Flib datasheet

PDFlib Products in the Real World

What is PDFlib?

PDFlib is the leading developer toolbox for generating and manipulating files in the Portable Document Format (PDF).

PDFlib's main targets are dynamic PDF creation on a Web server or any other server system, and to implement »Save as PDF« in existing applications. You can use PDFlib to dynamically create PDF documents from database contents, similar to dynamic Web pages

PDFlib has proven itself in a wide range of other use cases as well.

Application programmers need only decent graphics or print but put experience to be able to use PDFlib quickly. Since PDFlib frees you from the technicalities of the PDF file format, you can focus on acquiring the data and arranging text, graphics, and images on the page.

The PDFlib Product Family

The PDFlib product family comprises the following product

PDFlib offers all functions required to generate PDF documents with text, graphics, images, and interactive elements such as annotations or bookmarks.

PDFIIb+PDI includes all PDFIIb functions plus the PDF mport Library (PDI). With PDI you can open existing PDF documents and incorporate some pages into the PDFIIb output.

PDFlib Personalization Server (PPS) Includes PDFlib+PDI plus additional functions for variable data processing using PDFlib Blocks. PPS makes applications independent of any layout changes.

»Save as PDF« for Applications

I work with a software development company and want to implement a »Save as CDF« feature in our applications.

PDFlib easily integrates into all kinds of applications to enable reliable and high-quality PDF output. Many well-known developers of graphics programs, geographical information systems (GIS), prepress and DTP applications and from many other domains rely on PDFlib to add PDF generation features to their products.

Invoices for an Online Skop

How can I create PDE invoices dynamically in my online shop?

Dynamic Invoice generation is one of the most popular PDFlib conarios. The generated PDF invoices can be viewed in the Web browser, made available for separate download, or e-mailed to the user.

Use PDFIN to place transaction data (customer details, item list, prices, tc.) on a PDF page. Add images, such as a company logo, in a variety of image formats. Use PDFlib+PDI to incorporate existing PDF material, for example company stationery as background.

Majl Merge

How can I merge personal data into an existing PDF document to create mass mailings?

PDFlib+PDI imