input file to an output file.
- extracts the character set if it can be determined by the document reader.
- extracts file format information (document type, format, version, and so on) if available in the source document.
- extracts metadata if available in the source document. This program extracts all the metadata from the document, but only displays the first element of metadata.

Run FilterFileToFile on Windows

To run FilterFileToFile on Windows

1. In the FilterFileToFile.bat file, set INSTALL_DIR to the Filter installation directory.
2. Run the batch file in the directory *install\javaapi\sample*, where *install* is the path name of the Filter installation directory. Type the following:

```
filterfiletofile inputfile outputfile
```

where,

inputfile is the path and file name of the source file.

outputfile is the path and file name of the generated file. If a path is not specified, the file is output to the current directory.

Run FilterFileToFile on UNIX

To run FilterFileToFile on UNIX

1. In the FilterFileToFile.csh file, set MKENV to the platform you are running, either *_hpux11*, *_ilnx21*, *_rs6k43*, or *_sso126*.
2. In the FilterFileToFile.csh file, set INSTALL_DIR to the Filter installation directory.
3. Run the C shell file in the directory *install/javaapi/sample*, where *install* is the path name of the Filter installation directory. Type the following:

```
./FilterFileToFile.csh inputfileoutputfile
```

where,

inputfile is the path and file name of the source file.

outputfile is the path and file name of the generated file. If a path is not specified, the file is output to the current directory.

FilterStreamToStream

The FilterStreamToStream program filters an input stream to an output stream using Java API methods in Filter. It demonstrates the following functions:

- creates an input and an output stream. Filters the input stream to the output stream.
- extracts file format information (document type, format, version, and so on) if available in the source document.
- extracts metadata if available in the source document. This program extracts all the metadata from the document, but only displays the first element of metadata.

Run FilterStreamToStream on Windows

To run FilterStreamToStream on Windows

1. In the FilterStreamToStream.bat file, set INSTALL_DIR to the Filter installation directory.
2. Run the batch file in the directory `install\javaapi\sample`, where `install` is the path name of the Filter installation directory. Type the following:

```
filterstreamtostream inputfile
```

where,

`inputfile` is the path and file name of the source file.

The generated text is output to the current DOS prompt.

Run FilterStreamToStream on UNIX

To run FilterStreamToStream on UNIX

1. In the FilterStreamToStream.csh file, set MKENV to the platform you are running, either `_hpux11`, `_ilnx21`, `_rs6k43`, or `_sso126`.
2. In the FilterStreamToStream.csh file, set INSTALL_DIR to the Filter installation directory.
3. Run the C shell file in the directory `install/javaapi/sample`, where `install` is the path name of the Filter installation directory. Type the following:

```
./FilterStreamToStream.csh inputfile
```

where,

`inputfile` is the path and file name of the source file.

The generated text is output to the current console (standard out).

FilterFileToStream

The FilterFileToStream program filters an input file to an output stream using Java API methods in Filter.

Run FilterFileToStream on Windows

To run FilterFileToStream on Windows

1. In the FilterFileToStream.bat file, set INSTALL_DIR to the Filter installation directory.
2. Run the batch file in the directory install\javaapi\sample, where install is the path name of the Filter installation directory. Type the following:

```
filterfiletostream inputfile
```

where,

inputfile is the path and file name of the source file.

The generated text is output to the current DOS prompt.

Run FilterFileToStream on UNIX

To run FilterFileToStream on UNIX

1. In the FilterFileToStream.csh file, set MKENV to the platform you are running, either _hpux11, _ilnx21, _rs6k43, or _ssol26.
2. In the FilterFileToStream.csh file, set INSTALL_DIR to the Filter installation directory.
3. Run the C shell file in the directory install/javaapi/sample, where install is the path name of the Filter installation directory. Type the following:

```
./FilterFileToStream.csh inputfile
```

where,

inputfile is the path and file name of the source file.

The generated text is output to the current console (standard out).

FilterStreamToFile

The FilterStreamToFile program filters an input stream to an output file using Java API methods in Filter.

Run FilterStreamToFile on Windows

To run FilterStreamToFile on Windows

1. In the FilterStreamToFile.bat file, set INSTALL_DIR to the Filter installation directory.
2. Run the batch file in the directory install\javaapi\sample, where install is the path name of the Filter installation directory. Type the following:

```
filterstreamtofile inputfileoutputfile
```

where,

`inputfile` is the path and file name of the source file.

`outputfile` is the path and file name of the generated file. If a path is not specified, the file is output to the current directory.

Run FilterStreamToFile on UNIX

To run FilterStreamToFile on UNIX

1. In the `FilterStreamToFile.csh` file, set `MKENV` to the platform you are running, either `_hpux11`, `_ilnx21`, `_rs6k43`, or `_ssol26`.
2. In the `FilterStreamToFile.csh` file, set `INSTALL_DIR` to the Filter installation directory.
3. Run the C shell file in the directory `install/javaapi/sample`, where `install` is the path name of the Filter installation directory. Type the following:

```
./FilterStreamToFile.csh inputfileoutputfile
```

where,

`inputfile` is the path and file name of the source file.

`outputfile` is the path and file name of the generated file. If a path is not specified, the file is output to the current directory.

FilterFileByChunk

The `FilterFileByChunk` program filters an input file to an output file using the Java API method `doFilterChunk()`. The method filters an input source and returns one chunk of output data. The program calls the method repeatedly until the entire file is processed.

Run FilterFileByChunk on Windows

To run FilterFileByChunk on Windows

1. In the `FilterFileByChunk.bat` file, set `INSTALL_DIR` to the Filter installation directory.
2. Run the batch file in the directory `install\javaapi\sample`, where `install` is the path name of the Filter installation directory. Type the following:

```
filterfilebychunkinputfile outputfile
```

where,

`inputfile` is the path and file name of the source file.

`outputfile` is the path and file name of the generated file. If a path is not specified, the file is output to the current directory.

Run FilterFileByChunk on UNIX

To run FilterFileByChunk on UNIX

1. In the `FilterFileByChunk.csh` file, set `MKENV` to the platform you are running, either `_hpux11`, `_ilnx21`, `_rs6k43`, or `_ssol26`.
2. In the `FilterFileByChunk.csh` file, set `INSTALL_DIR` to the Filter installation directory.
3. Run the C shell file in the directory `install/javaapi/sample`, where `install` is the path name of the Filter installation directory. Type the following:

```
/FilterFileByChunk.csh inputfile outputfile
```

where,

`inputfile` is the path and file name of the source file.

`outputfile` is the path and file name of the generated file. If a path is not specified, the file is output to the current directory.

FilterStreamByChunk

The `FilterStreamByChunk` program filters an input stream to an output stream using the Java API method `doFilterChunk()`. The method filters an input source and returns one chunk of output data. The program calls the method repeatedly until the entire output buffer is processed.

Run FilterStreamByChunk on Windows

To run FilterStreamByChunk on Windows

1. In the `FilterStreamByChunk.bat` file, set `INSTALL_DIR` to the Filter installation directory.
2. Run the batch file in the directory `install\javaapi\sample`, where `install` is the path name of the Filter installation directory. Type the following:

```
filterstreambychunk inputfile outputfile
```

where,

`inputfile` is the path and file name of the source file.

`outputfile` is the path and file name of the generated file. If a path is not specified, the file is output to the current directory.

Run FilterStreamByChunk on UNIX

To run FilterStreamByChunk on UNIX

1. In the `FilterStreamByChunk.csh` file, set `MKENV` to the platform you are running, either `_hpux11`, `_ilnx21`, `_rs6k43`, or `_ssol26`.
2. In the `FilterStreamByChunk.csh` file, set `INSTALL_DIR` to the Filter installation directory.

3. Run the C shell file in the directory `install/javaapi/sample`, where `install` is the path name of the Filter installation directory. Type the following:

```
./FilterStreamByChunk.csh inputfile outputfile
```

where,

`inputfile` is the path and file name of the source file.

`outputfile` is the path and file name of the generated file. If a path is not specified, the file is output to the current directory.

Appendixes

This section lists supported formats, supported character sets, and redistributed files, and provides information on format detection and developing a custom document reader.

Appendix A: Supported Formats

This section lists information about the file formats that can be detected and processed (either filtered, converted, or displayed) by the KeyView suite of products. The KeyView suite includes KeyView Filter SDK, KeyView Export SDK, and KeyView Viewing SDK.

- [Supported Formats](#) 94
- [Supported Formats \(Detected\)](#) 118

Supported Formats

The tables in this section provide the following information:

- The file formats supported by the Filter API, Export API, Viewing API, and File Extraction API. The supported versions and the format's extension are also listed.

The formats listed in this section can also be detected by the KeyView format detection module (kwad). The [Supported Formats \(Detected\)](#) section lists formats that can be detected, but cannot be filtered, converted, or displayed.

- The file formats for which KeyView can detect and extract the character set and metadata information (properties such as title, author, and subject).

Even though a file format might be able to provide character set information, some documents might not contain character set information. Therefore, the document reader would not be able to determine the character set of the document. In this case, either the operating system code page or the character set specified in the API is used.

- The document reader used to filter each format.

Key to Support Tables

Symbol	Description
Y	The format is supported. You can extract metadata for this format. You can determine the character set for this format.
N	The format is not supported. You cannot extract metadata for this format. You cannot determine the character set for this format.
P	Partial metadata is extracted from this format. Some non-standard fields are not extracted.
T	Only text is extracted from this format. Formatting information is not extracted.
M	Only metadata (title, subject, author, and so on) is extracted from this format. Text and

Key to Support Tables, continued

Symbol	Description
	formatting information are not extracted.

Archive Formats

Supported Archive Formats

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
7-Zip	4.57	z7zsr, multiarcsr ¹	7Z	N	N	Y	Y	N	n/a	N
AD1	n/a	ad1sr	AD1	N	N	Y	Y	N	n/a	N
ARJ	n/a	multiarcsr	ARJ	N	N	N	Y	N	n/a	N
B1	n/a	b1sr	B1	N	N	Y	Y	N	n/a	N
BinHex	n/a	kvhqxsr	HQX	N	N	Y	Y	N	n/a	N
Bzip2	n/a	bzip2sr	BZ2	N	N	Y	Y	N	n/a	N
Expert Witness Compression Format (EnCase)	6	encasesr	E01, L01	N	N	Y	Y	N	n/a	N
	7	encase2sr	Lx01	N	N	Y	Y	N	n/a	N
GZIP	2	kvgzsr	GZ	N	N	N	Y	N	n/a	N
		kvgz	GZ	N	N	Y	N	N	n/a	N
ISO	n/a	isosr	ISO	N	N	Y	Y	N	n/a	N
Java Archive	n/a	unzip	JAR	N	N	Y	Y	N	n/a	N
Legato EMailXtender	n/a	emxsr	EMX	N	N	Y	Y	N	n/a	N

¹7zip is supported with the multiarcsr reader on some platforms for Extract.

Supported Archive Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Archive										
MacBinary	n/a	macbinsr	BIN	N	N	Y	Y	N	n/a	N
Mac Disk Copy Disk Image	n/a	dmgsr	DMG	N	N	Y	Y	N	n/a	N
Microsoft Backup File	n/a	bkfsr	BKF	N	N	Y	Y	N	n/a	N
Microsoft Cabinet format	1.3	cabsr	CAB	N	N	Y	Y	N	n/a	N
Microsoft Compiled HTML Help	3	chmsr	CHM	N	N	Y	Y	N	n/a	N
Microsoft Compressed Folder	n/a	lzhsr	LZH LHA	N	N	N	Y	N	n/a	N
PKZIP	through 9.0	unzip	ZIP	N	N	Y	Y	N	n/a	N
RAR archive	2.0 through 3.5	rarsr	RAR	N	N	N	Y	N	n/a	N
RAR5 archive	5	multiarcsr	RAR5	N	N	N	Y	N	n/a	N
Tape Archive	n/a	tarsr	TAR	N	N	Y	Y	N	n/a	N
UNIX Compress	n/a	kvzeesr	Z	N	N	N	Y	N	n/a	N
		kvzee	Z	N	N	Y	N	N	n/a	N
UUEncoding	all versions	uudsr	UUE	N	N	Y	Y	N	n/a	N
XZ	n/a	multiarcsr	XZ	N	N	N	Y	N	n/a	N

Supported Archive Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Windows Scrap File	n/a	olesr	SHS	N	N	N	Y	N	n/a	N
WinZip	through 10	unzip	ZIP	N	N	Y	Y	N	n/a	N

Binary Format**Supported Binary Formats**

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Executable	n/a	exesr	EXE	N	N	Y	N	N	n/a	N
Link Library	n/a	exesr	DLL	N	N	Y	N	N	n/a	N

Computer-Aided Design Formats**Supported CAD Formats**

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
AutoCAD Drawing	R13, R14, R15/2000, 2004, 2007, 2010, 2013	kpODArdr kpDWGrdr ¹	DWG	Y	Y ²	Y ³	N	Y	Y	N

¹On Windows platforms, kpODArdr is used for all versions up to 2007 and graphic rendering is supported; for later versions, only text extraction is supported through the kpDWGrdr or kpDXFrdr reader.

²On non-Windows platforms, graphic rendering is supported through the kpDWGrdr reader for versions R13, R14, R15, and R18 (2004); for other versions, only text extraction is supported.

³On non-Windows platforms, graphic rendering is supported through the kpDWGrdr reader for versions R13, R14, R15, and R18 (2004); for other versions, only text extraction is supported.

Supported CAD Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
AutoCAD Drawing Exchange	R13, R14, R15/2000, 2004, 2007, 2010, 2013	kpODArdr kpDXFrdr ¹	DXF	Y	Y ²	Y ³	N	Y	Y	N
CATIA formats	5	kpCATrdr	CAT ⁴	Y	N	N	N	Y	N	N
Microsoft Visio	4, 5, 2000, 2002, 2003, 2007, 2010 ⁵	vsdsr	VSD	Y	Y	Y	Y ⁶	Y	Y	N
		kpVSD2rdr	VSD, VSS VST	Y	Y	Y	N	Y	Y	N
	2013	ActiveX components	VSDM VSSM VSTM VSDX VSSX VSTX	N	N	Y ⁷	N	Y	N	N
		kpVSDXrdr	VSDM	Y	Y	Y ⁴	Y	Y	Y	N

¹On Windows platforms, kpODArdr is used for all versions up to 2007 and graphic rendering is supported; for later versions, only text extraction is supported through the kpDWGrdr or kpDXFrdr reader.

²On non-Windows platforms, graphic rendering is supported through the kpDXFrdr reader for versions R13, R14, R15, and R18 (2004); for other versions, only text extraction is supported.

³On Windows platforms, kpODArdr is used for all versions up to 2007 and graphic rendering is supported; for later versions, only text extraction is supported through the kpDWGrdr or kpDXFrdr reader.

⁴All CAT file extensions, for example CATDrawing, CATProduct, CATPart, and so on.

⁵Viewing and Export use the graphic reader, kpVSD2rdr for Microsoft Visio 2003, 2007, and 2010, and vsdsr for all earlier versions. Image fidelity in Viewing and Export is therefore only supported for versions 2003 and above. Filter uses the graphic reader kpVSD2rdr for Microsoft Visio 2003, 2007, and 2010, and vsdsr for all earlier versions.

⁶Extraction of embedded OLE objects is supported for Filter on Windows platforms only.

⁷Visio 2013 is supported in Viewing only, with the support of ActiveX components from the Microsoft Visio 2013 Viewer. Image fidelity is supported but other features, such as highlighting, are not.

Supported CAD Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
			VSSM VSTM VSDX VSSX VSTX							

Database Formats**Supported Database Formats**

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
dBase Database	III+, IV	dbfsr	DBF	Y	Y	Y	N	N	N	N
Microsoft Access	95, 97, 2000, 2002, 2003, 2007, 2010, 2013, 2016	mdbsr	MDB, ACCDB	Y	T	T	N	N	Y ¹	N
Microsoft Project	2000, 2002, 2003, 2007, 2010, 2013	mppsrs	MPP	Y	Y	Y	Y	Y	Y	N

Desktop Publishing**Supported Desktop Publishing Formats**

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Microsoft Publisher	98 to 2016	mspubsr	PUB	Y	T	T	Y	Y	Y	N

¹Charset is not supported for Microsoft Access 95 or 97.

Display Formats

Supported Display Formats

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Adobe PDF	1.1 to 1.7	pdfsr	PDF	Y	Y	N	Y ¹	Y	Y	N
		pdf2sr	PDF	N	Y	N	N	N	N	N
		kppdfldr	PDF	N	Y	Y	N	N	N	N
		kppdf2ldr ²	PDF	N	N	Y	N	N	N	N

Graphic Formats

Supported Graphic Formats

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Computer Graphics Metafile	n/a	kpcgmrdr ³	CGM	Y	Y	Y	N	N	N	N
CorelDRAW ⁴	through 9.0 10, 11, 12, X3	kpcdrdr	CDR	N	Y	Y	N	N	N	N

¹Includes support for extraction of subfiles from PDF Portfolio documents.

²kppdf2ldr is an alternate graphic-based reader that produces high-fidelity output but does not support other features such as highlighting or text searching.

³Files with non-partitioned data are supported.

⁴CDR/CDR with TIFF header.

Supported Graphic Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
DCX Fax System	n/a	kpdcxrdr	DCX	N	Y	Y	N	N	N	N
Digital Imaging & Communications in Medicine (DICOM)	n/a	dcmsr	DCM	M	N	N	N	Y	N	N
Encapsulated PostScript (raster)	TIFF header	kpepsrdr	EPS	N	Y	Y	N	N	N	N
Enhanced Metafile	n/a	kpemfrdr	EMF	Y	Y	Y	N	Y	N	N
GIF	87, 89	kpgifrdr	GIF	N	Y	Y	N	N	N	N
		gifsr		M	M	N	N	Y	N	N
JBIG2	n/a	kpJBIG2rdr	JBIG2	N	Y	Y	N	N	N	N
JPEG	n/a	kpjpgdr	JPEG	N	Y	Y	N	N	N	N
		jpgsr		M	M	N	N	Y	N	N
JPEG 2000	n/a	kpjp2000rdr	JP2, JPF, J2K, JPWL, JPX, PGX	N	Y	Y	N	N	N	N
		jp2000sr		M	M	N	N	Y	N	N
Lotus AMIDraw Graphics	n/a	kpsdwrdr	SDW	N	Y	Y	N	N	N	N
Lotus Pic	n/a	kppicrdr	PIC	Y	Y	Y	N	N	N	N
Macintosh Raster	2	kppctrdr	PIC PCT	N	Y	Y	N	N	N	N
MacPaint	n/a	kpmacrdr	PNTG	N	Y	Y	N	N	N	N

Supported Graphic Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Microsoft Office Drawing	n/a	kpmrdr	MSO	N	Y	Y	N	N	N	N
Omni Graffle	n/a	kpGFLrdr	GRAFFLE	Y	N	N	N	Y	Y	N
PC PaintBrush	3	kppcxrdr	PCX	N	Y	Y	N	N	N	N
Portable Network Graphics	n/a	kppngrdr	PNG	N	Y	Y	N	N	N	N
		pngsr	PNG	M	M	N	N	Y	N	N
SGI RGB Image	n/a	kpsgirdr	RGB	N	Y	Y	N	N	N	N
Sun Raster Image	n/a	kpsunrdr	RS	N	Y	Y	N	N	N	N
Tagged Image File	through 6.0 ¹	tifsr	TIFF	M	M	N	N	Y	N	N
		kptifdr	TIFF	N	Y	Y	N	N	N	N
Truevision Targa	2	kptrdr	TGA	N	Y	Y	N	N	N	N
Windows Animated Cursor	n/a	kpanirdr	ANI	N	Y	Y	N	N	N	N
Windows Bitmap	n/a	kpbmprdr	BMP	N	Y	Y	N	N	N	N
		bmpsr	BMP	M	M	N	N	Y	N	N
Windows Icon Cursor	n/a	kpicdr	ICO	N	Y	Y	N	N	N	N
Windows Metafile	3	kpwmfrdr	WMF	Y	Y	Y	N	N	N	N
WordPerfect Graphics 1	1	kpwpgrdr	WPG	N	Y	Y	N	N	N	N

¹The following compression types are supported: no compression, CCITT Group 3 1-Dimensional Modified Huffman, CCITT Group 3 T4 1-Dimensional, CCITT Group 4 T6, LZW, JPEG (only Gray, RGB and CMYK color space are supported), and PackBits.

Supported Graphic Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
WordPerfect Graphics 2	2, 7	kpwg2rdr	WPG	N	Y	Y	N	N	N	N

Mail Formats**Supported Mail Formats**

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Documentum EMC MF	n/a	msgsr	EMCMF	N	N	Y	Y	Y	Y	N
Domino XML Language ¹	n/a	dxlsr	DXL	N	N	Y	Y	Y	N	N
GroupWise FileSurf	n/a	gwfssr	GWFS	N	N	Y	Y	Y	N	N
Legato Extender	n/a	onmsr	ONM	N	N	Y	Y	Y	N	N
Lotus Notes database	4, 5, 6.0, 6.5, 7.0, 8.0	nsfsr	NSF	N	N	Y	Y	Y	N	N
Mailbox ²	Thunderbird 1.0,	mbxsr ³	MBX	N	N	T	Y	Y	Y	N

¹Supports non-encrypted embedded files only.

²KeyView supports MBX files created by Eudora Email and Mozilla Thunderbird. MBX files created by other common mail applications are typically filtered, converted, and displayed.

³This reader supports both clear signed and encrypted S/MIME. KeyView supports S/MIME for PST, EML, MBX, and MSG files.

Supported Mail Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
	Eudora 6.2									
Microsoft Entourage Database	2004	entsr	various	N	N	Y	Y	Y	Y	N
Microsoft Outlook	97, 2000, 2002, 2003, 2007, 2010, 2013, 2016	msgsr ¹	MSG, OFT	Y	T	T	Y	Y	Y ²	N
Microsoft Outlook DBX	5.0, 6.0	dbxsr	DBX	N	N	Y	Y	Y	Y	N
Microsoft Outlook Express	Windows 6 MacIntosh 5	emlsr ³	EML	Y	T	T	Y	Y	Y	N
		mbxsr ⁴	EML	N	N	T	Y	Y	Y	N
Microsoft Outlook iCalendar	1.0, 2.0	icssr	ICS, VCS	N	N	Y	Y	Y	Y	N
Microsoft Outlook for Macintosh	2011	olmsr	OLM	N	N	Y	Y	N	Y	N
Microsoft Outlook Offline Storage File	97, 2000, 2002, 2003, 2007, 2010,	pffsr ⁵	OST	N	N	Y	Y	Y	Y	N

¹This reader supports both clear signed and encrypted S/MIME. KeyView supports S/MIME for PST, EML, MBX, and MSG files.

²Returns "Unicode" character set for version 2003 and up, and "Unknown" character set for previous versions.

³This reader supports both clear signed and encrypted S/MIME. KeyView supports S/MIME for PST, EML, MBX, and MSG files.

⁴This reader supports both clear signed and encrypted S/MIME. KeyView supports S/MIME for PST, EML, MBX, and MSG files.

⁵The reader pffsr is available only on Windows and Linux.

Supported Mail Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
	2013									
Microsoft Outlook Personal Folder	97, 2000, 2002, 2003, 2007, 2010, 2013, 2016	pstsr ¹²	PST	N	N	Y	Y	Y	N	N
	97, 2000, 2002, 2003, 2007, 2010, 2013	pstnsr	PST	N	N	Y	Y	Y	Y	N
Microsoft Outlook vCard Contact	2.1, 3.0, 4.0	vcfsr	VCF	Y	Y	T	N	Y	N	N
Text Mail (MIME)	n/a	emlsr ³	various	Y	T	T	Y	Y	Y	N
		mbxsr ⁴	various	Y	T	T	Y	Y	Y	N
Transport Neutral Encapsulation Format	n/a	tnefsr	various	N	N	Y	Y	Y	Y	N

¹This reader supports both clear signed and encrypted S/MIME. KeyView supports S/MIME for PST, EML, MBX, and MSG files.

²Uses Microsoft Messaging Application Programming Interface (MAPI). Note that the native PST reader (*pstsr*) works only on Windows, and requires that you have Microsoft Outlook installed. As an alternative, the MAPI reader (*pstnsr*) runs on all platforms, and does not require Microsoft Outlook. For more information on using the native PST reader or the MAPI reader, see the sections 'Use the Native PST Reader (*pstnsr*)' and 'Use the MAPI Reader (*pstsr*)' in Chapter 3.

³This reader supports both clear signed and encrypted S/MIME. KeyView supports S/MIME for PST, EML, MBX, and MSG files.

⁴This reader supports both clear signed and encrypted S/MIME. KeyView supports S/MIME for PST, EML, MBX, and MSG files.

Multimedia Formats

Viewing SDK plays some multimedia files using the Windows Media Control Interface (MCI). MCI is a set of Windows APIs that communicate with multimedia devices.

Supported Multimedia Formats

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Advanced Systems Format	1.2	asfsr	ASF WMA WMV	N	N	N	N	Y	N	N
Audio Interchange File Format	n/a	MCI	AIFF	N	N	Y	N	N	N	N
		aiffr	AIFF	M	N	N	N	Y	N	N
Microsoft Wave Sound	n/a	MCI	WAV	N	N	Y	N	N	N	N
		riffr	WAV	M	N	N	N	Y	N	N
MIDI	n/a	MCI	MID	N	N	Y	N	N	N	N
MPEG-1 Audio layer 3	ID3 v1 and v2	MCI	MP3	N	N	Y	N	N	N	N
		mp3sr	MP3	M	M	Y	N	Y	N	N
MPEG-1 Video	2, 3	MCI	MPG	N	N	Y	N	N	N	N
MPEG-2 Audio	n/a	MCI	MPEGA	N	N	Y	N	N	N	N
MPEG-4 Audio	n/a	mpeg4sr	MP4 3GP	M	N	N	N	Y	N	N
NeXT/Sun Audio	n/a	MCI	AU	N	N	Y	N	N	N	N
QuickTime Movie	2, 3, 4	MCI	QT MOV	N	N	Y	N	N	N	N

Supported Multimedia Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Windows Video	2.1	MCI	AVI	N	N	Y	N	N	N	N

NOTE:

Depending on the default multimedia player installed on your computer, the View API might not be able to play some supported multimedia formats. To play multimedia files, the View API uses the Windows Media Control Interface (MCI) to communicate with the multimedia player installed on your computer. If the player does not play a multimedia file that is supported by the Viewing SDK, the View API cannot play the file.

If you cannot play a supported multimedia file by using the View API, install a different multimedia player or compressor/decompressor (codec) component.

Presentation Formats**Supported Presentation Formats**

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Apple iWork Keynote	2, 3, '08, '09	kplWPGdrdr	GZ	Y	Y	Y	N	Y	Y	N
Applix Presents	4.0, 4.2, 4.3, 4.4	kpagrdr	AG	Y	Y	Y	N	N	N	N
Corel Presentations	6, 7, 8, 9, 10, 11, 12, X3	kpshwrdr	SHW	Y	Y	Y	N	N	N	N
Extensible Forms Description Language	n/a	kpXFDLrdr	XFD XFDL	Y	Y	Y	N	Y	Y	N
Lotus Freelance Graphics	96, 97, 98, R9, 9.8	kpprzrdr	PRZ	Y	Y	Y	N	N	N	N
Lotus Freelance Graphics 2	2	kpprerdr	PRE	Y	Y	Y	N	N	N	N

Supported Presentation Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Macromedia Flash	through 8.0	swfsr	SWF	Y	Y	Y	N	N	Y ¹	N
Microsoft OneNote	2007, 2010, 2013, 2016	kpONErdr	ONE ONETOC2	Y	Y	Y	Y	N	Y	N
Microsoft PowerPoint Macintosh	98	kpp40rdr	PPT	Y	Y	Y	N	N	N	N
	2001, v.X, 2004	kpp97rdr	PPT PPS POT	Y	Y	Y	N	P	Y	N
Microsoft PowerPoint PC	4	kpp40rdr	PPT	Y	Y	Y	N	P	N	N
Microsoft PowerPoint Windows	95	kpp95rdr	PPT	Y	Y	Y	N	P	Y	N
Microsoft PowerPoint Windows	97, 2000, 2002, 2003	kpp97rdr	PPT PPS POT	Y	Y	Y	Y	P	Y	Y ²
Microsoft PowerPoint Windows XML	2007, 2010, 2013, 2016	kpppxrdr	PPTX PPTM POTX POTM PPSX PPSM PPAM	Y	Y	Y	Y	Y	Y	Y

¹The character set cannot be determined for versions 5.x and lower.²Slide footers are supported for Microsoft PowerPoint 97 and 2003.

Supported Presentation Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
OASIS Open Document Format	1, 2 ¹	kpodfrdr	SXD SXI ODG ODP	Y	Y	Y	Y ²	Y	Y	N
OpenOffice Impress, LibreOffice Impress	1 to 5	sosr	SXI SXP ODP	Y	T	T	N	Y	Y	N
StarOffice Impress	6, 7, 8, 9	sosr	SXI SXP ODP	Y	T	T	N	Y	Y	N

Spreadsheet Formats**Supported Spreadsheet Formats**

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Apple iWork Numbers	'08, '09	iwsssr	GZ	Y	Y	Y	N	Y	Y	N
	'13, '16	iwss13sr	NUMBERS	Y	T	T	N	N	Y	N
Applix Spreadsheets	4.2, 4.3, 4.4	assr	AS	Y	Y	Y	N	N	Y	N
Comma Separated Values	n/a	csvsr	CSV	Y	Y	Y	N	N	N	N

¹Generated by OpenOffice Impress 2.0, StarOffice 8 Impress, and IBM Lotus Symphony Presentation 3.0.

²Supported using the olesr embedded objects reader.

Supported Spreadsheet Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Corel Quattro Pro	5, 6, 7, 8	qpssr	WB2 WB3	Y	Y	Y	N	P	Y	N
	X4	qpwsr	QPW	Y	N	Y	N	P	Y	N
Data Interchange Format	n/a	difsr		Y	Y	Y	N	N	N	N
Lotus 1-2-3	96, 97, R9, 9.8	l123sr	123	Y	Y	Y	N	P	Y	N
Lotus 1-2-3	2, 3, 4, 5	wkssr	WK4	Y	Y	Y	N	N	Y	N
Lotus 1-2-3 Charts	2, 3, 4, 5	kpchtrdr	123	N	Y	Y	N	N	N	N
Microsoft Excel Charts	2, 3, 4, 5, 6, 7	kpchtrdr	XLS	N	Y	Y	N	N	N	N
Microsoft Excel Macintosh	98, 2001, v.X, 2004	xlssr	XLS	Y	Y	Y	Y ¹	Y	Y	N
Microsoft Excel Windows	2.2 through 2003	xlssr	XLS XLW XLT XLA	Y	Y	Y	Y ²	Y	Y	Y
Microsoft Excel Windows XML	2007, 2010, 2013, 2016	xlxsxr	XLSX XLTX XLSM XLTM XLAM	Y	Y	Y	Y	Y	Y	Y
Microsoft Excel Binary	2007, 2010,	xlsbsr	XLSB	Y	Y	Y	N	N	N	N

¹Supported using the embedded objects reader `olesr`.²Supported for versions 97 and higher using the embedded objects reader `olesr`.

Supported Spreadsheet Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Format	2013, 2016									
Microsoft Works Spreadsheet	2, 3, 4	mwssr	S30 S40	Y	Y	Y	N	N	Y	N
OASIS Open Document Format	1, 2 ¹	odfssr	ODS SXC STC	Y	Y	Y	Y ²	Y	Y	N
OpenOffice Calc, LibreOffice Calc	1 to 5	sosr	SXC ODS OTS	Y	T	T	N	Y	Y	N
StarOffice Calc	6, 7, 8, 9	sosr	SXC ODS	Y	T	T	N	Y	Y	N

Text and Markup Formats**Supported Text and Markup Formats**

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
ANSI	n/a	afsr	TXT	Y	Y	Y	N	N	N	N
ASCII	n/a	afsr	TXT	Y	Y	Y	N	N	N	N
HTML	3, 4	htmsr	HTM	Y	Y	Y	N	P	Y	N
Microsoft Excel Windows XML	2003	xmlsr	XML	Y	T	T	N	Y	Y	N

¹Generated by OpenOffice Calc 2.0, StarOffice 8 Calc, and IBM Lotus Symphony Spreadsheet 3.0.

²Supported using the embedded objects reader olesr.

Supported Text and Markup Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Microsoft Word Windows XML	2003	xmlsr	XML	Y	T	T	N	Y	Y	N
Microsoft Visio XML	2003	xmlsr	VDX VTX	Y	T	T	N	Y	Y	N
MIME HTML	n/a	mhtsr	MHT	Y	Y	Y	N	Y	Y	N
Rich Text Format	1 through 1.7	rtfsr	RTF	Y	Y	Y	N	P	Y	Y
Unicode HTML	n/a	unihtmsr	HTM	Y	Y	Y	N	Y	Y	N
Unicode Text	3, 4	unisr	TXT	Y	Y	Y	N	N	Y	N
XHTML	1.0	htmsr	HTM	Y	Y	Y	N	Y	Y	N
XML (generic)	1.0	xmlsr	XML	Y	T	T	N	Y	Y	N

Word Processing Formats**Supported Word Processing Formats**

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Adobe FrameMaker Interchange Format	5, 5.5, 6, 7	mifsr	MIF	Y	Y	Y	N	N	Y	N
Apple iChat Log	1, AV 2 AV 2.1, AV 3	ichatsr	ICHAT	Y	Y	Y	N	N	N	N

Supported Word Processing Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Apple iWork Pages	'08, '09	iwwpsr	GZ	Y	Y	Y	N	Y	Y	N
	'13, '16	iwwp13sr	PAGES	Y	T	T	N	N	N	N
Applix Words	3.11, 4, 4.1, 4.2, 4.3, 4.4	awsr	AW	Y	Y	Y	N	N	Y	Y
Corel WordPerfect Linux	6.0, 8.1	wp6sr	WPS	Y	Y	Y	N	P	Y	N
Corel WordPerfect Macintosh	1.02, 2, 2.1, 2.2, 3, 3.1	wpmsr	WPM	Y	Y	Y	N	N	Y	N
Corel WordPerfect Windows	5, 5.1	wosr	WO	Y	Y	Y	N	P	Y	Y
Corel WordPerfect Windows	6, 7, 8, 9, 10, 11, 12, X3	wp6sr	WPD	Y	Y	Y	N	P	Y	Y
DisplayWrite	4	dw4sr	IP	Y	Y	Y	N	N	Y	N
Folio Flat File	3.1	foliosr	FFF	Y	Y	Y	N	Y	Y	Y
Founder Chinese E-paper Basic	3.2.1	cebsr ¹	CEB	Y	N	N	N	N	N	N
Fujitsu Oasys	7	oa2sr	OA2	Y	Y	Y	N	P	N	N

¹This reader is only supported on Windows 32-bit platforms.

Supported Word Processing Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Haansoft Hangul	97	hwpsr	HWP	Y	N	N	N	N	Y	N
	2002, 2005, 2007, 2010	hwposr	HWP	Y	T	T	Y	Y	Y	N
Health level7	2.0	hl7sr	HL7	Y	Y	Y	N	Y	Y	N
IBM DCA/RFT (Revisable Form Text)	SC23-0758-1	dcasr	DC	Y	Y	Y	N	N	Y	N
JustSystems Ichitaro	8 through 2013	jtdsr	JTD	Y	Y	Y	N	P	N	Y
Lotus AMI Pro	2, 3	lasr	SAM	Y	Y	Y	N	P	Y	Y
Lotus AMI Professional Write Plus	2.1	lasr	AMI	Y	Y	Y	N	N	N	Y
Lotus Word Pro	96, 97, R9	lwpsr	LWP	Y	Y	Y	N	P	N	Y
Lotus SmartMaster	96, 97	lwpsr	MWP	Y	Y	Y	N	N	N	N
Microsoft Word Macintosh	4, 5, 6, 98	mbsr	DOC	Y	Y	Y	N	Y	N	Y
	2001, v.X, 2004	mw8sr	DOC DOT	Y	Y	Y	Y ¹	Y	Y	N
Microsoft Word PC	4, 5, 5.5, 6	mwsr	DOC	Y	Y	Y	N	N	N	Y
Microsoft Word Windows	1.0 and 2.0	misr	DOC	Y	Y	Y	N	N	N	Y

¹Supported using the embedded objects reader olesr.

Supported Word Processing Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Microsoft Word Windows	6, 7, 8, 95	mw6sr	DOC	Y	Y	Y	N	Y	Y	Y
Microsoft Word Windows	97, 2000, 2002, 2003	mw8sr	DOC DOT	Y	Y	Y	Y ¹	Y	Y	Y
Microsoft Word Windows XML	2007, 2010, 2013, 2016	mwxsr	DOCM DOCX DOTX DOTM	Y	Y	Y	Y	Y	Y	Y
Microsoft Works	1, 2, 3, 4	mswsr	WPS	Y	Y	Y	N	N	N	Y
Microsoft Works	6, 2000	msw6sr	WPS	Y	Y	Y	N	N	N	Y
Microsoft Windows Write	1, 2, 3	mwsr	WRI	Y	Y	Y	N	N	Y	N
OASIS Open Document Format	1, 2 ²	odfwpsr	ODT SXW STW	Y	Y	Y	Y ³	Y	Y	Y
Omni Outliner	v3, OPML, OOutline	oo3sr	OO3 OPML OOUTLINE	Y	Y	Y	N	N	Y	N
OpenOffice Writer, LibreOffice Writer	1 to 5	sosr	SXW ODT	Y	T	T	N	Y	Y	N

¹Supported using the embedded objects reader `olesr`.²Generated by OpenOffice Writer 2.0, StarOffice 8 Writer, and IBM Lotus Symphony Documents 3.0.³Supported using the embedded objects reader `olesr`.

Supported Word Processing Formats, continued

Format	Version	Reader	Extension	Filter	Export	View	Extract	Metadata	Charset	Header/Footer
Open Publication Structure eBook	2.0, 3.0	epubsr	EPUB	Y	Y	Y	N	Y	Y	N
StarOffice Writer	6, 7, 8, 9	sosr	SXW ODT	Y	T	T	N	Y	Y	N
Skype Log	3	skypesr	DBB	Y	Y	Y	N	N	N	N
WordPad	through 2003	rtfsr	RTF	Y	Y	Y	N	P	Y	N
XML Paper Specification	n/a	xpssr	XPS	Y	T	T	N	N	N	N
XyWrite	4.12	xywsr	XY4	Y	Y	Y	N	N	N	N
Yahoo! Instant Messenger	n/a	yimsr ¹	DAT	Y	Y	Y	N	N	N	N

¹To successfully use this reader, you must set the KV_YAHOO_ID environment variable to the Yahoo user ID. You can optionally set the KV_OTHER_YAHOO_ID environment variable to the other Yahoo user ID. If you do not set it, "Other" is used by default. If you enter incorrect values for the environment variables, erroneous data is generated.

Supported Formats (Detected)

The file formats listed in this section can be detected by the KeyView format detection module (*kwad*), but cannot be filtered, converted, or displayed. The detection module determines a file's format and reports the information to the developer's application.

The formats listed in [Supported Formats, on page 94](#) can be detected as well as filtered, exported, and viewed.

- Ability Office (SS, DB, GR, WP, COM)
- AC3 audio
- ACT
- Adobe FrameMaker
- Adobe FrameMaker Markup Language
- AES Multiplus Comm
- Aldus Freehand (Macintosh)
- Aldus PageMaker (DOS)
- Aldus PageMaker (Macintosh)
- Amiga IFF-8SVX sound
- Amiga MOD sound
- Apple Binary Property List
- Apple Double
- Apple iWork
- Apple Photoshop Document
- Apple Single
- Apple XML Property List
- Appleworks
- Applix Alis
- Applix Asterix
- Applix Graphics
- ARC/PAK Archive
- ASCII-armored PGP encoded
- ASCII-armored PGP Public Keyring
- ASCII-armored PGP signed
- AutoDesk Animator FLIC Animation
- AutoDesk Animator Pro FLIC Animation
- AutoDesk WHIP
- AutoShade Rendering
- B1 Archive
- BlackBerry Activation File

- CADAM Drawing
- CADAM Drawing Overlay
- CCITT Group 3 1-Dimensional (G31D)
- COMET TOP Word
- Confifer Software WavPack
- Convergent Tech DEF Comm.
- Corel Draw CMX
- cpio Archive (UNIX/VAX/SUN)
- CPT Communication
- Creative Voice (VOC) sound
- Curses Screen Image (UNIX/VAX/SUN)
- Data Point VISTAWORD
- DCX Fax
- DEC WPS PLUS
- DECdx
- Desktop Color Separation (DCS)
- Device Independent file (DVI)
- DG CEOwrite
- DG Common Data Stream (CDS)
- DIF Spreadsheet
- Digital Document Interchange Format (DDIF)
- Digital Imaging and Communications in Medicine (DICOM)
- Disk Doubler Compression
- EBCDIC Text
- eFax
- ENABLE
- ENABLE Spreadsheet (SSF)
- Envoy (EVY)
- Executable UNIX/VAX/SUN
- FileMaker (Macintosh)
- FPX format
- Framework
- Framework II
- Freehand 11
- FTP Session Data
- GEM Bit Image
- Ghost Disk Image
- Google SketchUp

- Graphics Environment Manager (GEM VDI)
- Harvard Graphics
- Hewlett Packard
- Honey Bull DSA101
- HP Graphics Language (HP-GL)
- HP Graphics Language (Plotter)
- HP PCL and PJP Languages
- HP Word PC
- IBM 1403 Line Printer
- IBM DCA-FFT
- IBM DCF Script
- Informix SmartWare II
- Informix SmartWare II Communication File
- Informix SmartWare II Database
- Informix SmartWare Spreadsheet
- Interleaf
- Java Class file
- JPEG File Interchange Format (JFIF)
- Keyhole Markup Language
- KW ODA G4 (G4)
- KW ODA G31D (G31)
- KW ODA Internal G32D (G32)
- KW ODA Internal Raw Bitmap (RBM)
- Lasergraphics Language
- Link Library UNIX/VAX/SUN
- Lotus Notes Bitmap
- Lotus Notes CDF
- Lotus Screen Cam
- Lyrinx
- Macromedia Director
- MacWrite
- MacWrite II
- MASS-11
- MATLAB MAT Format
- Micrografx Designer
- Microsoft Access 2007
- Microsoft Access 2007 Template
- Microsoft Common Object File Format (COFF)

- Microsoft Compiled HTML Help
- Microsoft Device Independent Bitmap
- Microsoft Document Imaging (MDI)
- Microsoft Excel 2007 Macro-Enabled Spreadsheet Template
- Microsoft Excel 2007 Spreadsheet Template
- Microsoft Exchange Server Database File
- Microsoft Object File Library
- Microsoft Office Drawing
- Microsoft Office Groove
- Microsoft Outlook Restricted Permission Message File
- Microsoft Windows Cursor (CUR) Graphics
- Microsoft Windows Group File
- Microsoft Windows Help File
- Microsoft Windows Icon (ICO)
- Microsoft Windows NT Event Log
- Microsoft Windows OLE 2 Encapsulation
- Microsoft Windows Vista Event Log
- Microsoft Word (UNIX)
- Microsoft Works (Macintosh)
- Microsoft Works Communication (Macintosh)
- Microsoft Works Communication (Windows)
- Microsoft Works Database (Macintosh)
- Microsoft Works Database (PC)
- Microsoft Works Database (Windows)
- Microsoft Works Spreadsheet (Macintosh)
- Microstation
- Milestone Document
- MORE Database Outliner (Macintosh)
- MPEG4 (ISO IEC MPEG4)
- MPEG-PS container with CDXA stream
- MS DOS Batch File format
- MS DOS Device Driver
- MultiMate 4.0
- Multiplan Spreadsheet
- Navy DIF
- NBI Async Archive Format
- NBI Net Archive Format
- Nero Encrypted File

- Netscape Bookmark file
- NeWS font file (SUN)
- NIOS TOP
- Nota Bene
- NURSTOR Drawing
- Object Module UNIX/VAX/SUN
- ODA/ODIF
- ODA/ODIF (FOD 26)
- Office Writer
- OLE DIB object
- OLIDIF
- Open PGP (new format packets)
- OS/2 PM Metafile Graphics
- PaperPort image file
- Paradox (PC) Database
- PC COM executable (detected in file mode only)
- PC Library Module
- PC Object Module
- PC True Type Font
- PCD Image
- PeachCalc Spreadsheet
- Persuasion Presentation
- PEX Binary Archive (SUN)
- PGP Compressed Data
- PGP Encrypted Data
- PGP Public Keyring
- PGP Secret Keyring
- PGP Signature Certificate
- PGP Signed and Encrypted Data
- PGP Signed Data
- Philips Script
- PKCS #12 (p12) Format
- Plan Perfect
- Portable Bitmap Utilities (PBM)
- Portable Greymap Utilities (PGM)
- Portable Pixmap Utilities (PPM)
- PostScript File
- PostScript Type 1 Font File

- PRIMEWORD
- Program Information File
- PTC Creo
- Q & A for DOS
- Q & A for Windows
- Quadratron Q-One (V1.93J)
- Quadratron Q-One (V2.0)
- Quark Xpress (Macintosh)
- QuickDraw 3D Metafile (3DMF)
- Real Audio
- RealLegal E-Transcript
- Reflex Database (R2D)
- RIFF Device Independent Bitmap
- RIFF MIDI
- RIFF Multimedia Movie
- SAMNA Word IV
- Samsung Electronics JungUm Global format
- SEG-Y Seismic Data format
- Serialized Object Format (SOF) Encapsulation
- SGML
- Simple Vector Format (SVF)
- SMTP document
- SolidWorks
- Sony WAVE64 format
- Star Office Calc Spreadsheet (versions 3-5)
- Star Office Impress Presentation (versions 3-5)
- Star Office Math (versions 3-5)
- Star Office Writer Text (versions 3-5)
- StuffIt Archive (Macintosh)
- SUN vfont definition
- SYLK Spreadsheet
- Symphony Spreadsheet
- Targon Word (V 2.0)
- Unigraphics NX
- Uniplex (V6.01)
- UNIX SHAR Encapsulation
- Usenet format
- Volkswriter

- Vorbis OGG format
- VRML
- VRML 2.0
- WANG PC
- Wang WITA
- WANG WPS Comm.
- Web ARChive (WARC)
- Windows C++ Object Storage
- Windows Journal
- Windows Micrografx Draw (DRW)
- Windows Palette
- Windows scrap file (SHS)
- Wireless Markup Language
- Word Connection
- WordMARC word processor
- WordPerfect General File
- WordStar
- WordStar 6.0
- WordStar 2000
- WriteNow
- Writing Assistant word processor
- X Bitmap (XBM)
- X Image
- X Pixmap (XPM)
- Xerox 860 Comm.
- Xerox DocuWorks
- Xerox Writer word processor
- Yahoo! Messenger chat log
- Zipped Keyhole Markup Language

Appendix B: Character Sets

This section provides information on the handling of character sets in the KeyView suite of products, which includes KeyView Filter SDK, KeyView Export SDK, and KeyView Viewing SDK.

- [Multibyte and Bidirectional Support](#) 125
- [Coded Character Sets](#) 133

Multibyte and Bidirectional Support

The KeyView SDKs can process files that contain multibyte characters. A multibyte character encoding represents a single character with consecutive bytes. KeyView can also process text from files that contain bidirectional text. Bidirectional text contains both Latin-based text which is read from left to right, and text that is read from right to left (Hebrew and Arabic).

[Multibyte and bidirectional support](#), below indicates which character encodings are supported by KeyView for each format.

Multibyte and bidirectional support

Format	Single-byte	Multibyte	Bidirectional
Archive			
7-Zip (7Z)	n/a	n/a	n/a
AD1 Evidence file	n/a	n/a	n/a
ADJ	n/a	n/a	n/a
B1	n/a	n/a	n/a
BinHex (Hqx)	n/a	n/a	n/a
Bzip2 (BZ2)	n/a	n/a	n/a
EnCase – Expert Witness Compression Format (E01)	n/a	n/a	n/a
GZIP (GZ)	n/a	n/a	n/a
ISO (ISO)	n/a	n/a	n/a
Java Archive (JAR)	n/a	n/a	n/a
Legato EMailXtender Archive (EMX)	n/a	n/a	n/a
MacBinary (BIN)	n/a	n/a	n/a
Mac Disk Copy Disk Image (DMG)	n/a	n/a	n/a
Microsoft Backup File (BKF)	n/a	n/a	n/a
Microsoft Cabinet format (CAB)	n/a	n/a	n/a

Multibyte and bidirectional support, continued

Format	Single-byte	Multibyte	Bidirectional
Microsoft Compiled HTML Help (CHM)	n/a	n/a	n/a
Microsoft Compressed Folder (LZH)	n/a	n/a	n/a
PKZip (ZIP)	n/a	n/a	n/a
Microsoft Outlook DBX (DBX)	Y	Y	Y
Microsoft Outlook Offline Storage File (OST)	Y	Y	Y
RAR Archive (RAR)	n/a	n/a	n/a
Tape Archive (TAR)	n/a	n/a	n/a
UNIX Compress (Z)	n/a	n/a	n/a
UUEncoding (UUE)	n/a	n/a	n/a
Windows Scrap File (SHS)	n/a	n/a	n/a
WinZip (ZIP)	n/a	n/a	n/a
Binary			
Executable (EXE)	n/a	n/a	n/a
Link Library (DLL)	n/a	n/a	n/a
Computer-aided Design			
AutoCAD Drawing (DWG)	Y	Y	Y
AutoCAD Drawing Exchange (DXF)	Y	Y	Y
CATIA formats (CAT)	Y	N	N
Microsoft Visio (VSD)	Y	Y	Y
Database			
dBase Database	Y	N	N
Microsoft Access (MDB)	Y	Y	N
Microsoft Project (MPP)	Y	Y	N
Desktop Publishing			
Microsoft Publisher	N	Y	N
Display			
Adobe Portable Document Format (PDF)	Y	Y ¹	Y
Graphics			
Computer Graphics Metafile (CGM)	Y	N	N

Multibyte and bidirectional support, continued

Format	Single-byte	Multibyte	Bidirectional
Corel DRAW (CDR)	n/a	n/a	n/a
DCX Fax System (DCX)	Y	N	N
DICOM – Digital Imaging and Communications in Medicine (DCM)	n/a	n/a	n/a
Encapsulated PostScript (EPS)	Y	N	N
Enhanced Metafile (EMF)	Y	Y	N
Graphic Interchange Format (GIF)	n/a	n/a	n/a
JBIG2	n/a	n/a	n/a
JPEG	n/a	n/a	n/a
JPEG 2000	n/a	n/a	n/a
Lotus AMIDraw Graphics (SDW)	n/a	n/a	n/a
Lotus Pic (PIC)	n/a	n/a	n/a
Macintosh Raster (PICT/PCT)	n/a	n/a	n/a
MacPaint (PNTG)	n/a	n/a	n/a
Microsoft Office Drawing (MSO)	n/a	n/a	n/a
Omni Graffiti (GRAFFLE)	Y	N	N
PC PaintBrush (PCX)	n/a	n/a	n/a
Portable Network Graphics (PNG)	n/a	n/a	n/a
SGI RGB Image (RGB)	n/a	n/a	n/a
Sun Raster Image (RS)	n/a	n/a	n/a
Tagged Image File (TIFF)	Y	N	N
Truevision Targa (TGA)	n/a	n/a	n/a
Windows Animated Cursor (ANI)	n/a	n/a	n/a
Windows Bitmap (BMP)	n/a	n/a	n/a
Windows Icon Cursor (ICO)	n/a	n/a	n/a
Windows Metafile (WMF)	Y	Y	N
WordPerfect Graphics 1 (WPG)	Y	N	N
WordPerfect Graphics 2 (WPG)	Y	N	N
Mail			
Documentum EMCDF Format	Y	Y	Y

Multibyte and bidirectional support, continued

Format	Single-byte	Multibyte	Bidirectional
Domino XML Language (DXL)	Y	Y	N
GroupWise FileSurf	Y	N	N
Legato Extender (ONM)	Y	Y	N
Lotus Notes database (NSF)	Y	Y	Y
Mailbox (MBX)	Y	Y	Y
Microsoft Entourage Database	Y	Y	Y
Microsoft Outlook (MSG)	Y	Y	Y
Microsoft Outlook Express (EML)	Y	Y	Y
Microsoft Outlook iCalendar	Y	Y	Y
Microsoft Outlook for Macintosh	Y	Y	Y
Microsoft Outlook Offline Storage File	Y	Y	Y
Microsoft Outlook Personal File Folders (PST)	Y	Y	Y
Microsoft Outlook vCard Contact	?	?	?
Text Mail (MIME)	Y	Y	Y
Transport Neutral Encapsulation Format	Y	Y	Y
Multimedia			
Advanced Systems Format (ASF)	n/a	n/a	n/a
Audio Interchange File Format (AIFF)	n/a	n/a	n/a
Microsoft Wave Sound (WAV)	n/a	n/a	n/a
MIDI (MID)	n/a	n/a	n/a
MPEG 1 Audio Layer 3 (MP3)	n/a	n/a	n/a
MPEG 1 Video (MPG)	n/a	n/a	n/a
MPEG 2 Audio (MPEGA)	n/a	n/a	n/a
MPEG 4 Audio (MP4)	n/a	n/a	n/a
NeXT/Sun Audio (AU)	n/a	n/a	n/a
QuickTime Movie (QT/MOV)	n/a	n/a	n/a
Windows Video (AVI)	n/a	n/a	n/a
Presentations			
Apple iWork Keynote (GZ)	Y	Y	N

Multibyte and bidirectional support, continued

Format	Single-byte	Multibyte	Bidirectional
Applix Presents (AG)	character set 1252 only	N	N
Corel Presentations (SHW)	character set 1252 only	N	N
Extensible Forms Description Language (XFD)	Y	Y	N
Lotus Freelance Graphics 2 (PRE)	character set 850 only	N	N
Lotus Freelance Graphics (PRZ)	Y	Japanese, Simple Chinese, Traditional Chinese, Thai only	N
Macromedia Flash (SWF)	Y	Y	N
Microsoft OneNote	Y	Y	N
Microsoft PowerPoint PC (PPT)	character set 1252 only	Traditional Chinese only	N
Microsoft PowerPoint Windows (PPT)	Y	Japanese, Simple Chinese, Traditional Chinese, Korean only	Hebrew only
Microsoft PowerPoint Macintosh (PPT)	Y	N	N
Microsoft PowerPoint Windows XML 2007 and 2010 (PPTX)	Y	Y	Y
OASIS Open Document (ODP)	Y	Y	N
OpenOffice Impress (ODP)	Y	Y	N
StarOffice Impress (ODP)	Y	Y	N
Spreadsheets			
Apple iWork Numbers (GZ)	Y	Y	N
Applix Spreadsheets (AS)	character set 1252 only	N	N
Comma Separated Values (CSV)	character set 1252 only	N	N
Corel Quattro Pro (QPW/WB3)	Y	N	N
Data Interchange Format (DIF)	Y	Y	Y ²

Multibyte and bidirectional support, continued

Format	Single-byte	Multibyte	Bidirectional
Lotus 1-2-3 (123)	Y	Y	Y
Lotus 1-2-3 (WK4)	Y	Y	N
Lotus 123 Charts (123)	Y	Y	N
Microsoft Excel Charts (XLS)	Y	Y	N
Microsoft Excel Macintosh (XLS)	Y	N	N
Microsoft Excel Windows (XLS)	Y	Y	Y ²
Microsoft Excel Windows XML 2007 (XLSX)	Y	Y	N
Microsoft Office Excel Binary Format (XLSB)	Y	Y	N
Microsoft Works Spreadsheet (S30/S40)	Y	N	N
OASIS Open Document (ODS)	Y	Y	N
OpenOffice Calc (ODS)	Y	Y	N
StarOffice Calc (ODS)	Y	Y	N
Text and Markup			
ANSI (TXT)	Y	Y	Y ²
ASCII (TXT)	Y	Y	Y ²
HTML (HTM)	Y	Y	Y ^{2, 3}
Microsoft Excel Windows XML 2003	Y	Y	Y
Microsoft Word for Windows XML 2003	Y	Y	Y
Microsoft Visio XML 2003	Y	Y	Y
Rich Text Format (RTF)	Y	Y	Y ³
Unicode HTML	Y	Y	Y ^{2, 3}
Unicode Text (TXT)	Y	Y	Y ²
XHTML	Y	Y	Y ³
XML	Y	Y	Y
Word Processing			
Adobe Maker Interchange Format (MIF)	character set 1252 only	N	N
Apple iChat Log (ICHAT)	Y	Y	N

Multibyte and bidirectional support, continued

Format	Single-byte	Multibyte	Bidirectional
Apple iWork Pages (GZ)	Y	Y	N
Applix Words (AW)	character set 1252 only	N	N
DisplayWrite (IP)	character set 500, 1026 only	N	N
Folio Flat File (FFF)	character set 1252 only	N	N
Founder Chinese E-paper Basic (CEB)	Y	Y	N
Fujitsu Oasys (OA2)	Y	Y	N
Hangul (HWP)	Y	Y	N
Health level7 (HL7)	Y	Y	Y
IBM DCA/RTF (DC)	character sets 500, 1026 only	N	N
JustSystems Ichitaro (JTD)	Y	Y	N
Lotus AMI Pro (SAM)	Y	Simple Chinese, Traditional Chinese, Japanese, Thai only	Y
Lotus AMI Professional Write Plus (AMI)	Y	Simple Chinese, Traditional Chinese, Japanese, Thai only	N
Lotus Word Pro (LWP)	Y	Y	Y ³
Lotus SmartMaster (MWP)	Y	Y	N
Microsoft Word PC (DOC)	character set 1252 only	N	N
Microsoft Word Windows V1-2 (DOC)	Y	N	N
Microsoft Word Windows V6, 7, 8, 95 (DOC)	Y	Y	Hebrew only ³
Microsoft Word Windows V97 through 2003 (DOC)	Y	Y	Y ³
Microsoft Word Windows XML 2007 and 2010 (DOCX)	Y	Y	Y ³
Microsoft Word Macintosh (DOC)	Y	N	Y ³
Microsoft Works (WPS)	Y	Japanese only	N

Multibyte and bidirectional support, continued

Format	Single-byte	Multibyte	Bidirectional
Microsoft Write (WRI)	Y	Japanese only	N
OASIS Open Document (ODT)	Y	Y	N
Omni Outliner (OO3)	Y	Y	N
OpenOffice Writer (ODT)	Y	Y	N
Open Publication Structure eBook (EPUB)	Y	Y	Y
StarOffice Writer (ODT)	Y	Y	N
Skype Log (DBB)	Y	Y (null-terminated charsets)	N
WordPad (RTF)	Y	Y	Y
WordPerfect Linux (WPS)	Y	N	N
WordPerfect Macintosh (WPS)	Y	N	N
WordPerfect Windows (WO)	Y	N	N
XML Paper Specification (XPS)	Y	Y	N
XYWrite Windows (XY4)	character set 1252 only	N	N
Yahoo! Instant Messenger (DAT)	Y	Y (null-terminated charsets)	N

1

Multibyte PDFs are supported, provided the PDF document is created by using either Character ID-keyed (CID) fonts, predefined CJK CMap files, or ToUnicode font encodings, and does not contain embedded fonts. See the Adobe website and the Adobe Acrobat documentation for more information. Any multibyte characters that are not supported are displayed using the replacement character. By default, the replacement character is a question mark (?).

To determine the type of font encodings that are used in a PDF, open the PDF in Adobe Acrobat, and select **File > Document Info > Fonts**. If the Encoding column lists Custom or Embedded encodings, you might encounter problems converting the PDF.

2

The text direction in the output file might not be correct.

3

In Export SDK, a bidirectional right-to-left (RTL) tag is extracted from this format and included in the direction element (<dir=RTL>) of the output.

Coded Character Sets

This section lists which character set you can use to specify the target character set. The coded character sets are enumerated in `kvtypes.h` and defined in the `Filter` class.

Code Character Sets

Coded Character Set	Description	Can be set as target charset?
KVCS_UNKNOWN	Unknown character set	N
KVCS_SJIS	Japanese (uses multibyte encoding), cp932	Y
KVCS_GB	Simplified Chinese (China, Singapore, Malaysia) cp936	Y
KVCS_BIG5	Traditional Chinese (Taiwan, Hong Kong, Macaw) cp950	Y
KVCS_KSC	Korean, cp949	Y
KVCS_1250	Windows Latin 2 (Central Europe)	Y
KVCS_1251	Windows Cyrillic (Slavic)	Y
KVCS_1252	Windows Latin 1 (ANSI)	Y
KVCS_1253	Windows Greek	Y
KVCS_1254	Windows Latin 5 (Turkish)	Y
KVCS_1255	Windows Hebrew	Y
KVCS_1256	Windows Arabic	Y
KVCS_1257	Windows Baltic Rim	Y
KVCS_1258	Windows Vietnamese	Y
KVCS_8859_1	ISO 8859-1 Latin 1 (Western Europe, Latin America)	Y
KVCS_8859_2	ISO 8859-2 Latin 2 (Central Eastern Europe)	Y
KVCS_8859_3	ISO 8859-3 Latin 3 (S.E. Europe)	Y
KVCS_8859_4	ISO 8859-4 Latin 4 (Scandinavia/Baltic)	Y
KVCS_8859_5	ISO 8859-5 Latin/Cyrillic	Y
KVCS_8859_6	ISO 8859-6 Latin/Arabic	Y
KVCS_8859_7	ISO 8859-7 Latin/Greek	Y
KVCS_8859_8	ISO 8859-8 Latin/Hebrew	Y

Code Character Sets, continued

Coded Character Set	Description	Can be set as target charset?
KVCS_8859_9	ISO 8859-9 Latin/Turkish	Y
KVCS_8859_14	ISO 8859-14	Y
KVCS_8859_15	ISO 8859-15	Y
KVCS_437	DOS Latin US	Y
KVCS_737	DOS Greek	Y
KVCS_775	DOS Baltic Rim	Y
KVCS_850	DOS Latin 1	Y
KVCS_851	DOS Greek	Y
KVCS_852	DOS Latin 2	Y
KVCS_855	DOS Cyrillic	Y
KVCS_857	DOS Turkish	Y
KVCS_860	DOS Portuguese	Y
KVCS_861	DOS Icelandic	Y
KVCS_862	DOS Hebrew	Y
KVCS_863	DOS Canadian French	Y
KVCS_864	DOS Arabic	Y
KVCS_865	DOS Nordic	Y
KVCS_866	DOS Cyrillic Russian	Y
KVCS_869	DOS Greek 2	Y
KVCS_874	Thai	Y
KVCS_PDFMACDOC	PDF MAC DOC	N
KVCS_PDFWINDOC	PDF WIN DOC	N
KVCS_STDENC	Adobe Standard Encoding	N
KVCS_PDFDOC	Adobe standard PDF character set	N
KVCS_037	EBCDIC code page 037	Y
KVCS_1026	EBCDIC code page 1026	Y

Code Character Sets, continued

Coded Character Set	Description	Can be set as target charset?
KVCS_500	EBCDIC code page 500	Y
KVCS_875	EBCDIC code page 875	Y
KVCS_LMBCS	Lotus multibyte character set Group 1 and Group 2	N
KVCS_UNICODE	Unicode, UCS-2	Y
KVCS_UTF16	16-bit Unicode transformation format	Y
KVCS_UTF8	8-bit Unicode transformation format	Y
KVCS_UTF7	7-bit Unicode transformation format	Y
KVCS_2022_JP	ISO 2022-JP, Japanese mail and news safe encoding (JIS-7)	N
KVCS_2022_CN	ISO 2022-CN, Chinese mail and news safe encoding	N
KVCS_2022_KR	ISO 2022-KR, Korean mail and news safe encoding	N
KVCS_WP6X	Word Perfect 6.x and higher character mapping	N
KVCS_10000	Western European (Macintosh)	Y
KVCS_KSC5601	Unified Hangul	Y
KVCS_GB2312	Simplified Chinese (China, Singapore, Hong Kong)	Y
KVCS_GB12345	Traditional Chinese (China) - analogue of GB2312	Y
KVCS_CNS11643	Traditional Chinese - Taiwan. Supplement to Big5	Y
KVCS_JIS0201	Japanese - contains ASCII character set (JIS-Roman)	N
KVCS_JIS0212	Japanese. Supplement to JIS0208.	Y
KVCS_EUC_JP	Japanese Extended UNIX Code	Y
KVCS_EUC_GB	Simplified Chinese Extended UNIX Code	Y
KVCS_EUC_BIG5	Traditional Chinese Extended UNIX Code	N
KVCS_EUC_KSC	Korean Extended UNIX Code	N
KVCS_424	EBCDIC Hebrew	N
KVCS_856	PC Hebrew (old)	N
KVCS_1006	IBM AIX Pakistan (Urdu)	N
KVCS_KOI8R	Cyrillic (Russian)	Y

Code Character Sets, continued

Coded Character Set	Description	Can be set as target charset?
KVCS_PDF_JAPAN1	Adobe-Japan1-2 character collection	N
KVCS_PDF_KOREA1	Adobe-Korea1-0 character collection	N
KVCS_PDF_GB1	Adobe-GB1-3 character collection	N
KVCS_PDF_CNS1	Adobe-CNS1-2 character collection	N
KVCS_2022_JP_8	ISO 2022-JP, Japanese mail and news safe encoding (JIS8)	N
KVCS_720	Arabic DOS-720	Y
KVCS_VISCII	Vietnamese VISCII	Y
KVCS_8859_10	ISO 8859-10 (Latin 6 Nordic)	Y ¹
KVCS_8859_13	ISO 8859-13 (Latin 7 Baltic)	Y ¹
KVCS_57002	ISCII Devanagari (x-iscii-de)	Y ¹
KVCS_57003	ISCII Bengali (x-iscii-be)	Y ¹
KVCS_57004	ISCII Tamil (x-iscii-ta)	Y ¹
KVCS_57005	ISCII Telugu (x-iscii-te)	Y ¹
KVCS_57006	ISCII Assamese (x-iscii-as)	Y ¹
KVCS_57007	ISCII Oriya (x-iscii-or)	Y ¹
KVCS_57008	ISCII Kannada (x-iscii-ka)	Y ¹
KVCS_57009	ISCII Malayalam (x-iscii-ma)	Y ¹
KVCS_57010	ISCII Gujarathi (x-iscii-gu)	Y ¹
KVCS_57011	ISCII Panjabi (x-iscii-pa)	Y ¹
KVCS_GB18030b2	Reserved for internal use	n/a
KVCS_GB18030	GB18030 (Chinese 4-byte character set)	Y
KVCS_8859_11	ISO 8859-11 (Thai)	Y
KVCS_8859_16	ISO 8859-16 (Latin-10 South-Eastern Europe)	Y
KVCS_ARABICMAC	Arabic Mac (x-mac-arabic)	Y
KVCS_KOI8U	Cyrillic (KOI8U Ukrainian)	Y
KVCS_HZGB2312	The 7-bit representation of GB 2312 / RFC 1842	n/a

1

The character set cannot be forced as output in Export SDK and Viewing SDK because the character set is not supported by the major browsers.

Appendix C: File Formats and Extensions

This section lists the KeyView file format numbers and their associated file extensions.

- [File Format and Extension Table](#) 138

File Format and Extension Table

This section lists the KeyView file format codes and the file extensions that they are most commonly associated with.

NOTE: This is not a complete list of file extensions. KeyView returns format codes based on file content, which cannot always be predicted from the file extension. Some file extensions might also be associated with multiple format numbers.

KeyView file formats and extensions

Format Name	Format Number	Format Description	Associated File Extension
AES_Multiplus_Comm_Fmt	1	Multiplus (AES)	PTF
ASCII_Text_Fmt	2	Text	
MSDOS_Batch_File_Fmt	3	MS-DOS Batch File	BAT
Applix_Alis_Fmt	4	APPLIX ASTERIX	AX
BMP_Fmt	5	Windows Bitmap	BMP
CT_DEF_Fmt	6	Convergent Technologies DEF Comm. Format	
Corel_Draw_Fmt	7	Corel Draw	CDR
CGM_ClearText_Fmt	8	Computer Graphics Metafile (CGM)	CGM ¹
CGM_Binary_Fmt	9	Computer Graphics Metafile (CGM)	CGM ¹
CGM_Character_Fmt	10	Computer Graphics Metafile (CGM)	CGM ¹
Word_Connection_Fmt	11	Word Connection	CN
COMET_TOP_Word_Fmt	12	COMET TOP	

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
CEOwrite_Fmt	13	CEOwrite	CW
DSA101_Fmt	14	DSA101 (Honeywell Bull)	
DCA_RFT_Fmt	15	DCA-RFT (IBM Revisable Form)	RFT
CDA_DDIF_Fmt	16	CDA / DDIF	
DG_CDS_Fmt	17	DG Common Data Stream (CDS)	CDS
Micrografx_Draw_Fmt	18	Windows Draw (Micrografx)	DRW
Data_Point_VistaWord_Fmt	19	Vistaword	
DECdx_Fmt	20	DECdx	DX
Enable_WP_Fmt	21	Enable Word Processing	WPF
EPSF_Fmt	22	Encapsulated PostScript	EPS ¹
Preview_EPSF_Fmt	23	Encapsulated PostScript	EPS ¹
MS_Executable_Fmt	24	MSDOS/Windows Program	EXE
G31D_Fmt	25	CCITT G3 1D	
GIF_87a_Fmt	26	Graphics Interchange Format (GIF87a)	GIF ¹
GIF_89a_Fmt	27	Graphics Interchange Format (GIF89a)	GIF ¹
HP_Word_PC_Fmt	28	HP Word PC	HW
IBM_1403_LinePrinter_Fmt	29	IBM 1403 Line Printer	I4
IBM_DCF_Script_Fmt	30	DCF Script	IC
IBM_DCA_FFT_Fmt	31	DCA-FFT (IBM Final Form)	IF
Interleaf_Fmt	32	Interleaf	
GEM_Image_Fmt	33	GEM Bit Image	IMG
IBM_Display_Write_Fmt	34	Display Write	IP

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
Sun_Raster_Fmt	35	Sun Raster	RAS
Ami_Pro_Fmt	36	Lotus Ami Pro	SAM
Ami_Pro_StyleSheet_Fmt	37	Lotus Ami Pro Style Sheet	
MORE_Fmt	38	MORE Database MAC	
Lyrix_Fmt	39	Lyrix Word Processing	
MASS_11_Fmt	40	MASS-11	M1
MacPaint_Fmt	41	MacPaint	PNTG
MS_Word_Mac_Fmt	42	Microsoft Word for Macintosh	DOC ¹
SmartWare_II_Comm_Fmt	43	SmartWare II	
MS_Word_Win_Fmt	44	Microsoft Word for Windows	DOC ¹
Multimate_Fmt	45	MultiMate	MM ¹
Multimate_Fnote_Fmt	46	MultiMate Footnote File	FNX ¹
Multimate_Adv_Fmt	47	MultiMate Advantage	
Multimate_Adv_Fnote_Fmt	48	MultiMate Advantage Footnote File	
Multimate_Adv_II_Fmt	49	MultiMate Advantage II	MM ¹
Multimate_Adv_II_Fnote_Fmt	50	MultiMate Advantage II Footnote File	FNX ¹
Multiplan_PC_Fmt	51	Multiplan (PC)	
Multiplan_Mac_Fmt	52	Multiplan (Mac)	
MS_RTF_Fmt	53	Rich Text Format (RTF)	RTF
MS_Word_PC_Fmt	54	Microsoft Word for PC	DOC ¹
MS_Word_PC_StyleSheet_Fmt	55	Microsoft Word for PC Style Sheet	DOC ¹
MS_Word_PC_Glossary_Fmt	56	Microsoft Word for PC Glossary	DOC ¹

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
MS_Word_PC_Driver_Fmt	57	Microsoft Word for PC Driver	DOC ¹
MS_Word_PC_Misc_Fmt	58	Microsoft Word for PC Miscellaneous File	DOC ¹
NBI_Async_Archive_Fmt	59	NBI Async Archive Format	
Navy_DIF_Fmt	60	Navy DIF	ND
NBI_Net_Archive_Fmt	61	NBI Net Archive Format	NN
NIOS_TOP_Fmt	62	NIOS TOP	
FileMaker_Mac_Fmt	63	Filemaker MAC	FP5, FP7
ODA_Q1_11_Fmt	64	ODA / ODIF	OD ¹
ODA_Q1_12_Fmt	65	ODA / ODIF	OD ¹
OLIDIF_Fmt	66	OLIDIF (Olivetti)	
Office_Writer_Fmt	67	Office Writer	OW
PC_Paintbrush_Fmt	68	PC Paintbrush Graphics (PCX)	PCX
CPT_Comm_Fmt	69	CPT	
Lotus_PIC_Fmt	70	Lotus PIC	PIC
Mac_PICT_Fmt	71	QuickDraw Picture	PCT
Philips_Script_Word_Fmt	72	Philips Script	
PostScript_Fmt	73	PostScript	PS
PRIMEWORD_Fmt	74	PRIMEWORD	
Quadratron_Q_One_v1_Fmt	75	Q-One V1.93J	Q1 ¹ , QX ¹
Quadratron_Q_One_v2_Fmt	76	Q-One V2.0	Q1 ¹ , QX ¹
SAMNA_Word_IV_Fmt	77	SAMNA Word	SAM

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
Ami_Pro_Draw_Fmt	78	Lotus Ami Pro Draw	SDW
SYLK_Spreadsheet_Fmt	79	SYLK	
SmartWare_II_WP_Fmt	80	SmartWare II	
Symphony_Fmt	81	Symphony	WR1
Targa_Fmt	82	Targa	TGA
TIFF_Fmt	83	TIFF	TIF, TIFF
Targon_Word_Fmt	84	Targon Word	TW
Uniplex_Ucalc_Fmt	85	Uniplex Ucalc	SS
Uniplex_WP_Fmt	86	Uniplex	UP
MS_Word_UNIX_Fmt	87	Microsoft Word UNIX	DOC ¹
WANG_PC_Fmt	88	WANG PC	
WordERA_Fmt	89	WordERA	
WANG_WPS_Comm_Fmt	90	WANG WPS	WF
WordPerfect_Mac_Fmt	91	WordPerfect MAC	WPM, WPD ¹
WordPerfect_Fmt	92	WordPerfect	WO, WPD ¹
WordPerfect_VAX_Fmt	93	WordPerfect VAX	WPD ¹
WordPerfect_Macro_Fmt	94	WordPerfect Macro	
WordPerfect_Dictionary_Fmt	95	WordPerfect Spelling Dictionary	
WordPerfect_Thesaurus_Fmt	96	WordPerfect Thesaurus	
WordPerfect_Resource_Fmt	97	WordPerfect Resource File	
WordPerfect_Driver_	98	WordPerfect Driver	

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
Fmt			
WordPerfect_Cfg_Fmt	99	WordPerfect Configuration File	
WordPerfect_Hyphenation_Fmt	100	WordPerfect Hyphenation Dictionary	
WordPerfect_Misc_Fmt	101	WordPerfect Miscellaneous File	WPD ¹
WordMARC_Fmt	102	WordMARC	WM, PW
Windows_Metafile_Fmt	103	Windows Metafile	WMF ¹
Windows_Metafile_NoHdr_Fmt	104	Windows Metafile (no header)	WMF ¹
SmartWare_II_DB_Fmt	105	SmartWare II	
WordPerfect_Graphics_Fmt	106	WordPerfect Graphics	WPG, QPG
WordStar_Fmt	107	WordStar	WS
WANG_WITA_Fmt	108	WANG WITA	WT
Xerox_860_Comm_Fmt	109	Xerox 860	
Xerox_Writer_Fmt	110	Xerox Writer	
DIF_SpreadSheet_Fmt	111	Data Interchange Format (DIF)	DIF
Enable_Spreadsheet_Fmt	112	Enable Spreadsheet	SSF
SuperCalc_Fmt	113	Supercalc	CAL
UltraCalc_Fmt	114	UltraCalc	
SmartWare_II_SS_Fmt	115	SmartWare II	
SOF_Encapsulation_Fmt	116	Serialized Object Format (SOF)	SOF

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
PowerPoint_Win_Fmt	117	PowerPoint PC	PPT ¹
PowerPoint_Mac_Fmt	118	PowerPoint MAC	PPT ¹
PowerPoint_95_Fmt	119	PowerPoint 95	PPT ¹
PowerPoint_97_Fmt	120	PowerPoint 97	PPT ¹
PageMaker_Mac_Fmt	121	PageMaker for Macintosh	
PageMaker_Win_Fmt	122	PageMaker for Windows	
MS_Works_Mac_WP_Fmt	123	Microsoft Works for MAC	
MS_Works_Mac_DB_Fmt	124	Microsoft Works for MAC	
MS_Works_Mac_SS_Fmt	125	Microsoft Works for MAC	
MS_Works_Mac_Comm_Fmt	126	Microsoft Works for MAC	
MS_Works_DOS_WP_Fmt	127	Microsoft Works for DOS	WPS ¹
MS_Works_DOS_DB_Fmt	128	Microsoft Works for DOS	WDB ¹
MS_Works_DOS_SS_Fmt	129	Microsoft Works for DOS	
MS_Works_Win_WP_Fmt	130	Microsoft Works for Windows	WPS ¹
MS_Works_Win_DB_Fmt	131	Microsoft Works for Windows	WDB ¹
MS_Works_Win_SS_Fmt	132	Microsoft Works for Windows	S30, S40
PC_Library_Fmt	133	DOS/Windows Object Library	
MacWrite_Fmt	134	MacWrite	

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
MacWrite_II_Fmt	135	MacWrite II	
Freehand_Fmt	136	Freehand MAC	
Disk_Doubler_Fmt	137	Disk Doubler	
HP_GL_Fmt	138	HP Graphics Language	HPGL
FrameMaker_Fmt	139	FrameMaker	FM, FRM
FrameMaker_Book_Fmt	140	FrameMaker	BOOK
Maker_Markup_Language_Fmt	141	Maker Markup Language	
Maker_Interchange_Fmt	142	Maker Interchange Format (MIF)	MIF
JPEG_File_Interchange_Fmt	143	Interchange Format	JPG, JPEG
Reflex_Fmt	144	Reflex	
Framework_Fmt	145	Framework	
Framework_II_Fmt	146	Framework II	FW3
Paradox_Fmt	147	Paradox	DB
MS_Windows_Write_Fmt	148	Windows Write	WRI
Quattro_Pro_DOS_Fmt	149	Quattro Pro for DOS	
Quattro_Pro_Win_Fmt	150	Quattro Pro for Windows	WB2, WB3
Persuasion_Fmt	151	Persuasion	
Windows_Icon_Fmt	152	Windows Icon Format	ICO
Windows_Cursor_Fmt	153	Windows Cursor	CUR
MS_Project_Activity_Fmt	154	Microsoft Project	MPP ¹
MS_Project_Resource_Fmt	155	Microsoft Project	MPP ¹

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
MS_Project_Calc_Fmt	156	Microsoft Project	MPP ¹
PKZIP_Fmt	157	ZIP Archive	ZIP
Quark_Xpress_Fmt	158	Quark Xpress MAC	
ARC_PAK_Archive_Fmt	159	PAK/ARC Archive	ARC, PAK
MS_Publisher_Fmt	160	Microsoft Publisher	PUB ¹
PlanPerfect_Fmt	161	PlanPerfect	
WordPerfect_Auxiliary_Fmt	162	WordPerfect auxiliary file	WPW
MS_WAVE_Audio_Fmt	163	Microsoft Wave	WAV
MIDI_Audio_Fmt	164	MIDI	MID, MIDI
AutoCAD_DXF_Binary_Fmt	165	AutoCAD DXF	DXF ¹
AutoCAD_DXF_Text_Fmt	166	AutoCAD DXF	DXF ¹
dBase_Fmt	167	dBase	DBF
OS_2_PM_Metafile_Fmt	168	OS/2 PM Metafile	MET
Lasergraphics_Language_Fmt	169	Lasergraphics Language	
AutoShade_Rendering_Fmt	170	AutoShade Rendering	
GEM_VDI_Fmt	171	GEM VDI	VDI
Windows_Help_Fmt	172	Windows Help File	HLP
Volkswriter_Fmt	173	Volkswriter	VW4
Ability_WP_Fmt	174	Ability	
Ability_DB_Fmt	175	Ability	
Ability_SS_Fmt	176	Ability	

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
Ability_Comm_Fmt	177	Ability	
Ability_Image_Fmt	178	Ability	
XyWrite_Fmt	179	XYWrite / Nota Bene	XY4
CSV_Fmt	180	CSV (Comma Separated Values)	CSV
IBM_Writing_Assistant_Fmt	181	IBM Writing Assistant	IWA
WordStar_2000_Fmt	182	WordStar 2000	WS2
HP_PCL_Fmt	183	HP Printer Control Language	PCL
UNIX_Exe_PreSysV_VAX_Fmt	184	Unix Executable (PDP-11/pre-System V VAX)	
UNIX_Exe_Basic_16_Fmt	185	Unix Executable (Basic-16)	
UNIX_Exe_x86_Fmt	186	Unix Executable (x86)	
UNIX_Exe_iAPX_286_Fmt	187	Unix Executable (iAPX 286)	
UNIX_Exe_MC68k_Fmt	188	Unix Executable (MC680x0)	
UNIX_Exe_3B20_Fmt	189	Unix Executable (3B20)	
UNIX_Exe_WE32000_Fmt	190	Unix Executable (WE32000)	
UNIX_Exe_VAX_Fmt	191	Unix Executable (VAX)	
UNIX_Exe_Bell_5_Fmt	192	Unix Executable (Bell 5.0)	
UNIX_Obj_VAX_Demand_Fmt	193	Unix Object Module (VAX Demand)	
UNIX_Obj_MS8086_Fmt	194	Unix Object Module (old MS 8086)	
UNIX_Obj_Z8000_Fmt	195	Unix Object Module (Z8000)	

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
AU_Audio_Fmt	196	NeXT/Sun Audio Data	AU
NeWS_Font_Fmt	197	NeWS bitmap font	
cpio_Archive_CRChdr_Fmt	198	cpio archive (CRC Header)	
cpio_Archive_CHRhdr_Fmt	199	cpio archive (CHR Header)	
PEX_Binary_Archive_Fmt	200	SUN PEX Binary Archive	
Sun_vfont_Fmt	201	SUN vfont Definition	
Curses_Screen_Fmt	202	Curses Screen Image	
UUEncoded_Fmt	203	UU encoded	UUE
WriteNow_Fmt	204	WriteNow MAC	
PC_Obj_Fmt	205	DOS/Windows Object Module	
Windows_Group_Fmt	206	Windows Group	
TrueType_Font_Fmt	207	TrueType Font	TTF
Windows_PIF_Fmt	208	Program Information File (PIF)	PIF
MS_COM_Executable_Fmt	209	PC (.COM)	COM
StuftIt_Fmt	210	StuftIt (MAC)	HQX
PeachCalc_Fmt	211	PeachCalc	
Wang_GDL_Fmt	212	WANG Office GDL Header	
Q_A_DOS_Fmt	213	Q & A for DOS	
Q_A_Win_Fmt	214	Q & A for Windows	JW
WPS_PLUS_Fmt	215	WPS-PLUS	WPL
DCX_Fmt	216	DCX FAX Format(PCX images	DCX
OLE_Fmt	217	OLE Compound Document	OLE
EBCDIC_Fmt	218	EBCDIC Text	
DCS_Fmt	219	DCS	

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
UNIX_SHAR_Fmt	220	SHAR	SHAR
Lotus_Notes_BitMap_Fmt	221	Lotus Notes Bitmap	
Lotus_Notes_CDF_Fmt	222	Lotus Notes CDF	CDF
Compress_Fmt	223	Unix Compress	Z
GZ_Compress_Fmt	224	GZ Compress	GZ ¹
TAR_Fmt	225	TAR	TAR
ODIF_FOD26_Fmt	226	ODA / ODIF	F26
ODIF_FOD36_Fmt	227	ODA / ODIF	F36
ALIS_Fmt	228	ALIS	
Envoy_Fmt	229	Envoy	EVY
PDF_Fmt	230	Portable Document Format	PDF
BinHex_Fmt	231	BinHex	HQX
SMTP_Fmt	232	SMTP	SMTP
MIME_Fmt	233	MIME ²	EML, MBX
USENET_Fmt	234	USENET	
SGML_Fmt	235	SGML	SGML
HTML_Fmt	236	HTML	HTM ¹ , HTML ¹
ACT_Fmt	237	ACT	ACT
PNG_Fmt	238	Portable Network Graphics (PNG)	PNG
MS_Video_Fmt	239	Video for Windows (AVI)	AVI
Windows_Animated_Cursor_Fmt	240	Windows Animated Cursor	ANI
Windows_CPP_Obj_Storage_Fmt	241	Windows C++ Object Storage	
Windows_Palette_Fmt	242	Windows Palette	PAL
RIFF_DIB_Fmt	243	RIFF Device Independent Bitmap	

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
RIFF_MIDI_Fmt	244	RIFF MIDI	RMI
RIFF_Multimedia_Movie_Fmt	245	RIFF Multimedia Movie	
MPEG_Fmt	246	MPEG Movie	MPG, MPEG ¹
QuickTime_Fmt	247	QuickTime Movie, MPEG-4 Audio	MOV, QT, MP4
AIFF_Fmt	248	Audio Interchange File Format (AIFF)	AIF, AIFF
Amiga_MOD_Fmt	249	Amiga MOD	MOD
Amiga_IFF_8SVX_Fmt	250	Amiga IFF (8SVX) Sound	IFF
Creative_Voice_Audio_Fmt	251	Creative Voice (VOC)	VOC
AutoDesk_Animator_FLI_Fmt	252	AutoDesk Animator FLIC	FLI
AutoDesk_AnimatorPro_FLC_Fmt	253	AutoDesk Animator Pro FLIC	FLC
Compactor_Archive_Fmt	254	Compactor / Compact Pro	
VRML_Fmt	255	VRML	WRL
QuickDraw_3D_Metafile_Fmt	256	QuickDraw 3D Metafile	
PGP_Secret_Keyring_Fmt	257	PGP Secret Keyring	
PGP_Public_Keyring_Fmt	258	PGP Public Keyring	
PGP_Encrypted_Data_Fmt	259	PGP Encrypted Data	
PGP_Signed_Data_Fmt	260	PGP Signed Data	
PGP_SignedEncrypted_Data_Fmt	261	PGP Signed and Encrypted Data	

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
PGP_Sign_Certificate_Fmt	262	PGP Signature Certificate	
PGP_Compressed_Data_Fmt	263	PGP Compressed Data	
PGP_ASCII_Public_Keyring_Fmt	264	ASCII-armored PGP Public Keyring	
PGP_ASCII_Encoded_Fmt	265	ASCII-armored PGP encoded	PGP ¹
PGP_ASCII_Signed_Fmt	266	ASCII-armored PGP encoded	PGP ¹
OLE_DIB_Fmt	267	OLE DIB object	
SGL_Image_Fmt	268	SGL Image	RGB
Lotus_ScreenCam_Fmt	269	Lotus ScreenCam	
MPEG_Audio_Fmt	270	MPEG Audio	MPEGA
FTP_Software_Session_Fmt	271	FTP Session Data	STE
Netscape_Bookmark_File_Fmt	272	Netscape Bookmark File	HTM ¹
Corel_Draw_CMX_Fmt	273	Corel CMX	CMX
AutoDesk_DWG_Fmt	274	AutoDesk Drawing (DWG)	DWG
AutoDesk_WHIP_Fmt	275	AutoDesk WHIP	WHP
Macromedia_Director_Fmt	276	Macromedia Director	DCR
Real_Audio_Fmt	277	Real Audio	RM
MSDOS_Device_Driver_Fmt	278	MSDOS Device Driver	SYS
Micrografx_Designer_Fmt	279	Micrografx Designer	DSF

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
SVF_Fmt	280	Simple Vector Format (SVF)	SVF
Applix_Words_Fmt	281	Applix Words	AW
Applix_Graphics_Fmt	282	Applix Graphics	AG
MS_Access_Fmt	283	Microsoft Access	MDB ¹
MS_Access_95_Fmt	284	Microsoft Access 95	MDB ¹
MS_Access_97_Fmt	285	Microsoft Access 97	MDB ¹
MacBinary_Fmt	286	MacBinary	BIN
Apple_Single_Fmt	287	Apple Single	
Apple_Double_Fmt	288	Apple Double	
Enhanced_Metafile_Fmt	289	Enhanced Metafile	EMF
MS_Office_Drawing_Fmt	290	Microsoft Office Drawing	
XML_Fmt	291	XML	XML ¹
DeVice_Independent_Fmt	292	DeVice Independent file (DVI)	DVI
Unicode_Fmt	293	Unicode	UNI
Lotus_123_Worksheet_Fmt	294	Lotus 1-2-3	WK1 ¹
Lotus_123_Format_Fmt	295	Lotus 1-2-3 Formatting	FM3
Lotus_123_97_Fmt	296	Lotus 1-2-3 97	WK1 ¹
Lotus_Word_Pro_96_Fmt	297	Lotus Word Pro 96	LWP ¹
Lotus_Word_Pro_97_Fmt	298	Lotus Word Pro 97	LWP ¹
Freelance_DOS_Fmt	299	Lotus Freelance for DOS	
Freelance_Win_Fmt	300	Lotus Freelance for Windows	PRE

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
Freelance_OS2_Fmt	301	Lotus Freelance for OS/2	PRS
Freelance_96_Fmt	302	Lotus Freelance 96	PRZ ¹
Freelance_97_Fmt	303	Lotus Freelance 97	PRZ ¹
MS_Word_95_Fmt	304	Microsoft Word 95	DOC ¹
MS_Word_97_Fmt	305	Microsoft Word 97	DOC ¹
Excel_Fmt	306	Microsoft Excel	XLS ¹
Excel_Chart_Fmt	307	Microsoft Excel	XLS ¹
Excel_Macro_Fmt	308	Microsoft Excel	XLS ¹
Excel_95_Fmt	309	Microsoft Excel 95	XLS ¹
Excel_97_Fmt	310	Microsoft Excel 97	XLS ¹
Corel_Presentations_Fmt	311	Corel Presentations	XFD, XFDL
Harvard_Graphics_Fmt	312	Harvard Graphics	
Harvard_Graphics_Chart_Fmt	313	Harvard Graphics Chart	CH3, CHT
Harvard_Graphics_Symbol_Fmt	314	Harvard Graphics Symbol File	SY3
Harvard_Graphics_Cfg_Fmt	315	Harvard Graphics Configuration File	
Harvard_Graphics_Palette_Fmt	316	Harvard Graphics Palette	
Lotus_123_R9_Fmt	317	Lotus 1-2-3 Release 9	
Applix_Spreadsheets_Fmt	318	Applix Spreadsheets	AS
MS_Pocket_Word_Fmt	319	Microsoft Pocket Word	PWD, DOC ¹
MS_DIB_Fmt	320	MS Windows Device Independent Bitmap	
MS_Word_2000_Fmt	321	Microsoft Word 2000	DOC ¹

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
Excel_2000_Fmt	322	Microsoft Excel 2000	XLS ¹
PowerPoint_2000_Fmt	323	Microsoft PowerPoint 2000	PPT
MS_Access_2000_Fmt	324	Microsoft Access 2000	MDB ¹ , MPP ¹
MS_Project_4_Fmt	325	Microsoft Project 4	MPP ¹
MS_Project_41_Fmt	326	Microsoft Project 4.1	MPP ¹
MS_Project_98_Fmt	327	Microsoft Project 98	MPP ¹
Folio_Flat_Fmt	328	Folio Flat File	FFF
HWP_Fmt	329	HWP(Arae-Ah Hangul)	HWP
ICHITARO_Fmt	330	ICHITARO V4-10	
IS_XML_Fmt	331	Extended or Custom XML	XML ¹
Oasys_Fmt	332	Oasys format	OA2, OA3
PBM_ASC_Fmt	333	Portable Bitmap Utilities ASCII Format	
PBM_BIN_Fmt	334	Portable Bitmap Utilities Binary Format	
PGM_ASC_Fmt	335	Portable Greymap Utilities ASCII Format	
PGM_BIN_Fmt	336	Portable Greymap Utilities Binary Format	PGM
PPM_ASC_Fmt	337	Portable Pixmap Utilities ASCII Format	
PPM_BIN_Fmt	338	Portable Pixmap Utilities Binary Format	
XBM_Fmt	339	X Bitmap Format	XBM
XPM_Fmt	340	X Pixmap Format	XPM
FPX_Fmt	341	FPX Format	FPX
PCD_Fmt	342	PCD Format	PCD

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
MS_Visio_Fmt	343	Microsoft Visio	VSD
MS_Project_2000_Fmt	344	Microsoft Project 2000	MPP ¹
MS_Outlook_Fmt	345	Microsoft Outlook	MSG, OFT
ELF_Relocatable_Fmt	346	ELF Relocatable	O
ELF_Executable_Fmt	347	ELF Executable	
ELF_Dynamic_Lib_Fmt	348	ELF Dynamic Library	SO
MS_Word_XML_Fmt	349	Microsoft Word 2003 XML	XML ¹
MS_Excel_XML_Fmt	350	Microsoft Excel 2003 XML	XML ¹
MS_Visio_XML_Fmt	351	Microsoft Visio 2003 XML	VDX
SO_Text_XML_Fmt	352	StarOffice Text XML	SXW ¹ , ODT ¹
SO_Spreadsheet_XML_Fmt	353	StarOffice Spreadsheet XML	SXC ¹ , ODS ¹
SO_Presentation_XML_Fmt	354	StarOffice Presentation XML	SXI ¹ , SXP ¹ , ODP ¹
XHTML_Fmt	355	XHTML	XML ¹
MS_OutlookPST_Fmt	356	Microsoft Outlook PST	PST
RAR_Fmt	357	RAR	RAR
Lotus_Notes_NSF_Fmt	358	IBM Lotus Notes Database NSF/NTF	NSF
Macromedia_Flash_Fmt	359	SWF	SWF
MS_Word_2007_Fmt	360	Microsoft Word 2007 XML	DOCX, DOTX
MS_Excel_2007_Fmt	361	Microsoft Excel 2007 XML	XLSX, XLTX
MS_PPT_2007_Fmt	362	Microsoft PPT 2007 XML	PPTX, POTX, PPSX
OpenPGP_Fmt	363	OpenPGP Message Format (with new	PGP

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
		packet format)	
Intergraph_V7_DGN_Fmt	364	Intergraph Standard File Format (ISFF) V7 DGN (non-OLE)	DGN ¹
MicroStation_V8_DGN_Fmt	365	MicroStation V8 DGN (OLE)	DGN ¹
MS_Word_Macro_2007_Fmt	366	Microsoft Word Macro 2007 XML	DOCM, DOTM
MS_Excel_Macro_2007_Fmt	367	Microsoft Excel Macro 2007 XML	XLSM, XLTM, XLAM
MS_PPT_Macro_2007_Fmt	368	Microsoft PPT Macro 2007 XML	PPTM, POTM, PPSM, PPAM
LZH_Fmt	369	LHA Archive	LZH, LHA
Office_2007_Fmt	370	Office 2007 document	XLSB
MS_XPS_Fmt	371	Microsoft XML Paper Specification (XPS)	XPS
Lotus_Domino_DXL_Fmt	372	IBM Lotus representation of Domino design elements in XML format	DXL
ODF_Text_Fmt	373	ODF Text	ODT ¹ , SXW ¹ , STW
ODF_Spreadsheet_Fmt	374	ODF Spreadsheet	ODS ¹ , SXC ¹ , STC
ODF_Presentation_Fmt	375	ODF Presentation	SXD ¹ , SXI ¹ , ODG ¹ , , ODP ¹
Legato_Extender_ONM_Fmt	376	Legato Extender Native Message ONM	ONM
bin_Unknown_Fmt	377	n/a	
TNEF_Fmt	378	Transport Neutral Encapsulation Format (TNEF)	various
CADAM_Drawing_Fmt	379	CADAM Drawing	CDD
CADAM_Drawing_Overlay_Fmt	380	CADAM Drawing Overlay	CDO
NURSTOR_	381	NURSTOR Drawing	NUR

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
Drawing_Fmt			
HP_GLP_Fmt	382	HP Graphics Language (Plotter)	HPG
ASF_Fmt	383	Advanced Systems Format (ASF)	ASF
WMA_Fmt	384	Window Media Audio Format (WMA)	WMA
WMV_Fmt	385	Window Media Video Format (WMV)	WMV
EMX_Fmt	386	Legato EMailXtender Archives Format (EMX)	EMX
Z7Z_Fmt	387	7 Zip Format(7z)	7Z
MS_Excel_Binary_2007_Fmt	388	Microsoft Excel Binary 2007	XLSB
CAB_Fmt	389	Microsoft Cabinet File (CAB)	CAB
CATIA_Fmt	390	CATIA Formats (CAT*)	CAT ³
YIM_Fmt	391	Yahoo Instant Messenger History	DAT ¹
ODF_Drawing_Fmt	392	ODF Drawing	SXD ¹ , SX ¹ , ODG ¹
Founder_CEB_Fmt	393	Founder Chinese E-paper Basic (ceb)	CEB
QPW_Fmt	394	Quattro Pro 9+ for Windows	QPW
MHT_Fmt	395	MHT format ²	MHT
MDI_Fmt	396	Microsoft Document Imaging Format	MDI
GRV_Fmt	397	Microsoft Office Groove Format	GRV
IWWP_Fmt	398	Apple iWork Pages format	PAGES, GZ ¹
IWSS_Fmt	399	Apple iWork Numbers format	NUMBERS, GZ ¹
IWPG_Fmt	400	Apple iWork Keynote format	KEY, GZ ¹
BKF_Fmt	401	Windows Backup File	BKF
MS_Access_2007_Fmt	402	Microsoft Access 2007	ACCDB
ENT_Fmt	403	Microsoft Entourage Database Format	
DMG_Fmt	404	Mac Disk Copy Disk Image File	

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
CWK_Fmt	405	AppleWorks File	
OO3_Fmt	406	Omni Outliner File	OO3
OPML_Fmt	407	Omni Outliner File	OPML
Omni_Graffle_XML_File	408	Omni Graffle XML File	GRAFFLE
PSD_Fmt	409	Photoshop Document	PSD
Apple_Binary_PList_Fmt	410	Apple Binary Property List format	
Apple_iChat_Fmt	411	Apple iChat format	
OOUTLINE_Fmt	412	OOutliner File	OOUTLINE
BZIP2_Fmt	413	Bzip 2 Compressed File	BZ2
ISO_Fmt	414	ISO-9660 CD Disc Image Format	ISO
DocuWorks_Fmt	415	DocuWorks Format	XDW
RealMedia_Fmt	416	RealMedia Streaming Media	RM, RA
AC3Audio_Fmt	417	AC3 Audio File Format	AC3
NEF_Fmt	418	Nero Encrypted File	NEF
SolidWorks_Fmt	419	SolidWorks Format Files	SLDASM, SLDPRT, SLDDRW
XFDL_Fmt	420	Extensible Forms Description Language	XFDL, XFD
Apple_XML_PList_Fmt	421	Apple XML Property List format	
OneNote_Fmt	422	OneNote Note Format	ONE
Dicom_Fmt	424	Digital Imaging and Communications in Medicine	DCM
EnCase_Fmt	425	Expert Witness Compression Format (EnCase)	E01, L01, Lx01
Scrap_Fmt	426	Shell Scrap Object File	SHS
MS_Project_2007_Fmt	427	Microsoft Project 2007	MPP ¹

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
MS_Publisher_98_Fmt	428	Microsoft Publisher 98/2000/2002/2003/2007/	PUB ¹
Skype_Fmt	429	Skype Log File	DBB
HL7_Fmt	430	Health level7 message	HL7
MS_OutlookOST_Fmt	431	Microsoft Outlook OST	OST
Epub_Fmt	432	Electronic Publication	EPUB
MS_OEDBX_Fmt	433	Microsoft Outlook Express DBX	DBX
BB_Activ_Fmt	434	BlackBerry Activation File	DAT ¹
DiskImage_Fmt	435	Disk Image	
Milestone_Fmt	436	Milestone Document	MLS, ML3, ML4, ML5, ML6, ML7, ML8, ML9
E_Transcript_Fmt	437	RealLegal E-Transcript File	PTX
PostScript_Font_Fmt	438	PostScript Type 1 Font	PFB
Ghost_DiskImage_Fmt	439	Ghost Disk Image File	GHO, GHS
JPEG_2000_JP2_File_Fmt	440	JPEG-2000 JP2 File Format Syntax (ISO/IEC 15444-1)	JP2, JPF, J2K, JPWL, JPX, PGX
Unicode_HTML_Fmt	441	Unicode HTML	HTM ¹ , HTML ¹
CHM_Fmt	442	Microsoft Compiled HTML Help	CHM
EMCMF_Fmt	443	Documentum EMCMF format	EMCMF
MS_Access_2007_Tmpl_Fmt	444	Microsoft Access 2007 Template	ACCDT
Jungum_Fmt	445	Samsung Electronics Jungum Global document	GUL
JBIG2_Fmt	446	JBIG2 File Format	JB2, JBIG2
EFax_Fmt	447	eFax file	EFX
AD1_Fmt	448	AD1 Evidence file	AD1
SketchUp_Fmt	449	Google SketchUp	SKP

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
GWFS_Email_Fmt	450	Group Wise File Surf email	GWFS
JNT_Fmt	451	Windows Journal format	JNT
Yahoo_yChat_Fmt	452	Yahoo! Messenger chat log	YCHAT
PaperPort_MAX_File_Fmt	453	PaperPort image file	MAX
ARJ_Fmt	454	ARJ (Archive by Robert Jung) file format	ARJ
RPMSG_Fmt	455	Microsoft Outlook Restricted Permission Message	RPMSG
MAT_Fmt	456	MATLAB file format	MAT, FIG
SGY_Fmt	457	SEG-Y Seismic Data format	SGY, SEGY
CDXA_MPEG_PS_Fmt	458	MPEG-PS container with CDXA stream	MPG ¹
EVT_Fmt	459	Microsoft Windows NT Event Log	EVT
EVTX_Fmt	460	Microsoft Windows Vista Event Log	EVTX
MS_OutlookOLM_Fmt	461	Microsoft Outlook for Macintosh format	OLM
WARC_Fmt	462	Web ARChive	WARC
JAVACLASS_Fmt	463	Java Class format	CLASS
VCF_Fmt	464	Microsoft Outlook vCard file format	VCF
EDB_Fmt	465	Microsoft Exchange Server Database file format	EDB
ICS_Fmt	466	Microsoft Outlook iCalendar file format	ICS, VCS
MS_Visio_2013_Fmt	467	Microsoft Visio 2013	VSDX, VSTX, VSSX
MS_Visio_2013_Macro_Fmt	468	Microsoft Visio 2013 macro	VSDM, VSTM, VSSM
ICHITARO_Compr_Fmt	469	ICHITARO Compressed format	JTDC
IWWP13_Fmt	470	Apple iWork 2013 Pages format	IWA

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
IWSS13_Fmt	471	Apple iWork 2013 Numbers format	IWA
IWPG13_Fmt	472	Apple iWork 2013 Keynote format	IWA
XZ_Fmt	473	XZ archive format	XZ
Sony_WAVE64_Fmt	474	Sony Wave64 format	W64
Conifer_WAVPACK_Fmt	475	Conifer Wavpack format	WV
Xiph_OGG_VORBIS_Fmt	476	Xiph Ogg Vorbis format	OGG
MS_Visio_2013_Stencil_Fmt	477	MS Visio 2013 stencil format	VSSX
MS_Visio_2013_Stencil_Macro_Fmt	478	MS Visio 2013 stencil Macro format	VSSM
MS_Visio_2013_Template_Fmt	479	MS Visio 2013 template format	VSTX
MS_Visio_2013_Template_Macro_Fmt	480	MS Visio 2013 template Macro format	VSTM
Borland_Reflex_2_Fmt	481	Borland Reflex 2 format	R2D
PKCS_12_Fmt	482	PKCS #12 (p12) format	P12, PFX
B1_Fmt	483	B1 format	B1
ISO_IEC_MPEG_4_Fmt	484	ISO/IEC MPEG-4 format	MP4
RAR5_Fmt	485	RAR5 Format	RAR5
Unigraphics_NX_Fmt	486	Unigraphics (UG) NX CAD Format	PRT
PTC_Creo_Fmt	487	PTC Creo CAD Format	ASM, PRT
KML_Fmt	488	Keyhole Markup Language	KML
KMZ_Fmt	489	Zipped Keyhole Markup Language	KMZ
WML_Fmt	490	Wireless Markup Language	WML

KeyView file formats and extensions, continued

Format Name	Format Number	Format Description	Associated File Extension
SO_Text_Fmt	492	Star Office Writer Text	SDW, SGL, VOR
SO_Spreadsheet_Fmt	493	Star Office Calc Spreadsheet	SDC
SO_Presentation_Fmt	494	Star Office Impress Presentation	SDD, SDA
SO_Math_Fmt	495	Star Office Math	SMF

1

This file extension can return more than one format number.

2

MHT, EML, and MBX files might return either format 2, 233, or 395, depending on the text in the file. In general, files that contain fields such as To, From, Date, or Subject are considered to be email messages; files that contain fields such as content-type and mime-version are considered to be MHT files; and files that do not contain any of those fields are considered to be text files.

3

All CAT file extensions, for example CATDrawing, CATProduct, CATPart, and so on.

Appendix D: Extract and Format Lotus Notes Subfiles

This section describes how to create XML templates to alter the appearance of extracted Lotus mail note subfiles so that they maintain the look and feel of the original notes.

- [Overview](#) 163
- [Customize XML Templates](#) 163
- [Template Elements and Attributes](#) 165
- [Date and Time Formats](#) 169

Overview

KeyView uses the NSF reader, nsfsr, to extract Lotus database files, and places Lotus mail notes in subfiles. The NSF reader uses a set of default XML templates to extract the notes and apply formatting, thereby approximating the look and feel of the original notes.

In some cases, you might need to customize the XML templates, for instance if your notes contain custom data. In such cases, you can modify the existing XML templates or create your own.

During extraction, the NSF reader loads all XML files in the `NSFtemplates` directory and its subdirectories (except for the `NSFtemplates\images` directory, which is reserved for images). During initialization, the KeyView XML parser verifies the XML templates. If the templates contain any invalid XML, elements, or attributes, initialization fails and errors are recorded in the `nsfsr.log` file.

Customize XML Templates

XML templates are enabled by default. In most cases, the default templates should be sufficient; however, you can customize them or create your own as required.

To customize XML templates for Lotus note extraction

1. Modify the template files in the following directory.
`install\OS\bin\NSFtemplates`
The `main.xml` file must exist in the `NSFtemplates` directory. It is the top-level template file that extracts all subfiles, usually by calling other templates.
2. Make sure that any modifications or additional XML files conform to the supported elements and attributes described in [Template Elements and Attributes, on page 165](#).
3. Extract the Lotus database file.

Use Demo Templates

For testing purposes, you can extract notes by using a set of demo templates, which are provided to demonstrate the proper usage of all the XML elements and attributes, because the default templates do not use all the XML elements.

The demo templates are available at:

install\OS\bin\NSFtemplates

To use the demo XML templates

1. In the `formats.ini` file, set the following parameter.

```
[nsfsr]  
UseDemoTemplate=1
```

2. In the `main.xml` file, uncomment the following section.

```
<ifini name="UseDemoTemplate" text="1">  
  <call file="demo.xml"/>  
  <quit/>  
</ifini>
```

Use Old Templates

For testing purposes, you can extract notes by using legacy templates, which produce MHTML output. You can generate similar output by disabling the XML templates, but using the old templates enables you to see the XML code and compare it to the standard and demo templates.

To use the old XML templates

1. In the `formats.ini` file, set the following parameter.

```
[nsfsr]  
UseOldTemplate=1
```

2. In the `main.xml` file, uncomment the following section.

```
<ifini name="UseOldTemplate" text="1">  
  <call file="default_old.xml">  
  <quit>  
</ifini>
```

Disable XML Templates

For testing purposes, you can disable XML templates; KeyView extracts the notes in MHTML format. You can compare the MHTML output directly by the NSF reader with the MHTML output indirectly by the NSF reader through the XML templates.

To disable XML templates

- 1. In the formats.ini file, set the following parameter.

```
[nsfsr]  
ExtractByTemplate=0
```

Template Elements and Attributes

This section lists the valid XML elements and attributes that you can use when creating or modifying templates. See the demo templates for examples.

Conditional Elements

The following table lists the valid conditional elements.

Conditional elements

Element	Description
<keyview>	The KeyView XML template container ("root") element
<if*>	<p>If the condition from the comparison is true, process the XML. Conditions can be nested up to 25 levels deep.</p> <p>Attributes</p> <ul style="list-style-type: none">• name. (Required) The name of the main item to compare to item or text.• item. (Required if no text) The name of the item to compare to the item specified by name.• text. (Required if no item) The text to compare to the item specified by name.
<ifex>, <ifnx>	<p>If name item exists and has a text value or not.</p> <p>The Notes item might have a value that cannot be converted to text, such as an image.</p>
<ifeq>, <ifne>, <iflt>, <ifle>, <ifgt>, <ifge>	<p>Respectively, if text ==, !=, <, >, <=, >, >=.</p> <p>Text comparison uses a case-insensitive string compare.</p>
<iftdeq>, <iftdne>, <iftdlt>, <iftdle>, <iftdgt>, <iftdge>	<p>Respectively, if time/date ==, !=, <, >, <=, >, >=.</p> <p>Time/date comparison converts dates to text in local time using the Notes default, TZFMT_NEVER, because Notes also sometimes converts fields to text internally. For example:</p> <p>text="06/30/2005 02:52:04 PM"</p>
<iftzeq>, <iftzne>	<p>Respectively, if the time zone equals or does not equal the comparison text, for example CDT, EST, and so on.</p>

Conditional elements, continued

Element	Description
<ifini>	If the value of the INI option specified in name equals the text value.
<else>	If the condition from the last <if> or <switch> was false, process XML.
<switch>	<p>If a name value exists, process XML.</p> <p>Attributes</p> <ul style="list-style-type: none"> name. (Required) The name of the main item to compare in <case> subelements.
<case>	<p>If the comparison condition is true, process XML, then stop processing the rest of <switch>.</p> <p>Attributes</p> <ul style="list-style-type: none"> text. (Required) The text to compare to the name item of <switch>.
<default>	If all <case> conditions were false, process XML. This element must be the last element in <switch>, after all the <case> elements. Any <case> elements after the <default> element are ignored.
<for>	<p>If a name value exists, process XML. Process for each part of the name item.</p> <p>Attributes</p> <ul style="list-style-type: none"> name. (Required) The name of the main item. max. (Optional) The maximum index to process. By default, all are processed.
<index>	Output <for> loop index (1-based). <index> is only valid within a <for> element.

Control Elements

The following table lists the valid control elements.

Control Elements

Element	Description
<call>	<p>Call another XML template. You can nest templates up to 10 levels deep.</p> <p>Attributes</p> <ul style="list-style-type: none"> file. (Required) The template file name. This name must be unique.
<log>	<p>Log message to the NSF log file.</p> <p>Attributes</p> <ul style="list-style-type: none"> text. (Required) The text to log.

Control Elements, continued

Element	Description
	<ul style="list-style-type: none">• type. (Optional) The type of log message. The following values are valid:<ul style="list-style-type: none">◦ ERROR◦ WARN◦ INFO◦ DIAG (the default option)◦ DEBUG◦ DUMP
<quit>	<p>Stop processing the template. Exits without error.</p> <p>Attributes</p> <ul style="list-style-type: none">• text. (Optional) The text to log.• type. (Optional) The type of log message. See <log>, on the previous page.
<stop>	<p>Stop processing the template. Exits with an ERROR log message.</p> <p>Attributes</p> <ul style="list-style-type: none">• text. (Required) The text to log.

Data Elements

The following table lists the valid data elements.

Data elements

Element	Description
<text>	<p>Output text.</p> <p>Attributes</p> <ul style="list-style-type: none">• name. (Required if there is no parent) The name of the item to output.
<rich>	<p>Output rich text (MHTML). Images are output in the next part or parts of the MHTML, after the first <HTML> part.</p> <p>Attributes</p> <ul style="list-style-type: none">• name. (Required if there is no parent) The name of the item to output.
<body>	<p>Output the message body in rich text (MHTML). As with <rich>, above, images are output in the next part or parts of the MHTML.</p>
<form>	<p>Output the message form (usually \$Body field) in rich text (MHTML).</p> <p>Attributes</p> <ul style="list-style-type: none">• name. (Required if there is no parent) The name of the item to output.

Data elements, continued

Element	Description
<addr>	<p>Output an address.</p> <p>Attributes</p> <ul style="list-style-type: none"> • <code>name</code>. (Required if there is no parent) The name of the item to output. • <code>type</code>. (Optional) The type of address to output. Set this attribute to <code>CN</code> (Common Name), which is the only supported type.
<name>	<p>Output the name of the last <code>name</code> item, or in other words the current main item. The item must exist.</p>
<format>	<p>Set the default format for <date> and <date_kv>. This element does not set the <text> format. See Date and Time Formats, on the next page for a list of all Notes and KeyView date and time formats and integer values.</p> <p>Attributes</p> <ul style="list-style-type: none"> • <code>format</code>. (Optional. Omit to reset to defaults) The Notes and KeyView date and time format. You can set the following formats: <ul style="list-style-type: none"> ◦ <code>TD=int</code>. The Time Date format (<code>TDFMT_*</code>) ◦ <code>TS=int</code>. The Time Show format (<code>TSFMT_*</code>) ◦ <code>TT=int</code>. The Time Time format (<code>TTFMT_*</code>) ◦ <code>TZ=int</code>. The Time Zone format (<code>TZFMT_*</code>) ◦ <code>KV=int</code>. The KeyView date and time format <p>where <code>int</code> is an integer value that corresponds to the desired format.</p> <p>Separate multiple formats with commas. For example:</p> <p><code>format="TD=0,TS=2,TT=1,TZ=1,KV=55"</code></p>
<date>	<p>Output a Notes date.</p> <p>Attributes</p> <ul style="list-style-type: none"> • <code>name</code>. (Required if there is no parent) The name of the item to output. • <code>format</code>. (Optional) See <format>, above. You can set the following values: <ul style="list-style-type: none"> ◦ <code>TD</code> ◦ <code>TS</code> ◦ <code>TT</code> ◦ <code>TZ</code>
<date_kv>	<p>Output a KeyView date.</p> <p>Attributes</p> <ul style="list-style-type: none"> • <code>name</code>. (Required if there is no parent) The name of the item to output. • <code>format</code>. (Optional) See <format>, above. You can set the following values: <ul style="list-style-type: none"> ◦ <code>TZ</code>

Data elements, continued

Element	Description
	<ul style="list-style-type: none"> ◦ KV
<time>	<p>Output a time range, for example 1 hour, 30 minutes.</p> <p>Attributes</p> <ul style="list-style-type: none"> • name. (Required if there is no parent) The item name of the start date or time. • item. (Required) The item name of the end date or time.
<zone>	<p>Output a Notes time zone mnemonic, for example MST.</p> <p>Attributes</p> <ul style="list-style-type: none"> • name. (Required if there is no parent) The name of date item to output.
<zone_UTC>	<p>Output a time zone as UTC, for example (UTC-06:00).</p>
<logo>	<p>Output the mail header logo.</p> <p>The image link is included in the output; the actual image is output to a different part of the MHTML subfile.</p>
<image>	<p>Output an image.</p> <p>The image link is included in the output; the actual image is output to the MHTML next part, as with <rich>, on page 167 and <body>, on page 167.</p>
<image_uri>	<p>Output an image URI, in quotation marks. The actual image is output to a different part of the MHTML subfile.</p> <p>Attributes</p> <ul style="list-style-type: none"> • link. (Required if there is no file) The image link, such as a form or title name. For example: • link="StdNotesLtr0" • file. (Required if there is no link) The name of the image file. The file must exist in the ../../templates/images directory. For example: • file="boxcheck.gif"

Date and Time Formats

This section lists the supported Notes and KeyView date and time formats for use with <format>, <date>, and <date_kv>.

Lotus Notes Date and Time Formats

This section lists supported Lotus Notes date and time formats, and the integer values that specify each one.

Lotus Notes date and time formats

Format	Integer Value	Description
TDFMT_FULL	0	(The Notes default) Year, month, and day
TDFMT_CPARTIAL	1	Month and day, year if not this year
TDFMT_PARTIAL	2	Month and day
TDFMT_DPARTIAL	3	Year and month
TDFMT_FULL4	4	Four-digit year, month, and day
TDFMT_CPARTIAL4	5	Month and day, four-digit year if not this year
TDFMT_DPARTIAL4	6	Four-digit year and month
TTFMT_FULL	0	(Notes default) Hour, minute, and second
TTFMT_PARTIAL	1	Hour and minute
TTFMT_HOUR	2	Hour
TZGMT_NEVER	0	(Notes default) All time zones are converted to the current time zone
TZGMT_SOMETIMES	1	Show only when outside the current time zone
TZGMT_ALWAYS	2	Show for all time zones
TSFMT_DATE	0	Date
TSFMT_TIME	1	Time
TSFMT_DATETIME	2	(The Notes default) Date and time
TSFMT_CDATETIME	4	Date and time, or time today or time yesterday

KeyView Date and Time Formats

This section lists KeyView date and time formats. The KeyView formats use the following syntax:

Month Month = full month name
 Mon = abbreviated month name
 m = month (number)
 mm = two-digit month (leading 0)

Weekday	Weekday = full weekday name
	Wday = abbreviated weekday name
Year	yy = two-digit year
	yyyy = four-digit year
Day	d = day (number)
	dd = two-digit day (leading 0)
Time	h = 12-hour
	H = 24-hour
	m = minutes
	s = seconds
	P = AM/PM
	p = am/pm
Separators	_ = space
	c = comma
	s = slash
	a = dash
	o = dot

KeyView date and time formats

Format	Output	Integer Value
12-Hour and 24-Hour Time Formats		
KVDTF_P	P	1
KVDTF_P_hmm	P h:mm	2
KVDTF_hmm_P	h:mm P	3
KVDTF_P_hhmm	P hh:mm	4
KVDTF_hhmm_P	hh:mm P	5
KVDTF_P_hmmss	P h:mm:ss	6
KVDTF_hmmss_P	h:mm:ss P	7
KVDTF_P_hhmmss	P hh:mm:ss	8
KVDTF_hhmmss_P	hh:mm:ss P	9
KVDTF_Hmm	H:mm	10
KVDTF_HHmm	HH:mm	11

KeyView date and time formats, continued

Format	Output	Integer Value
KVDTF_mmss	mm:ss	12
KVDTF_Hmmss	H:mm:ss	13
KVDTF_HH:mmss	HH:mm:ss	14
Numerical Date Formats with Slashes		
KVDTF_mmsdd	mm/dd	15
KVDTF_msdsyy	m/d/yy	16
KVDTF_mmsddsy	mm/dd/yy	17
KVDTF_mmsddsy	mm/dd/yyyy	18
KVDTF_ddsmm	dd/mm	19
KVDTF_ddsmmsy	dd/mm/yy	20
KVDTF_ddsmmsy_Hmm	dd/mm/yy H:mm	21
KVDTF_ddsmm_P_hmm	dd/mm P h:mm	22
KVDTF_ddsmm_hmm_P	dd/mm h:mm P	23
KVDTF_ddsmm_P_hhmm	dd/mm P hh:mm	24
KVDTF_ddsmm_hhmm_P	dd/mm hh:mm P	25
KVDTF_ddsmmsy_P_hmm	dd/mm/yy P h:mm	26
KVDTF_ddsmmsy_hmm_P	dd/mm/yy h:mm P	27
KVDTF_ddsmmsy_P_hmmss	dd/mm/yy P h:mm:ss	28
KVDTF_ddsmmsy_hmmss_P	dd/mm/yy h:mm:ss P	29
KVDTF_ddsmmsy_P_hhmmss	dd/mm/yy P hh:mm:ss	30
KVDTF_ddsmmsy_hhmmss_P	dd/mm/yy hh:mm:ss P	31
KVDTF_yysmmsdd_P_hhmmss	yy/mm/dd P hh:mm:ss	32
KVDTF_yysmmsdd_hhmmss_P	yy/mm/dd hh:mm:ss P	33
KVDTF_msdsy_Hmm	m/d/yy H:mm	34
KVDTF_mmsddsy_Hmm	mm/dd/yy H:mm	35
KVDTF_msdsy_P_hmm	m/d/yy P h:mm	36
KVDTF_msdsy_hmm_P	m/d/yy h:mm P	37

KeyView date and time formats, continued

Format	Output	Integer Value
KVDTF_mmsddsy_hmm_P	mm/dd/yy h:mm P	38
KVDTF_mmsdd_P_hhmm	mm/dd P hh:mm	39
KVDTF_mmsdd_hhmm_P	mm/dd hh:mm P	40
KVDTF_mmsddsy_P_hhmmss	mm/dd/yy P hh:mm:ss	41
KVDTF_mmsddsy_hhmmss_P	mm/dd/yy hh:mm:ss P	42
KVDTF_ms_d	m/d	43
KVDTF_yysm	yy/m	44
KVDTF_yysmm	yy/mm	45
KVDTF_yysmsd	yy/m/d	46
KVDTF_yysmmsdd	yy/mm/dd	47
KVDTF_yyyysmmsdd	yyyy/mm/dd	48
Numerical Date Formats with Dashes		
KVDTF_ddammayy	dd-mm-yy	49
KVDTF_mmadd	mm-dd	50
KVDTF_mmayy	mm-yy	51
KVDTF_yyammadd	yy-mm-dd	52
KVDTF_yyyammadd	yyyy-mm-dd	53
KVDTF_yyyammaddaHHmmss	yyyy-mm-dd-HH:mm:ss	54
Numerical Date Formats with Dots		
KVDTF_yyomod	yy.m.d	55
KVDTF_yyommodd	yy.mm.dd	56
KVDTF_mod	m.d	57
KVDTF_mmodd	mm.dd	58
Numerical and String Date Formats with Dashes, Commas, and Spaces		
KVDTF_ddaMon	dd-Mon	59
KVDTF_daMonayy	d-Mon-yy	60
KVDTF_ddaMonayy	dd-Mon-yy	61

KeyView date and time formats, continued

Format	Output	Integer Value
KVDTF_ddaMonayyyy	dd-Mon-yyyy	62
KVDTF_Mon	Mon	63
KVDTF_Monayy	Mon-yy	64
KVDTF_Monayyyy	Mon-yyyy	65
KVDTF_Monaddayy	Mon-dd-yy	66
KVDTF_yyammadd_P_hhmmss	yy-mm-dd P hh:mm:ss	67
KVDTF_mmadd_P_hhmm	mm-dd P hh:mm	68
KVDTF_Mon_yy	Mon yy	69
KVDTF_Monc_yy	Mon, yy	70
KVDTF_Month	Month	71
KVDTF_Monthayy	Month-yy	72
KVDTF_Month_yy	Month yy	73
KVDTF_Monthc_yy	Month, yy	74
KVDTF_Monthayyyy	Month-yyyy	75
KVDTF_Month_yyyy	Month yyyy	76
KVDTF_Monthc_yyyy	Month, yyyy	77
KVDTF_Mon_dc_yyyy	Mon d, yyyy	78
KVDTF_d_Monc_yyyy	d Mon, yyyy	79
KVDTF_yyyy_Mon_d	yyyy Mon d	80
KVDTF_Month_dc_yyyy	Month d, yyyy	81
KVDTF_d_Monthc_yyyy	d Month, yyyy	82
KVDTF_yyyy_Month_d	yyyy Month d	83
Weekday Date Formats		
KVDTF_wday	wday	84
KVDTF_Weekday	Weekday	85
KVDTF_wdayc_Mon_dc_yyyy	wday, Mon d, yyyy	86
KVDTF_Weekdayc_Month_dc_yyyy	Weekday, Month d, yyyy	87

KeyView date and time formats, continued

Format	Output	Integer Value
KVDTF_Weekdayc_d_Monthc_yyyy	Weekday, d Month, yyyy	88

Appendix E: File Format Detection

This section describes how file formats are detected in Filter SDK.

• Introduction	176
• Extract Format Information	176
• Determine Format Support	176
• Translate Format Information	179
• Determine a Document Reader	180
• Category Values in <code>formats.ini</code>	180

Introduction

The KeyView format detection module (`kwad`) detects a file's format, and reports the information to the API, which in turn reports the information to the developer's application. If the detected format is supported by the KeyView SDK, the detection module also loads the appropriate structured access layer and document reader for further processing. For a list of supported formats, see [Supported Formats, on page 94](#).

Extract Format Information

You can extract format information from a document by using one of the `getDocFormatInfo` methods. These methods extract the major format, file class, version, and document attributes, and populate the `DocFormatInfo` class. They return the format information as a string. The format information that you can extract is listed in the header file `adinfo.h`.

For information on how to translate the extracted format information, see [Translate Format Information, on page 179](#).

Determine Format Support

After the file format is extracted, the detection module uses the `formats.ini` file to determine whether the format is supported by KeyView, and the appropriate structured access layer and reader to load.

The `formats.ini` file is in the directory `install\OS\bin`, where `install` is the path name of the Filter installation directory and `OS` is the name of the operating system. It contains the following information:

- Coded format information. To translate this information, see [Translate Format Information, on page 179](#).
- The reader associated with each format. See [Determine a Document Reader, on page 180](#).
- Configuration parameters.
- Locale settings for internal use.

Example formats.ini file entries

```
123=mw
152=xyw
178=wp6
189=mw6
2=af
200=pdf
205=mb
210=htm
251=htm
```

NOTE:

The `formats.ini` file applies to all formats except graphics. Detection of graphics formats is handled by an internal module named KeyView Picture Interchange Format (KPIF).

Refine Detection of Text Files

During text detection, KeyView analyses the first 1 kB and last 1 kB of data in a document. If less than 10% of that data consists of non-ASCII characters, KeyView detects the document as a text file.

However, depending on the type of documents you are working with, the default settings might not provide the desired level of accuracy. Configuration flags enable you to change the amount of data to read at the end of a file, the percentage of non-ASCII characters permitted in a text file, and whether to use or ignore the file extension to determine the document format.

Change the Amount of File Data to Read

During file detection, KeyView reads characters from the beginning and end of a file—by default, it reads the first and last 1,024 bytes of data. Large text files might contain many irrelevant characters at the end of a file, so KeyView might not accurately detect the file format. You can set a configuration flag to increase the amount of data to read from the end of a file during detection.

To change the amount of data to read during detection

- In the `formats.ini` file, set the following flag in the `detection_flags` section:

```
[detection_flags]
non_ascii_chars_end_block_size=kB
```

where *kB* is the number of kilobytes to read from the end of the file, from 0 to 10. The default value is 1.

NOTE:

The file size must be greater than the value specified in the flag. If the flag value is greater than the file size, KeyView does not use the flag.

Change the Percentage of Allowed Non-ASCII Characters

By default, if less than 10% of the analyzed data in a document consists of non-ASCII characters, it is detected as a text file. Depending on the type of files that you are working with, changing the default percentage might increase detection accuracy.

To change the percentage of non-ASCII characters allowed in text files

- In the `formats.ini` file, set the following flag in the `detection_flags` section:

```
[detection_flags]
non_ascii_chars_in_text=N
```

where *N* is the percentage of non-ASCII characters to allow in text files. Files that contain a lower percentage of non-ASCII characters than *N* are detected as text files. The default value is 10.

Allow Consecutive NULL Bytes in a Text File

By default, if a document contains consecutive NULL bytes, it is not detected as text. Depending on the type of files that you are working with, changing the default might increase detection accuracy.

To allow consecutive NULL bytes of ASCII characters in text files

In the `formats.ini` file, set the following flag in the `detection_flags` section:

```
[detection_flags]
ascii_allow_null_bytes=1
```

The default value is 0 (do not allow consecutive NULL bytes).

Use the File Extension for Detection

Sometimes KeyView detects certain file formats, such as CSV, as ASCII because of the content of the documents. In such cases, you can configure KeyView to use the file extension to determine the document format. Using the file extension can improve detection of formats such as CSV, but might not detect text files successfully if they have incorrect file extensions.

To use the file extension for ASCII files during detection

- In the `formats.ini` file, set the following flag in the `detection_flags` section:

```
[detection_flags]
use_extension_for_ascii=1
```

The default is 0 (do not use the file extension).

Translate Format Information

Format information can include file attributes in the following categories:

- Major format
- File class
- Minor format
- Major version
- Minor version

Not all categories are required. Many formats only include major format and file class, or major format only.

The format information has the following structure:

```
MajorFormat.FileClass.MinorFormat.MajorVersion.MinorVersion
```

For example:

```
81.2.0.9.0
```

Each number in the format information represents a file attribute. The entry `81.2.0.9.0` represents a Lotus 1-2-3 Spreadsheet file version 9.0, where

81= Lotus 1-2-3 Spreadsheet (major format)

2 = Spreadsheet (file class)

0 = not defined (minor format)

9 = 9 (major version)

0 = 0 (minor version)

This example applies to the `formats.ini` file. When extracting format information using the `getDocFormatInfo` methods, the same format is represented as `294.2.9.0`.

NOTE:

The format values returned from `getDocFormatInfo` differ from those in `formats.ini` because the former defines a unique ID for each major format, while the latter uses a major version, minor version, and minor format to distinguish between formats.

Distinguish Between Formats

The `DocFormatInfo` class provides a unique ID for each major format. For example, a call to `getDocFormatInfo` would return `351.1.0` for a Microsoft Word XML format. The major format 351 is unique to this format.

Unlike `DocFormatInfo`, the `formats.ini` file distinguishes between formats by using the major version number. For example, in the `formats.ini` file, a Microsoft Word 2003 XML format is defined as `285.1.0.100.0`. The major format 285 and file class 1 are the same values for generic XML. The major version 100 distinguishes the format as Microsoft Word 2003 XML.

The major version is used to specify the following formats:

- Microsoft Office 2003 XML. This format has the same major format and file class as generic XML (285.1). It is distinguished from generic XML by using the following major versions:
 - Word: 100
 - Excel: 101
 - Visio: 110
- The XHTML format has the same major format and file class as HTML (210.1). It is distinguished from HTML by using the major version 100.

Determine a Document Reader

The format detection module uses the `formats.ini` file to determine whether a format is supported, and to determine the reader to use to parse a format. The entries in the `formats.ini` file list each format's coded value, and an abbreviation for the format's reader.

The reader abbreviation is a truncated version of the reader's library name. Adding "sr" to the end of an abbreviation creates the name of the reader. For example, this example entry specifies that a Lotus 1-2-3 Spreadsheet file version 9.0 is parsed by the Lotus 1-2-3 filter, 1123sr:

```
81.2.0.9.0=1123
```

[List of Required Files for Redistribution, on page 199](#) lists the readers provided with KeyView.

Category Values in formats.ini

This section lists the possible category values for format information in the `formats.ini` file. The corresponding values for format information extracted by a call to `getDocFormatInfo` are listed in the header file `adinfo.h`.

- [Major Formats](#)
- [File Classes](#)
- [Minor Formats](#)

Major Formats

Number	Format	File Class
1	AES Multiplus Comm Format	Word processor
2	ASCII File word processor/MS DOS Batch File format	Word processor
3	Applix Asterix	Word processor
4	Microsoft Windows Bitmap image (BMP)	Raster image
5	Convergent Tech DEF Comm. format	Word processor
6	Corel Draw (CDR)	Vector graphic

Major Formats, continued

Number	Format	File Class
7	Keyword COM.FILE (KSIF)	
8	Computer Graphics Metafile (CGM)	Vector graphic
9	Word Connection	Word processor
10	COMET TOP Word	Word processor
11	DG CEOWrite	Word processor
12	Honey Bull DSA101	Word processor
13	IBM DCA-RFT	Word processor
14	DDIF	Word processor
15	Dummy File (Internal)	
16	DG Common Data Stream (CDS)	Word processor
17	Dummy Print File (Internal)	
18	Windows Micrografx Draw (DRW)	Vector graphic
19	Data Point VISTAWORD	Word processor
20	DECdx	Word processor
21	Enable	Word processor
22	Encapsulated PostScript (EPS)	Raster image
23	DOS/Windows Executable (EXE, DLL)	Executable
24	CCITT Group 3 1-Dimensional (G31D)	Raster image
25	Graphics Interchange format (GIF)	Raster image
26	Hewlett Packard	Word processor
27	IBM 1403 Line Printer	Word processor
28	IBM DCF Script	Word processor
29	IBM DCA-FFT	Word processor
30	Interleaf	Word processor
31	GEM Bit Image	Raster image
32	IBM Display Write 4	Word processor
33	Raster Graphics	Raster image

Major Formats, continued

Number	Format	File Class
34	Keywords PICL	
35	Lotus AMI Pro	Word processor
36	MORE Database Outliner (Mac)	Outline/planning
37	Lyrix	Word processor
38	MASS-11	Word processor
39	MacPaint	Raster image
40	Microsoft Word Mac	Word processor
41	Informix SmartWare II Communication File	Communications
42	Microsoft Word for Windows	Word processor
43	MultiMate 4.0	Word processor
44	Multiplan Spreadsheet	Spreadsheet
45	Microsoft Rich Text Format (RTF)	Word processor
46	Microsoft Word 5.0 (PC)	Word processor
47	NBI Async Archive Format	Word processor
48	Navy DIF	Word processor
49	NBI Net Archive Format	Word processor
50	NIOS TOP	Word processor
51	FileMaker (Mac)	Database
52	ODA/ODIF	Word processor
53	OLIDIF	Word processor
54	Keyword OSM	
55	Office Writer	Word processor
56	PC Paint Brush Graphics (PCX)	Raster image
57	CPT Communication Format	Word processor
58	Lotus PIC	Vector graphic
59	Macintosh Quick Draw Picture Format (PICT)	Raster image
60	Philips Script	Word processor

Major Formats, continued

Number	Format	File Class
61	PostScript File	Vector graphic
62	PRIMEWORD	Word processor
63	Quadratron Q-One (V1.93J)	Word processor
64	Quadratron Q-One (V2.0)	Word processor
65	SAMNA Word IV	Word processor
66	Lotus AMI Pro Draw (SDW)	Raster image
67	SYLK Spreadsheet	Spreadsheet
68	Informix SmartWare II	Word processor
69	Symphony Spreadsheet	Spreadsheet
70	Truevision Targa	Raster image
71	Tagged Image File (TIFF)	Raster image
72	Targon Word (V 2.0)	Word processor
73	Uniplex Ucalc Spreadsheet	Spreadsheet
74	Uniplex (V6.01)	Word processor
75	Microsoft Word (UNIX)	Word processor
76	WANG PC	Word processor
77	WordERA (V 1.0)	Word processor
78	WANG WPS Comm. format	Word processor
79	WordPerfect Mac	Word processor
80	WordPerfect 5.2	Word processor
81	Lotus 1-2-3 Spreadsheet	Spreadsheet
82	WordMARC word processor	Word processor
83	Microsoft Windows Metafile (WMF) Graphics	Raster image
84	Informix SmartWare II Database	Database
85	WordPerfect Graphics V1.0 (WPG)	Raster image
86	WordPerfect	Word processor
87	WordStar	Word processor

Major Formats, continued

Number	Format	File Class
88	Wang WITA	Word processor
89	Xerox 860 Comm. format	Word processor
90	Microsoft Excel Spreadsheet	Spreadsheet
91	Xerox Writer word processor	Word processor
92	DIF Spreadsheet	Spreadsheet
93	ENABLE Spreadsheet	Spreadsheet
94	Supercalc Spreadsheet	Spreadsheet
95	Ultracalc Spreadsheet	Spreadsheet
96	Informix SmartWare Spreadsheet	Spreadsheet
97	Serialized Object Format (SOF) Encapsulation format	Encapsulation
98	Microsoft PowerPoint (PC)	Presentation
99	Microsoft PowerPoint (Mac)	Presentation
100	Aldus PageMaker (Mac)	Desktop Publishing
101	Aldus PageMaker (DOS)	Desktop Publishing
103	Microsoft Works (Mac)	Word processor
104	Microsoft Works Database (Mac)	Database
105	Microsoft Works Spreadsheet (Mac)	Spreadsheet
106	Microsoft Works Communication (Mac)	Communication
107	Microsoft Works (PC)	Word processor
108	Microsoft Works Database (PC)	Database
109	Microsoft Works Spreadsheet (PC)	Spreadsheet
111	PC Library Module	Library module
112	MacWrite	Word processor
113	MacWrite II	Word processor
114	Aldus Freehand Mac	Vector graphic
115	Disk Doubler Compression format	Encapsulation
116	HP Graphics Language (HP-GL)	Vector graphic

Major Formats, continued

Number	Format	File Class
117	Adobe Maker Interchange Format (MIF)	Desktop Publishing
118	JPEG File Interchange Format (JFIF)	Raster image
119	Reflex Database	Database
120	Framework II	Mixed format
121	Paradox (PC) Database	Database
123	Microsoft Windows Write	Word processor
124	Quattro Pro Spreadsheet (DOS)	Spreadsheet
126	Persuasion Presentation	Presentation
127	Corel Presentation	Presentation
128	Microsoft Windows Icon Format (ICO) Graphics	Raster image
129	Microsoft Project	Time scheduling
131	Harvard Graphics	Desktop publishing
132	Zip Archive Format	Encapsulation
133	Microsoft Windows Cursor (CUR) Graphics	Raster image
134	Quark Express (Mac)	Desktop publishing
135	ARC/PAK Archive format	Encapsulation
136	Adobe FrameMaker	Desktop publishing
137	Microsoft Publisher	Desktop publishing
138	Plan Perfect	Time scheduling
139	WordPerfect General File Format	Miscellaneous
140	Lotus Freelance	Presentation
141	Microsoft Wave Sound File	Sound
142	MIDI Sound File	Sound
143	AutoCAD DXF Graphics	Vector graphic
144	dBase Database	Database
145	OS/2 PM Metafile Graphics	Vector graphic
146	Lasergraphics Language	Vector graphic

Major Formats, continued

Number	Format	File Class
147	AutoShade Rendering File Format	Vector graphic
148	Graphics Environment Manager (GEM VDI)	Vector graphic
149	Microsoft Windows Help File	Miscellaneous
150	Volkswriter	Word processor
151	Ability Office (SS, DB, GR, WP, COM)	
152	XyWrite/Nota Bene	Word processor
153	Comma Separated Values (CSV)	Spreadsheet
154	Writing Assistant word processor	Word processor
155	WordStar 2000	Word processor
156	WordStar 6.0	Word processor
157	HP Printer Control Language (PCL)	Vector graphic
158	(UNIX/VAX/SUN) Executable	Executable
159	(UNIX/VAX/SUN) Object Module	Object module
160	(UNIX/VAX/SUN) Link Library	Library module
161	NeXT SUN Audio Data	Sound
162	NeWS font file (SUN)	Font
163	cpio Archive Format (UNIX/VAX/SUN)	Encapsulation
164	PEX Binary Archive (SUN)	Encapsulation
165	SUN vfont definition	Font
166	Curses Screen Image (UNIX/VAX/SUN)	Raster image
167	UU Encoded Encryption File	Encapsulation
168	WriteNow	Word processor
169	PC Object Module	Object module
170	Microsoft Windows Group File	Miscellaneous
171	PC True Type Font	Font
172	Program Information File	Miscellaneous
173	PC COM executable file	Executable

Major Formats, continued

Number	Format	File Class
174	Adobe FrameMaker Markup Language	Desktop publishing
175	Stuff It Archive (Mac)	Encapsulation
176	PeachCalc Spreadsheet	Spreadsheet
177	Wang Office GDL Header Encapsulation	Encapsulation
178	WordPerfect 6.0	Word processor
179	Q & A for DOS	Word processor
180	Q & A for Windows	Word processor
181	DEC WPS PLUS	Word processor
182	DCX Fax format	Fax
183	Microsoft Windows OLE 2 Encapsulation	Encapsulation
184	Quattro Pro for Windows	Spreadsheet
185	Keyword Viewer Markup Format	
186	EBCDIC Text	Word processor
187	DCS	Word processor
188	Microsoft Excel Spreadsheet 95, 2000	Spreadsheet
189	Microsoft Word for Windows 95	Word processor
190	UNIX SHAR Encapsulation	Encapsulation
191	Lotus Notes Bitmap	Raster image
192	UNIX Compress Encapsulation	Encapsulation
193	Lotus Notes CDF	Word processor
194	UNIX TAR Encapsulation	Encapsulation
195	WordPerfect Graphics V2.0 (WPG2)	Raster image Vector graphic
196	ODA/ODIF (FOD 26)	Word processor
197	ALIS	Word processor
198	GZ Compress Encapsulation	Encapsulation
199	Envoy (EVY)	Word processor

Major Formats, continued

Number	Format	File Class
200	Adobe Portable Document Format (PDF)	Word processor
201	KW ODA Internal Raw Bitmap (RBM)	Raster image
202	KW ODA G4 (G4)	Raster image
203	KW ODA G31D (G31)	Raster image
204	KW ODA Internal G32D (G32)	Raster image
205	Microsoft Word for Mac V 4.x/5.x	Word processor
206	BinHex 4.0 encoded file	Encapsulation
207	SMTP document	Encapsulation
208	MIME format - Microsoft Outlook Express (EML)/Mailbox (MBX)	Encapsulation
209	SGML document	Word processor
210	HTML document XHTML ¹	Word processor
211	ACT Format	Word processor
212	Microsoft PowerPoint 95	Presentation
213	Portable Network Graphics (PNG)	Raster image
214	Video for Windows	Movie
215	Windows Animated Cursor	Raster image
216	Windows C++ Object Storage	Mixed format
217	Windows Palette	Raster image
218	RIFF Device Independent Bitmap	Raster image
219	RIFF MIDI	Sound
220	RIFF Multimedia Movie	Movie
221	MPEG Movie	Movie
222	QuickTime Movie	Movie
223	Audio Interchange File Format (AIFF) Sound	Sound
224	Amiga MOD Sound	Sound
225	Amiga IFF (8SVX) Sound	Sound

Major Formats, continued

Number	Format	File Class
226	Creative Voice (VOC) Sound	Sound
227	Microsoft Works (Windows)	Word processor
228	Microsoft Works Spreadsheet (Windows)	Spreadsheet
229	AutoDesk Animator FLIC Animation	Animation
230	AutoDesk Animator Pro FLIC Animation	Animation
231	Microsoft Works Database (Windows)	Database
232	Microsoft Works Communication (Windows)	Communications
233	Compactor / Compact Pro Archive	Encapsulation
234	VRML	Vector graphic
235	QuickDraw 3D Metafile (3DMF)	Vector graphic
236	PGP Secret Keyring	Encapsulation
237	PGP Public Keyring	Encapsulation
238	PGP Encrypted Data	Encapsulation
239	PGP Signed Data	Encapsulation
240	PGP Signed and Encrypted Data	Encapsulation
241	PGP Signature Certificate	Encapsulation
242	ASCII-armored PGP Public Keyring	Encapsulation
243	ASCII-armored PGP encoded	Encapsulation
244	ASCII-armored PGP signed	Encapsulation
245	OLE DIB object	Raster image
246	PGP Compressed Data	Encapsulation
247	SGI Image	Raster image
248	Lotus Screen Cam	Animation
249	MPEG Audio	Sound
250	FTP Session Data	Communications
251	Netscape Bookmark file	Word processor
252	Corel Draw CMX	Vector image

Major Formats, continued

Number	Format	File Class
253	AutoCAD Drawing (DWG)	Vector graphic
254	AutoDesk WHIP	Vector graphic
255	Macromedia Director	Animation
256	Real Audio	Sound
257	MS DOS Device Driver	Executable
258	Micrografx Designer	Vector graphic
259	Simple Vector format (SVF)	Vector graphic
260	WordPerfect Office document (WPD)	
261	Applix Words	Word processor
262	Applix Graphics	Presentation
263	Microsoft Access	Database
264	Usenet format	Word processor
265	MacBinary	Encapsulation
266	Apple Single	Encapsulation
267	Apple Double	Encapsulation
268	Lotus Word Pro	Word processor
269	Microsoft Word 97, 2000	Word processor
270	Enhanced Window Metafile	Vector graphic
271	Microsoft Office Drawing	Vector graphic
272	Microsoft PowerPoint 97, 2000	Presentation
273	Extended or Custom XML	Word processor
274	Device Independent file (DVI)	Vector graphic
275	Unicode	Word processor
276	Framework	Mixed
277	KPIF Chart Stream	
278	Applix Spreadsheet	Spreadsheet
279	Microsoft Device Independent Bitmap	Raster image

Major Formats, continued

Number	Format	File Class
280	KeyView GPF Filter	
281	Microsoft Project 98, 2000, 2002	Time scheduling
282	Folio Flat file	Word processor
283	HWP (Arae-Ah Hangul)	Word processor
284	JustSystems Ichitaro	Word processor
285	Generic XML format Microsoft Office 2003 XML format ²	Word processor
286	Fujitsu Oasys	Word processor
287	Portable Bitmap Utilities (PBM)	Raster image
288	Portable Greymap Utilities (PGM)	Raster image
289	Portable Pixmap Utilities (PPM)	Raster image
290	X Bitmap (XBM)	Raster image
291	X Pixmap (XPM)	Raster image
292	X Image	Raster image
293	PCD Image	Raster image
294	Microsoft Visio	Presentation
295	Microsoft Outlook (MSG)	Encapsulation
296	XHTML document	Word processor
297	Microsoft Outlook Personal Folders file (PST)	Encapsulation
298	WinRAR Compressed Archive format (RAR)	Encapsulation
299	Lotus Notes Database (NSF) Legato Extender ONM	Encapsulation
300	Macromedia Flash	Word processor
301	Microsoft Word 2007 (XML format)	Word processor
302	Microsoft Excel 2007 (XML format)	Spreadsheet
303	Microsoft PowerPoint 2007 (XML format)	Presentation
304	Open PGP (new format packets only)	Encapsulation

Major Formats, continued

Number	Format	File Class
305	Intergraph version 7 DGN	Vector graphic
306	Microstation version 8 DGN	Vector graphic
307	Microsoft Word 2007 Macro	Word processor
308	Microsoft Excel 2007 Macro	Spreadsheet
309	Microsoft PowerPoint Macro	Presentation
310	Microsoft Compression folder (LZH)	Encapsulation
311	Office 2007 Document	Miscellaneous
312	XML Paper Specification	Word processor
313	Lotus Domino Extensible Language	Encapsulation
314	OASIS Open Document (ODT)	Word processor
315	OASIS Open Document (ODS)	Spreadsheet
316	OASIS Open Document (ODP)	Presentation
317	Legato EMailXtender Native Message	Word Processor
319	Transfer Neutral Encapsulation Format (TNEF)	Encapsulation
320	CADAM Drawing	Vector graphic
321	CADAM Drawing Overlay	Vector graphic
322	NURSTOR Drawing	Vector graphic
323	HP Graphics Language (Plotter)	Vector graphic
324	Advanced Systems Format	Miscellaneous
325	Windows Media Audio Format	Sound
326	Windows Media Video Format	Movie
327	Legato EMailXtender Archive	Encapsulation
328	7-Zip	Encapsulation
329	Microsoft Office 2007 Excel Binary Format	Spreadsheet
330	Microsoft Cabinet File	Encapsulation
331	CATIA formats	Vector graphic
332	Yahoo! Instant Messenger	Word processor

Major Formats, continued

Number	Format	File Class
333	Founder Chinese E-paper Basic	Word processor
334	Corel Quattro Pro X4	Spreadsheet
335	MIME HTML	Word processor
336	Microsoft Document Imaging Format	Raster image
337	Microsoft Office Groove File Format	Word processor
338	Apple iWorks Pages	Word processor
339	Apple iWorks Numbers	Spreadsheet
340	Apple iWorks Keynote	Presentation
341	Microsoft Backup File	Encapsulation
342	Microsoft Access 2007	Database
343	Microsoft Entourage Database	Encapsulation
344	Mac Disk Copy Disk Image File	Encapsularion
345	Appleworks File	Word processor
346	Omni Outliner (OO3) File	Word processor
347	Omni Outliner (OPML) File	Word processor
348	Omni Graffle XML File	Vector graphic
349	Apple Photoshop Document	Raster image
350	Apple Binary Property List	Miscellaneous
351	Apple iChat Format	Word processor
352	Omni Outliner (OOUTLINE) File	Word processor
353	Bzip 2 Compressed File	Encapsulation
354	ISO-9660 CD Disc Image Format	Encapsulation
355	Xerox DocuWorks	Word processor
356	RealMedia Streaming Media	Movie
357	AC3 Audio File Format	Sound
358	Nero Encrypted File	Encapsulation
359	SolidWorks	Vector graphic

Major Formats, continued

Number	Format	File Class
362	UniGraphics NX	Vector graphic
366	Extensible Forms Description Language	Presentation
367	Apple XML Property List	Miscellaneous
368	OneNote Note Format	Presentation
370	Digital Imaging and Communications in Medicine (DICOM)	Raster image
371	Expert Witness Compression Format	Encapsulation
372	Shell Scrap Object File	Encapsulation
373	Microsoft Project 2007	Time scheduling
374	Microsoft Publisher 98–	Desktop publishing
375	Skype Log File	Word processor
376	Lotus Notes Bitmap Format (DXL embedded images)	Raster image
377	Health level7 message	Word processor
378	Microsoft Outlook Offline Storage File	Encapsulation
379	Open Publication Structure eBook	Word processor
380	Microsoft Outlook Express DBX	Encapsulation
381	BlackBerry Activation File	Word processor
382	Disk Image	Encapsulation
383	Milestone	Raster Image
384	RealLegal E-Transcript File	Word processor
385	PostScript Type 1 Font	Font
386	Ghost Disk Image File	Encapsulation
387	JPEG-2000 JP2 File Format Syntax (ISO/IEC 15444-1)	Raster Image
388	Unicode HTML	Word processor
389	Microsoft Compiled HTML Help	Encapsulation
390	Documentum EMCMP	Encapsulation
393	JBIG2 File	Raster image
395	AD1 Evidence file	Encapsulation

Major Formats, continued

Number	Format	File Class
397	Group Wise File Surf email	Encapsulation
402	ARJ	Encapsulation
409	Microsoft Outlook for Macintosh	Encapsulation
412	Microsoft Outlook vCard Contact	Word processor
414	Microsoft Outlook iCalendar	Encapsulation
418	Apple iWork 2013 Pages	Word processor
419	Apple iWork 2013 Numbers	Spreadsheet
420	Apple iWork 2013 Keynote	Presentation
421	XZ	Encapsulation
427	B1	Encapsulation
428	MP4	Movie
429	Rar5	Encapsulation
430	PTC Creo	Vector graphic
431	Keyhole Markup Language	
432	Zipped Keyhole Markup Language	
433	Wireless Markup Language	
435	Star Office Writer Text	
436	Star Office Calc Spreadsheet	
437	Star Office Impress Presentation	
438	Star Office Math	

1 If the major version is 100, the file format is XHTML.

2 The major version determines whether the Microsoft Office XML file is a Word, Excel or Visio document. The major version for each format is as follows:

Word: 100

Excel: 101

Visio: 110

File Classes

Attribute Number	File Class
0	No file class

File Classes, continued

Attribute Number	File Class
01	Word processor
02	Spreadsheet
03	Database
04	Raster image
05	Vector graphic
06	Presentation
07	Executable
08	Encapsulation
09	Sound
10	Desktop publishing
11	Outline/planning
12	Miscellaneous
13	Mixed format
14	Font
15	Time scheduling
16	Communications
17	Object module
18	Library module
19	Fax
20	Movie
21	Animation

Minor Formats

Attribute Number	Minor Format
00	Minor format not defined
01	Standard
02	Book

Minor Formats, continued

Attribute Number	Minor Format
03	Chart
04	Macro
05	Text
06	Binary
07	PC
08	Windows
09	DOS
10	Macintosh
11	RGB
12	TIFF
13	IFF
14	Experimental
15	Format Information
16	RLE
17	Symbol
18	Old
19	Footnote
20	Style
21	Palette
22	Configuration
23	Activity
24	Resource
25	Calculation
26	Glossary
27	Spelling
28	Thesaurus
29	Hyphenation

Minor Formats, continued

Attribute Number	Minor Format
30	Miscellaneous
31	UNIX
32	VAX
33	Driver
34	Archive

Appendix F: List of Required Files for Redistribution

This section lists the Filter files that can be redistributed in your applications under the licensing agreement. These files are in the directory *install\OS\bin*, where *install* is the path name of the Filter installation directory and *OS* is the name of the operating system.

NOTE:

On Windows systems, the libraries are .dll files. On UNIX systems, the libraries are .so, .a, or .sl files.

Core Files

The following core files can be redistributed with your application.

File	Description
formats.ini	Initialization file. For more information on this file, see Determine Format Support, on page 176 .
FilterDotNet.*	Required by .NET API.
KeyViewFilter.*	Required by the Java API.
kpifcnvt.*	For presentation graphics, converts from one picture format to another.
kpifutil.*	Utility for handling the internal picture interchange format for presentation graphics.
kvxtract.*	File Extraction API.
kvfilter.*	Filter API.
kvolefio.*	Embedded OLE object writer.
kvutil.*	Internal KeyView utility functions.
kvxpgsa.*	Interface between presentation readers and kvfilter. Required to extract metadata from AutoCAD files.
kvxssa.*	Interface between spreadsheet readers and kvfilter.
kvxwpsa.*	Interface between word processing readers and kvfilter.
kwad.*	File auto-recognition module.
txtcnv.*	Converter for document token stream.

Support Files

The following support files can be redistributed with your application.

File	Description
bentofio.*	Required by 1123sr and kpprzrdr.
cbmap.map	Character mappings for Adobe Portable Document Format (PDF).
chartbls.ux	Character mappings.
chmdl1.*	Required by chmsr.
kppng.*	Required for ZLIB decompression.
kvxconfig.ini	Contains element extraction settings for XML files.
kvoop.*	Required for out-of-process filtering.
kvthread.*	Required for multithreaded out-of-process filtering.
kv.lic	Contains license information for KeyView products. This file is opened and validated when a KeyView API is used.
MSVCP60.*	Microsoft Visual C++ Runtime library V6.0.
msvcrt.*	Microsoft Visual C Runtime library.
wpmap.*	Extended character mapping for WordPerfect and Corel Presentation.
xmlsh.*	Contains a library of content handlers for each XML file type. Required by the Expat XML parser.

Document Readers

The following readers can be redistributed with your application.

File	Description
ad1sr.*	AD1 Evidence file reader
afsr.*	ASCII reader
aiffsr.*	Audio Interchange Format File (AIFF) reader
asfsr.*	Advanced Systems Format reader
assr.*	Applix Spreadsheet reader
awsr.*	Applix Word reader

File	Description
b11sr.*	B1 archive reader
bkfsr.*	Microsoft Backup File reader
bmpsr.*	Windows bitmap (BMP) reader
bzip2sr.*	Bzip2 reader
cabsr.*	Microsoft Cabinet format reader
cebsr.*	Founder Chinese E-paper Basic reader
chmsr.*	Microsoft Compiled HTML Help reader
csvsr.*	Comma-Separated Values reader
dbfsr.*	dBase Database reader
dbxsr.*	Microsoft Outlook Express DBX reader
dcasr.*	Document Content Architecture/Revisable Form Text (DCA/RFT) reader
dcmsr.*	Digital Imaging and Communications in Medicine (DICOM) reader
difsr.*	Data Interchange Format reader
dmgsr.*	Mac Disk Copy Disk Image File reader
dw4sr.*	DisplayWrite reader
dx1sr.*	Domino XML Language reader
em1sr.*	Microsoft Outlook Express (EML) reader. This is used to filter EML files when the MBX reader is not licensed.
emxsr.*	Legato EMailXtender (EMX) reader
encasesr.*	Expert Witness Compression Format (EnCase) v6 reader
encase2sr.*	Expert Witness Compression Format (EnCase) v7 reader
entsr.*	Microsoft Entourage Database Format reader
epubsr.*	Open Publication Structure eBook reader
foliosr.*	Folio Flat File reader
gifsr.*	Graphics Interchange Format (GIF) reader
gwfssr.*	GroupWise FileSurf reader
h17sr.*	Health level7 reader (metadata only)
htmsr.*	HTML and XHTML reader

File	Description
hwpsr.*	Hangul 97 reader
hwposr.*	Hangul 2002, 2005, 2007 reader
ichatsr.*	Apple iChat Log reader
icssr.*	Microsoft Outlook iCalendar reader
isosr.*	ISO-9660 CD Disc Image Format reader
iwwpsr.*	Apple iWork Pages reader
iwsssr.*	Apple iWork Numbers reader
jp2000sr.*	JPEG 2000 metadata reader
jpgsr.*	JPEG metadata reader
jtdsr.*	JustSystems Ichitaro reader
kpagrdr.*	Applix Presentations reader
kpCATrdr.*	CATIA format reader
kpcgmrdr.*	Computer Graphics Metafile reader
kpDWGrdr.*	AutoCAD Drawing format reader
kpDXFrdr.*	AutoCAD Drawing Exchange format reader
kpemfrdr.*	Enhanced Metafile reader
kpGFLrdr.*	Omni Graffle reader
kpgifrdr.*	Graphic Interchange Format (GIF) reader
kpIWGrdr.*	Apple iWork Keynote reader
kpmssrdr.*	Microsoft Office Drawing Objects (office 97, 2000, and XP) reader
kpODArdr.*	AutoCAD reader (Windows only)
kpodfrdr.*	Oasis Open Document Format presentation (ODP) reader
kpONErdr.*	Microsoft OneNote reader
kpp40rdr.*	Microsoft PowerPoint PC 4.0 and PowerPoint Mac reader
kpp95rdr.*	Microsoft PowerPoint 95 reader
kpp97rdr.*	Microsoft PowerPoint 97 and higher reader
kppctrdr.*	Macintosh Quick Draw Picture (PICT) reader
kppicrdr.*	Pictor PC Paint (PIC) reader

File	Description
kpppxrdr.*	Microsoft PowerPoint XML reader 2007
kpprerdr.*	Lotus Freelance Graphics for Windows V2.0 reader
kpprzrdr.*	Lotus Freelance Graphics 96/97/98 reader
kpshwrdr.*	Corel Presentations reader
kpugrdr.*	Unigraphics (UG) NX reader
kpgw2rdr.*	WordPerfect Graphics 2 reader
kpwmfrdr.*	Windows Metafile reader
kpwpgrdr.*	WordPerfect Graphics 1 reader
kpXFDLrdr.*	Extensible Forms Description Language reader
kvgzsr.*	GZIP reader
kvhqxsr.*	BinHex reader
kvzeesr.*	UNIX Compress reader
l123sr.*	Lotus 123 v96/97/98 reader
lasr.*	Lotus AMI Pro reader
ltbenn30.dll	Lotus Word Pro support (supported on Windows x86 platform only)
ltscsn10.dll	Lotus Word Pro support (supported on Windows x86 platform only)
lwpapin.dll	Lotus Word Pro support (supported on Windows x86 platform only)
lwppann.dll	Lotus Word Pro support (supported on Windows x86 platform only)
lwpsr.dll	Lotus Word Pro reader (supported on Windows x86 platform only)
lzhshr.*	Microsoft Compression Folder reader
macbinsr.*	MacBinary reader
mbsr.*	Microsoft Word Macintosh reader
mbxsr.*	Mailbox (MBX) and Microsoft Outlook Express (EML) reader ¹
mdbsr.*	Microsoft Access reader
mhtsr.*	MIME HTML reader
mifsr.*	Adobe Maker Interchange reader

¹This reader is an advanced feature and is sold and licensed separately from KeyView Filter SDK. See [License Information, on page 17](#)

File	Description
misr.*	Microsoft Word 2 reader
mp3sr.*	MP3 reader for metadata extraction reader
mpeg4sr.*	MPEG-4 Audio file reader
mppsр.*	Microsoft Project reader
msgsr.*	Microsoft Outlook (MSG) reader
mspubsr.*	Microsoft Publisher reader
msw6sr.*	Microsoft Works 6 and 2000 reader
mswsr.*	Microsoft Works V1 and 2 reader
multiarcsr	ARJ Reader
mw6sr.*	Microsoft Word 95 reader
mw8sr.*	Microsoft Word 97, 2000, and XP reader
mwsr.*	Microsoft Word for DOS and Microsoft Write reader
mwssr.*	Microsoft Works Spreadsheet reader
mwxsr.*	Microsoft Word 2007 XML reader
nsfsr.*	Lotus Notes database reader 1
oa2sr.*	Fujitsu Oasys reader
odfsssr.*	Oasis Open Document Format spreadsheets (ODS) reader
odfwpsr.*	Oasis Open Document Format word processing (ODS) reader
olesr.*	Embedded OLE object reader
olmsr.*	Microsoft Outlook for Macintosh reader
onmsr.*	Legato EMailXtender Native Message reader
oo3sr.*	Omni Outliner reader
pdfsr.*	Adobe Portable Document Format file (PDF) reader
pffsr.*	Microsoft Outlook Offline Storage File reader
pngsr.*	Portable Network Graphics (PNG) reader
pstsr.dll	Microsoft Outlook Personal Folders file MAPI-based reader (supported on Windows platform only) 1
pstnsr.*	Microsoft Outlook Personal Folders file native reader 1

File	Description
qpssr.*	Corel Quattro Pro spreadsheet reader
qpwsr.*	Corel Quattro Pro version X4 spreadsheet reader
rarsr.*	RAR Archive reader
riffsr.*	Microsoft WAVE reader
rtfsr.*	Microsoft Rich Text reader
skypesr.*	Skype log file reader
sosr.*	StarOffice/OpenOffice reader
sunadsr.*	Sun Audio Data reader
swfsr.*	Macromedia Flash reader
tarsr.*	Tape archive reader
tifsr.*	TIFF reader (metadata only)
tnefsr.*	Transfer Neutral Encapsulation Format
unihtmlsr.*	Unicode HTML reader
unisir.*	Unicode reader
unzip.*	Zip file reader
utf8sr.*	UTF-8 reader
uudsr.*	UUEncoding reader
vcfsr.*	Microsoft Outlook vCard Contact reader
vsdsr.*	Microsoft Visio reader
wkssr.*	Lotus 123 v2.0 through 5.0 reader
wosr.*	WordPerfect 5.x reader
wp6sr.*	WordPerfect 6.0 through 10.0 reader
wpmsr.*	WordPerfect for Macintosh reader
xlsbsr.*	Microsoft Office 2007 Excel Binary Format reader
xlssr.*	Microsoft Excel reader
xlsxsr.*	Microsoft Excel 2007 XML reader
xmlsr.*	Generic XML reader
xpssr.*	XML Paper Specification reader

File	Description
xywsr.*	XYWrite reader
yimsr.*	Yahoo! Instant Messenger reader
z7zsr.*	7-Zip reader

Appendix G: Develop a Custom Reader

This section describes how to develop a reader for a format not supported by KeyView.

• Introduction	207
• How to Write a Custom Reader	208
• Development Tips	218
• Functions	219

Introduction

The Filter SDK enables you to write custom readers for formats not directly supported by KeyView. A reader is required to parse the file format and generate a KeyView token stream, which represents the content and format of the document. Filter can then use this token stream to generate a text version of the original document. The readers interact with a structured access layer and a writer to generate a text file in Filter, an HTML file in HTML Export, an XML file in XML Export, and a near-to-original view of the document in the Viewing SDK.

The complexity of a custom reader depends on the file format used by the source document type. A simple reader extracts only the textual content, but ignores formatting and all other non-textual content. Readers of increasing complexity must address one or more of the following:

- formatting (including fonts, foreground and background colors, paragraph borders and shading, character and paragraph styles)
- tables
- lists
- headers
- footers
- footnotes
- endnotes
- graphics
- bookmarks to internal links
- hyperlinks to external documents or webpages
- other structures, such as a table of contents or index

Even a simple reader might have to parse the following components of a document:

- word processing commands or tags
- encrypted or encoded text
- multiple character sets
- text modified, but retained within the file
- text displayed in an order other than its physical occurrence within the source file

It is very important to fully understand the file specification for the file format used by the document. This is essential in determining how to parse the source file and generate a token stream that accurately and effectively represents the original document.

Within Filter, the custom reader must interact with a structured access layer and the format detection API, which in turn interacts with the top-level API. For a description of the Filter architecture, see [Architectural Overview, on page 21](#).

The custom reader must have a module definition file (*.def) that defines the exported API function calls. In addition, the `formats.ini` file must be modified to identify the custom reader and its associated format detection function.

See the source code for the sample custom reader (`utf8sr`), which parses plain text files encoded in UTF-8. The source code is in the directory `install/samples/utf8sr`, where `install` is the path name of the Filter installation directory.

How to Write a Custom Reader

Two include files define the requirements for a custom reader: `kvcfsr.h` and `kvtoken.h`. The definitions of the KeyView tokens are in `kvtoken.h`. For more information on tokens, see [Token Buffer, on the next page](#). The file `kvcfsr.h` defines two structures: `TPrereaderInterface` and `adTPDocInfo`.

The `TPrereaderInterface` structure defines the API functions implemented by the custom reader. For basic readers, only the first four functions must be implemented. These functions are called by the structured access layer to parse the source file and generate the token stream.

All readers must be threadsafe. This means that global variables must not be used. To pass information between functions, it is necessary to define a "global" context structure that stores all information required throughout the life of the DLL. The initial parameter of all but one of the `TPrereaderInterface` functions is a pointer to a global context structure defined for the custom reader.

The `adTPDocInfo` structure defines the information required for the format detection API, which associates the custom reader with the required file format.

Naming Conventions

Use the following naming conventions for functions and files:

- The initial letters of the custom reader file name should identify the file format being parsed. For example, `pdf` for Adobe PDF files, `rtf` for RTF files, and `xls` for Microsoft Excel files. In the examples in this appendix, this is represented by `xxx`.
- The name of the shared library must end with the letters `sr`.
- The name of the exported functions in the module definition file must be `xxxGetReaderInterface` and `xxxsrAutoDet`.

NOTE:

The letters `sr` are excluded from `xxxGetReaderInterface`, but are included in `xxxsrAutoDet`.

Basic Steps

The basic steps for developing a custom reader are as follows.

To develop a custom reader

1. Design the global context structure.
2. Write the basic API functions:

- `xxxAllocateContext()`
- `xxxInitDoc()`
- `xxxFillBuffer()`
- `xxxFreeContext()`
- `xxxCharSet()`
- `xxxsrAutoDet()`

From within the `xxxFillBuffer()` function, it is necessary to call other functions that repeatedly read a chunk of a source file, parse the chunk, and generate a token stream until the entire source file is processed.

3. Map all but the last function to the `TPReaderInterface` structure.
4. Write the module definition file (`*.def`), exporting the reader interface and format detection functions.
5. Modify the `formats.ini` file to identify the custom reader and its associated format detection function. See `xxxsrAutoDet()`, on page 219. For example, the following lines would be added to the `[Formats]` section of the `formats.ini` file for the UTF-8 reader:

```
456.1.0.0=utf8
[CustomFilters]
1=utf8sr
```

Token Buffer

Filter technology parses the native file structure to generate an intermediate stream called a *token buffer*. The token buffer consists of multiple sequences of tokens, which are defined in `kvtoken.h` and listed below.

```
#define KVT_TEXT          0x00 /* PutText() */
#define KVT_PARAINFO      0x01 /* SetParaInfo() */
#define KVT_SETTABS       0x02 /* SetTabs() */
#define KVT_TAB           0x03 /* Tab() */
#define KVT_MODE          0x04 /* SetMode() */
#define KVT_PARASPACE     0x05 /* SetParaSpace() */
#define KVT_ROWDEFN       0x06 /* DefineRow(), EndTable() */
#define KVT_COLUMNS       0x07 /* StartColumns(), etc. */
#define KVT_CELLSTART     0x08 /* NextCell() */
#define KVT_BITMAP        0x09 /* Reserved for annotations. */
#define KVT_PAGEOBJ       0x0A /* PutHeader(), PrintPage(), etc.*/
```

```
#define KVT_NOOP          0x0B /* Just skip a BYTE. */
#define KVT_PAGE_BREAK   0x0C /* PageBreak() */
#define KVT_PARA_BREAK   0x0D /* ParaEnd() */
#define KVT_LINE_BREAK   0x0E /* LineBreak() */
#define KVT_SET_FONT     0x0F /* SetFont() */
#define KVT_PAGE         0x10 /* SetPageInfo() */
#define KVT_HOTSPOT      0x11 /* StartHotSpot() */
#define KVT_LINESPACE    0x12 /* SetLineSpacing() */
#define KVT_COLOR        0x13 /* VESetTextColor(),VESetBkColor()*/
#define KVT_PICTURE      0x14 /* PutPicture() */
#define KVT_CELLMERGE    0x15 /* MergeCells() */
#define KVT_RULE         0x16 /* HorzRule() */
#define KVT_PATTERN      0x17 /* StartPattern(), etc. */
#define KVT_BORDER       0x18 /* StartParaBorder(), etc. */
#define KVT_HEADING      0x19 /* PutParaHeading() */
#define KVT_LISTING      0x1A /* StartList(), etc. */
#define KVT_CHARSET      0x1B /* SetCharSet() */
#define KVT_STYLE        0x1C /* PutCharStyle(), PutParaStyle()*/
#define KVT_BIDI         0x1D /* Set Bidirectional text */
#define KVT_LOCALE       0x1E /* Set locale of a document */
#define KVT_ZONE         0x1F /* StartZone(), EndZone() */
#define KVT_POSITION     0x20 /* SetPosition(), etc. */
#define KVT_AUTOREC      0x21 /* Reserved for Internal Use */
#define KVT_METADATA     0x22 /* Rsserved for Internal Use */
#define KVT_BYTEORDER    0x23 /* SetByteOrder() */
#define KVT_PARASPACEAUTO 0x24 /* SetParaSpaceAuto() */
#define KVT_ATTACH       0x25 /* PutAttachment() */
#define KVT_TOCPRINTIMAGE 0x26 /* StartTOCPrintImage(), etc. */
#define KVT_STREAM       0x27 /* PutStream(),Reserved */
#define KVT_REVISIONMARK 0x28 /* StartRevisionMark(),
EndRevisionMark(), SetRMAuthor(), SetRMDateTime() */
#define KVT_DOCXTRINFO   0x29 /* SetDocXtrInfo() */
#define KVT_PCTEMDFT     0x30 /* SetPctEmdFt() */
```

A token is a single-byte identifier that corresponds to attributes in a document. Each token has one or more associated macros that provide detailed information about an attribute. Many of these tokens define components of the document, such as page margins, line indentation, and foreground and background color. Collectively, these are referred to as the *state* of the document. This state changes as the document is parsed.

Macros

Some of the macros are simple while others are complicated. An example of a simple macro is ParaEnd (pcBuf) which terminates the current paragraph.

```
#define ParaEnd(pcBuf)      \
{                           \
    *pcBuf++ = KVT_PARA_BREAK; \
}
```

```
        KVT_PUTINT(pcBuf, KVTSIZE_PARA_BREAK);        \
    }
}
```

In Filter SDK, this generates an 0x0d, 0x0a pair of bytes on a Windows machine. In HTML Export this can generate a <p style="..."> element, depending on the value of other paragraph attributes.

One of the more complicated macros is PutPictureEx().

```
#define PutPictureEx(pcBuf, lpszKey, cx, cy, flags,        \
    scaleHeight, scaleWidth,                             \
    cropFromL, cropFromT, cropFromR, cropFromB,         \
    anchorHorizontal, anchorVertical, offsetX, offsetY)\
{                                                        \
    PutPic(pcBuf, lpszKey, cx, cy, flags,               \
    scaleHeight, scaleWidth,                           \
    cropFromL, cropFromT, cropFromR, cropFromB,         \
    anchorHorizontal, anchorVertical, offsetX, offsetY, \
    180, 0, 180, 0, -1, 0, 0, 0, 0)                    \
}
```

You can generate a representation of the token stream by running `filtertest.exe` with the `-d` command-line option. This stream does not include the tokens generated for headers or footers. The `filtertest.exe` is in the directory `install\samples\utf8\bin`, where `install` is the path name of the Filter installation directory.

Reader Interface

All custom readers use the reader interface defined in `kvcfsr.h`. The members of this structure are:

```
fpAllocateContext()
fpInitDoc()
fpFillBuffer()
fpFreeContext()
fpHotSpothit()
fpGetSummaryInfo()
fpOpenStream()
fpCloseStream()
fpGetURL()
fpGetCharSet()
```

NOTE:

`fpHotSpothit()` and `fpGetURL()` are currently reserved and must be NULL.

Function Flow

The structured access layer calls the functions as follows:

1. `fpAllocateContext()` is called and returns a pointer to the global context structure.

2. After further processing within the structured access layer, `fpInitDoc()` is called. This function performs all required initialization for the global context structure and then returns control to the structured access layer.
3. After further processing within the structured access layer, the `fpFillBuffer()` function is called repeatedly until the document is completely parsed.
4. Finally, `fpFreeContext()` is called. This function frees all memory allocated within the custom reader and then returns control to the structured access layer.

Related Topics

- [Functions, on page 219](#)

Example Development of `fffFillBuffer()`

The following is an example of how the `fpFillBuffer()` function in `foliosr` could be developed. The example demonstrates how the code changes as limitations of the implementation are identified. With each implementation, code revisions are shown in bold.

Implementation 1—`fpFillBuffer()` Function

```
/******  
*Function: fffFillBuffer()  
*Summary: Read fff input from stream and parse into kvtoken.h codes  
*****/  
int pascal _export fffFillBuffer(  
    void    *pCfContext,  
    BYTE    *pcBuf,  
    UINT    *pnBufOut,  
    int     *pnPercentDone,  
    UINT     cbBufOutMax )  
{  
    BOOL bRetVal;  
    TPffffGlobals *pContext = (TPffffGlobals *)pCfContext;  
    pContext->pcBufOut = pcBuf;  
    fffReadSourceFile(pContext);  
    bRetVal = fffProcessBuffer(pContext, pcBuf);  
    *pnPercentDone = (int)(pContext->unTotalBytesProcessed *  
        (UINT)100 / pContext->unFileSize);  
    *pnBufOut = (UINT)(pContext->pcBufOut - pcBuf);  
    return (bRetVal ? KVERR_Success : KVERR_General);  
}
```

The parameters in `fffFillBuffer()` are as follows:

Parameter	In/Out	Description
<code>pCfContext</code>	In	A pointer to the context structure of the custom reader.
<code>pcBuf</code>	In/Out	A pointer to the token output buffer.

Parameter	In/Out	Description
pnBufOut	Out	A pointer to the number of bytes written to the output buffer.
pnPercentDone	Out	A pointer to the percentage complete.
cbBufOutMax	In	The maximum number of bytes that the token output buffer can hold.

Structure of Implementation 1

1. The local variable `pContext` is set to the address of the `pCFContext` void pointer, cast to a pointer to the global context structure for the reader. This provides access to all members of this structure.
2. After setting the `pContext` variable, a call is made to read the source file.
3. Next, a call is made to `fffProcessBuffer()`. The second parameter in the call is a pointer to the token output buffer. If this call fails, usually because of memory allocation errors, it returns `FALSE`.
4. The percentage complete is calculated.
5. The number of BYTES written to the token output buffer is calculated. This is based on the value of `pContext->pcBufOut`, which is increased each time a token is written to the buffer.
6. The function returns to the structured access layer.
7. Subsequent calls to `fffFillBuffer()` are made by the structured access layer until the percentage complete is 100.

Problems with Implementation 1

- There is a limit to the size of the token output buffer, typically 4 KB. If `fffProcessBuffer()` generates a token stream larger than this, there is a memory overflow. If `fffProcessBuffer()` generates a small token stream and the entire file has not been read, the output token buffer is underutilized.
- It might not be possible to process the entire input buffer from the source file because of boundary conditions. An example of a "boundary condition" is when the input buffer terminates part way through a control sequence in the original document. Another file read operation is required before the complete control sequence can be parsed.
- This function might be interrupted by other calls from the structured access layer to process headers, footers, footnotes, and endnotes, or to retrieve the document summary information. This can cause values of variables in the global context to change, and the source file to be repositioned.

Implementation 2—Processing a Large Token Stream

Implementation 2 addresses the problem of processing a token stream that is larger than the output buffer size limit.

```

/*****
* Function:   fffFillBuffer()
* Summary:   Read fff input from stream and parse into kvtoken.h codes
*****/
int pascal _export fffFillBuffer(

```

```
void      *pCFContext,
BYTE      *pcBuf,
UINT      *pnBufOut,
int       *pnPercentDone,
UINT      cbBufOutMax )
{
    BOOL bRetVal = TRUE;
    TPfffGlobals *pContext = (TPfffGlobals *)pCFContext;
    pContext->pcBufOut = pcBuf;
    pContext->cbBufOutMax = 9 * cbBufOutMax / 10; /* Process the portion of the
fff file that is in the input buffer but do * not return from the fffFillBuffer()
function unless the output buffer is * at least 90% full. If any of the memory
allocations fail during the * execution of fffProcessBuffer(), bRetVal will be
set to FALSE, resulting * in this conversion failing "gracefully".
*/
    do
    {
        if( pContext->bBufOutFull )
        {
            pContext->bBufOutFull = FALSE;
        }
        else
        {
            fffReadSourceFile(pContext);
        }
        bRetVal = fffProcessBuffer(pContext, pcBuf);
        *pnPercentDone = (int)(pContext->unTotalBytesProcessed *
(UINT)100 / pContext->unFileSize);
    }while( bRetVal && !pContext->bBufOutFull && *pnPercentDone < 100 );
    *pnBufOut = (UINT)(pContext->pcBufOut - pcBuf);
    return (bRetVal ? KVERR_Success : KVERR_General);
}
```

Structure of Implementation 2

1. cbBufOutMax is used to set pContext->cbBufOutMax. This is used in fffProcessBuffer() to monitor how full the token output buffer becomes as the source file is processed.
2. When the source file input buffer has been processed, fffProcessBuffer() returns, and the percentage complete is calculated.
3. If the token output buffer is not filled to a value greater than pContext->cbBufOutMax, pContext->bBufOutFull remains set to FALSE, and if the percentage complete is less than 100, the do-while loop is re-entered without returning from this function to the structured access layer. There is another call to fffReadSourceFile(), followed by fffProcessBuffer().

4. When the token output buffer is filled to a value greater than `pContext->cbBufOutMax`, `pContext->bBufOutFull` is set to `TRUE`. In this case, the `do-while` loop ends, the number of bytes written to the token output buffer is calculated, and control returns to the structured access layer.
5. The structured access layer continues to make calls to `ffffFillBuffer()` until the entire source file is processed.
6. Each time the structured access layer calls `ffffFillBuffer()`, another empty token output buffer is provided for the custom reader to use.
7. If the previous call to `ffffFillBuffer()` exited because the previous token output buffer exceeded allowable capacity, `pContext->bBufOutFull` is reset to `FALSE` and no call is made to read the next buffer from the input source file.

Problems with Implementation 2

- It might not be possible to process the entire input buffer from the source file because of boundary conditions.
- This function might be interrupted by other calls from the structured access layer to process headers, footers, footnotes, or endnotes, or to retrieve the document summary information. This can cause values of variables in the global context to change, and the source file to be repositioned.

Boundary Conditions

A boundary condition can result from many situations arising from input file processing. For example, the input buffer might end with an incomplete command. In Folio flat files, this could be an incomplete element. In other word processing documents, a boundary condition might result from an incomplete control sequence, a split double-byte character, or a partial UTF-7 or UTF-8 sequence. These can be handled jointly by `ffffProcessBuffer()`, which must detect the boundary condition, and `ffffReadSourceFile()`.

The following example shows partial code used in `ffffReadSourceFile()`:

```
/*
 *
 * Function:    ffffReadSourceFile()
 *
 */
int pascal ffffReadSourceFile(TPffffGlobals *pContext)
{
    int nBytes;
    /* Transfer remaining data to beginning of buffer prior to next read */
    if( pContext->nResidualBytes )
    {
        memcpy(pContext->cInputBuf, pContext->pcBufIn, pContext->nResidualBytes);
    }
    /* Read from file, without over-writing any text from the previous buffer */
    nBytes = (*pContext->pIO->kwReadFunc)(pContext->pIO,
        pContext->cInputBuf + pContext->nResidualBytes,
        BUFFERSIZE - pContext->nResidualBytes);
    /* Update input buffer control parameters */
}
```

```
pContext->unTotalBytesRead += (UINT)nBytes;
pContext->pcBufIn = pContext->cInputBuf;
pContext->pcBufInMax = pContext->pcBufIn + pContext->nResidualBytes + nBytes;
pContext->nResidualBytes = 0;
return nBytes;
}
```

If `ffffProcessBuffer()` is unable to process the entire input source file buffer, it sets the value for `pContext->nResidualBytes`. When the next call to `ffffReadSourceFile()` is made, any residual bytes are copied to the beginning of the input source file buffer, and the number of bytes to be read is reduced to make sure that this buffer does not overflow.

A good way to test the code for boundary conditions is to vary the size of `BUFFERSIZE` and make sure that the results remain consistent.

NOTE:

With `ReadSourceFile()`, the source file can be read by calls to retrieve header or footer information. If this occurs, the value for `pContext->unTotalBytesRead` is incorrect.

Implementation 3—Interrupting Structured Access Layer Calls

Implementation 3 addresses the problem of boundary conditions and interrupting calls from the structured access layer.

```
/* *****
* Function:   ffffFillBuffer()
* Summary:    Read fff input from stream and parse into kvtoken.h codes
* *****/
int pascal _export ffffFillBuffer(
    void      *pCfContext,
    BYTE      *pcBuf,
    UINT      *pnBufOut,
    int       *pnPercentDone,
    UINT      cbBufOutMax )
{
    double dTotalBytesProcessed, dFileSize;
    BOOL bRetVal = TRUE;
    TPffffGlobals *pContext = (TPffffGlobals *)pCfContext;
    pContext->pcBufOut = pcBuf;
    pContext->cbBufOutMax = 9 * cbBufOutMax / 10;
    /* Process the portion of the fff file that is in the input buffer but do
    * not return from the ffffFillBuffer() function unless the output buffer is
    * at least 90% full. If any of the memory allocations fail during the
    * execution of ffffProcessBuffer(), bRetVal will be set to FALSE, resulting
    * in this conversion failing "gracefully". */
    do
    {
        if( pContext->bBufOutFull )
        {
            pContext->bBufOutFull = FALSE;

```



```
    }
    else
    {
        fffReadSourceFile(pContext);
    }
    bRetVal = fffProcessBuffer(pContext, pcBuf);
    if( pContext->bHeaderCompleted )

{
    *pnPercentDone = 100;
    pContext->bHeaderCompleted = FALSE;
}
else if( pContext->bFooterCompleted )

{
    *pnPercentDone = 100;
    pContext->bFooterCompleted = FALSE;
}
else

{
    if( pContext->unTotalBytesProcessed >= pContext->unFileSize )
    {
        *pnPercentDone = 100;
    }
    else if( pContext->unFileSize < FFF_MAX_ULONG )
    {
        *pnPercentDone = (int)(pContext->unTotalBytesProcessed *
(UINT)100 / pContext->unFileSize);
    }
    else

{
        dTotalBytesProcessed = pContext->unTotalBytesProcessed;
        dFileSize = pContext->unFileSize;
        *pnPercentDone = (int)(dTotalBytesProcessed * 100 / dFileSize);
    }
}
}while( bRetVal && !pContext->bBufOutFull && *pnPercentDone < 100 );
*pnBufOut = (UINT)(pContext->pcBufOut - pcBuf);
return (bRetVal ? KVERR_Success : KVERR_General);
}
```

Structure of Implementation 3

- The most significant change in Implementation 3 is the addition of the code that checks whether the processing of the header or footer is complete. The variables for `pContext->bHeaderCompleted` and `pContext->bFooterCompleted` are set to TRUE in `fffProcessBuffer()` when a header or footer is

processed and the end of that portion of the document is reached.

- The other piece of code added in Implementation 3 is unique to `foliosr`. Folio files can be 50 MB or larger. Therefore, an unsigned integer is too small to accurately calculate the percentage complete. If the file size exceeds `FFF_MAX_ULONG`, which is defined as `(UINT)(0xFFFFFFFF / 0x64)`, the doubles are used for that calculation.
- Prior to returning, the token output buffer is as full as possible and never overflows. The minimum number of calls is made.

Development Tips

- Avoid unnecessary initialization.

The context variable is allocated in `fpAllocateContext()`. This structure must be immediately `memset()` to zero. This sets all `BOOL` values to `FALSE`, all pointers to `NULL`, and all integers to 0. Only non-zero, non-`NULL` and `BOOL`s that must be `TRUE` need to be initialized. This is best done in `fpInitDoc()`.

- Know where you are in the input source file.

If you are processing headers, footers, notes, or (in the case of `rtfsr`) tables, you must be able to reposition the file pointer as required.

- Check buffer boundaries continuously.

Whenever you advance through the buffer, you need to know whether there is enough of the input stream to completely process the current command. If not, you need to append the next section of the input file before continuing.

- Strive for a "clean" token stream.

Use `filtertest` with the `-d` command-line option to generate a *token* version of the document. If there are redundant tokens, the reader is producing an inefficient token stream. You can keep the token stream free from redundancies by storing the state of the document and then applying the changes only when content is encountered. Content can be text, tabs, or picture objects. The `filtertest.exe` is in the directory `install\samples\utf8\bin`, where `install` is the path name of the Filter installation directory.

- Avoid large `switch()` statements whenever possible. They make both development and debugging more complicated than necessary. If there is a fixed set of commands, consider using a hash table that enables you to quickly identify a pointer to the function that handles that command.
- Filtering document metadata is a separate process.

Remember that `fpGetSummaryInfo()` is a completely separate process from the rest of your code. It creates its own context variable structure. It does not have to call `fpFillBuffer()`.

- Use caution when processing headers, footers, and notes.

If you need to process these items, the structured access layer calls `fpOpenStream()` and `fpCloseStream()`. It is critical that you save the state of your document and the file pointer position prior to returning from `fpOpenStream()`. Prior to returning from `fpCloseStream()`, you must restore the file pointer and the previous state of your document.

- Test your code.

The structured access layer for each SDK is unique. Test your code in Filter SDK, Export SDK, and Viewing SDK.

Functions

This section describes the functions used by custom readers to manage the source file and generate token streams required to convert a document.

xxxxsrAutoDet()

This function analyzes the source document and determines whether the detected file format requires the custom reader. It is called only when the [CustomFilters] section of the `formats.ini` file contains an entry identifying the complete file name of the custom reader. For more information on the `formats.ini` file, see [File Format Detection, on page 176](#).

Syntax

```
Bool pascal _export xxxxsrAutoDet(  
    adTPDocInfo    *pTPDocInfo,  
    KPTPIOobj      *pIO)
```

Arguments

`pTPDocInfo` A pointer to the `adTPDocInfo` structure provided by the structured access layer.

`pIO` A pointer to the I/O stream object for the document processed.

Returns

- TRUE if the file format matches that of the custom reader.
- FALSE if the file format does not match that of the custom reader.

Discussion

- Typically, only the first 1 KB of the file is read into a buffer and analyzed to determine if it matches the file format of the custom reader. If a match is determined, the following four members of the `adTPDocInfo` structure must be assigned before returning TRUE:

<code>adClass</code>	Must be set to 1.
<code>adFormat</code>	A numerical value assigned to this reader in the [Formats] section of the <code>formats.ini</code> file.
<code>descStr</code>	A string describing the file format.
<code>mMnmemStr</code>	The initial part of the custom reader file name with the "sr" excluded.

- If the return value is `TRUE`, the custom reader is used to parse the file and generate the token stream.
- If the return value is `FALSE`, all other readers in the `[CustomFilters]` section of the `formats.ini` file are tried. If no match is found, the file detection process continues checking for the formats supported by Filter SDK.
- The entry in the `[Formats]` section of the `formats.ini` file should be of the form `aaa.bbb.ccc.ddd`, where `aaa` is the value used for the `adFormat` parameter, `bbb` is the value of the file class, `ccc` is the value of the minor format, and `ddd` is the value of the major version.

xxxAllocateContext()

This function allocates a global memory block for a data context. A handle to this memory is returned to the structured access layer. The structured access layer passes this handle back to all reader entry points.

Syntax

```
void * pascal _export xxxAllocateContext(  
    void                *pSALContext,  
    LPARAM (pascal *fp)(void *,  
    UINT                LPARAM),  
    Bool                *pbOpenDoc,  
    TPVAPIServices      *pVapi,  
    DWORD               dwFlags)
```

Arguments

<code>pSALContext</code>	A pointer to the global data context structure of the structured access layer.
<code>fp</code>	A pointer to a structure of callback functions supported by the structured access layer.
<code>pbOpenDoc</code>	You must set this <code>BOOL</code> value to <code>TRUE</code> if the allocation of memory for the global data context structure is successful.
<code>pVapi</code>	A pointer to a structure providing memory management and character conversion functions. Because this functionality is proprietary to HPE, <code>TPVAPIServices</code> is redefined as <code>void</code> in <code>kvcfsr.h</code> .
<code>dwFlags</code>	Run-time flags controlled by the structured access layer.

Returns

- Upon success, a pointer to the global data context structure for the custom reader. This pointer is passed back to all other custom reader entry points.
- Upon error, a `NULL` pointer. This causes the structured access layer to shut down the process.

Discussion

The global context structure should be `memset()` to zero in this function.

xxxFreeContext()

This function terminates an instance of the custom reader.

Syntax

```
int pascal _export xxxFreeContext(void *pCfContext)
```

Arguments

`pCfContext` A pointer to the global context structure for the custom reader.

Returns

- Upon success, `KVERR_Success`.
- Upon error, a non-zero error code.

Discussion

All memory that still remains allocated within the custom reader must be freed within this function.

xxxInitDoc()

This function initializes non-zero, non-null members of `pContext`.

Syntax

```
int pascal _export xxxInitDoc(  
    void          *pCfContext,  
    adDocDesc     *pAutoInfo,  
    long          lcbFileSize,  
    KPTPIOobj     *pIO )
```

Arguments

`pCfContext` A pointer to the global context structure for the custom reader.

`pAutoInfo` A pointer to an `adDocDesc` structure defined in `kwautdef`.

`lcbFileSize` The length of the source file in bytes.

pIo A pointer to a KPTPIOobj structure defined in kvioobj.h.

Returns

- Upon success, KVERR_Success.
- Upon error, a non-zero error code. This causes the structured access layer to shut down the process.

Discussion

- For custom readers, the pAutoInfo variable can be ignored.
- If the structured access layer has determined the length of the source file, that value is provided by the lcbFileSize parameter. If it is zero, the file size must be determined in this function.
- The pointer pIo provides access to file management functions defined in kvioobj.h.
- In this function, all non-zero, non-NULL members of the global context structure should be initialized.

xxxFillBufferO

This function controls parsing of the source file and generation of tokens defined in kvtoken.h.

Syntax

```
int  pascal _export xxxFillBuffer(  
    void      *pCFContext,  
    BYTE      *pcBuf,  
    UINT      *pnBufOut,  
    int       *pnPercentDone,  
    UINT      cbBufOutMax)
```

Arguments

pCFContext	A pointer to the global context structure for the custom reader.
pcBuf	A pointer to a memory buffer to which the tokens are written.
pnBufOut	A pointer to a variable that specifies the actual number of bytes written to the token buffer.
pnPercentDone	A pointer to a variable that specifies the percentage completed of the file parsing.
cbBufOutMax	A pointer to a variable that specifies the maximum number of bytes written to the token buffer.

Returns

- Upon success, `KVERR_Success`.
- Upon error, a non-zero error code. This causes the structured access layer to shut down the process.

Discussion

- Calls are made to read and parse the source file within this function.
- This function is called repeatedly by the structured access layer until either the return value is `FALSE` or the percentage complete is 100.
- The actual number of bytes written to the token buffer must not exceed the value of `cbBufOutMax`.

xxxGetSummaryInfo()

This function is required to extract document summary information.

Syntax

```
int  pascal _export xxxGetSummaryInfo(  
    void                *pCfContext,  
    KVSummaryInfoEx     *pInfo,  
    BOOL                bFreeInfo)
```

Arguments

`pCfContext` A pointer to the global context structure for the custom reader.

`pInfo` A pointer to a `KVSummaryInfoEx` structure defined in `kvtypes.h`.

`bFreeInfo` A `BOOL` value indicating whether to free memory allocated for summary information.

Returns

- Upon success, `KVERR_Success`.
- Upon error, a non-zero error code.

Discussion

This function uses an instance of the global context structure that is different from the one used by all other reader interface functions.

This function can call the same functions used by `xxxFillBuffer()` or can be completely independent.

For more information, see [Extract Metadata, on page 58](#).

xxxOpenStreamO

This function is required when initiating processing of peripheral elements such as document headers, footers, footnotes, and endnotes.

Syntax

```
int pascal _export xxxOpenStream(  
    void      *pCFContext,  
    int       type,  
    int       nOrdinal)
```

Arguments

- | | |
|-------------------------|---|
| <code>pCFContext</code> | A pointer to the global context structure for the custom reader. |
| <code>type</code> | An integer identifying a specific header, footer, footnote, or endnote. Options are defined in <code>kvcfsr.h</code> . |
| <code>nOrdinal</code> | An integer identifying a specific header, footer, footnote, or endnote. See the associated macros in <code>kvtoken.h</code> . |

Returns

- Upon success, `KVERR_Success`.
- Upon error, a non-zero error code.

Discussion

A call to this function results in a call to `xxxFillBuffer()`. The function `xxxFillBuffer()` provides a new empty output buffer and a new token stream input buffer to process the alternate stream for peripheral elements. In this alternate stream, paragraph and character style properties are likely different from the main body. Therefore, as the document is parsed, the existing values from the main body must be saved. When the processing of the alternate stream is completed and processing of the main body resumes, these values must be restored in `xxxCloseStream()`.

xxxCloseStreamO

This function is required when terminating processing for document headers, footers, footnotes, and endnotes.

Syntax

```
int pascal _export xxxCloseStream(  
    void      *pCFContext,
```



```
int type)
```

Arguments

pCfContext A pointer to the global context structure for the custom reader.

type An integer identifying a specific header, footer, footnote, or endnote. Options are defined in `kvcfsr.h`.

Returns

- Upon success, `KVERR_Success`.
- Upon error, a non-zero error code.

Discussion

Prior to exiting this function, the previously saved values in the global context structure must be restored. This ensures that processing of the main body resumes with the correct document state.

xxxCharSet()

This function identifies the character encoding used within the source document.

Syntax

```
KVCharSet pascal _export xxxCharSet(  
    void *pCfContext,  
    BOOL *bMSBLSB)
```

Arguments

pCfContext A pointer to the global context structure for the custom reader.

bMSBLSB The `BOOL` value required for Unicode text. Set this argument to `TRUE` for Big Endian and `FALSE` for Little Endian.

Returns

One of the enumerated values defined in the `KVCharSet` structure of `kvtypes.h`.

Discussion

If the custom reader can determine the character encoding of the document, the corresponding enumerated value is returned. If the character encoding cannot be determined, `KVCS_UNKNOWN` is returned.

Appendix H: Password Protected Files

This section lists supported password-protected container and non-container files and describes how to open them.

- [Supported Password Protected File Types](#)226
- [Open Password Protected Container Files](#)227
- [Filter Password Protected Files](#)227

Supported Password Protected File Types

The following table lists the password-protected file types that KeyView supports.

Key to support table

Symbol	Description
Y	Format is supported.
N	Format is not supported.
S	Support for viewing subfiles.
V	Support for viewing content.
P	Password required.
C	Password and certificate or User ID file required.

Supported password-protected file types

File Type	Version	Filter	Export	Extract	View	Credentials
PST (Windows)	n/a	N	N	Y	S	P
PST (non-Windows) ¹	n/a	N	N	Y	S	N
ZIP	n/a	N	N	Y	S	P
7-Zip	n/a	N	N	Y	S	P
RAR	n/a	N	N	Y	S	P
SMIME in MSG, EML, MBX	n/a	N	N	Y	N	C
Lotus Notes NSF	n/a	N	N	Y	N	C

¹The native PST reader, `pstnsr`, does not require credentials to open password-protected PST files that use compressible encryption.

Supported password-protected file types, continued

File Type	Version	Filter	Export	Extract	View	Credentials
Adobe PDF	n/a	Y	Y	Y	V	P
Microsoft Office	97-2003 2007 2010	Y	Y	Y	V	P

Open Password Protected Container Files

This section describes how to extract password-protected container files by using the Java API. The following guidelines apply to specific file types.

- **Lotus Notes NSF files.** If you are running a Notes client with an active user connected to a Domino server, you must specify the user's password as a credential regardless of whether the NSF files you are opening are protected. This enables KeyView to access the Notes client and the Lotus Notes API. If the Notes client is not running with an active user, KeyView does not require credentials to access the client.
- **PST files.** To open password-protected PST files that use high encryption (Microsoft Outlook 2003 only), you must use the MAPI-based PST reader (`pstsr`). The native PST reader (`pstnsr`) returns the error message `KVERR_PasswordProtected` if a PST is encrypted with high encryption.

To open container files

- Set the credential information to an `ExtOpenDocConfig` object, and pass it to the `extOpenDocument` method. For example:

```
odconfig = new ExtOpenDocConfig();
odconfig.setPassword(m_password);
extContextID = m_objFilter.extOpenDocument(inFile, odconfig);
```

Filter Password Protected Files

This section describes how to filter password-protected non-container files with the Java API.

To filter password-protected files

- Use the `setSourcePassword(java.lang.String pwd)` method. For example:

```
objFilter.setSourcePassword(pwd);
```

where `pwd` is a null-terminated string of 255 characters or fewer.

Send documentation feedback

If you have comments about this document, you can [contact the documentation team](#) by email. If an email client is configured on this system, click the link above and an email window opens with the following information in the subject line:

Feedback on Filter SDK Java Programming Guide (KeyView 11.5)

Add your feedback to the email and click **Send**.

If no email client is available, copy the information above to a new message in a web mail client, and send your feedback to autonomytpfeedback@hpe.com.

We appreciate your feedback!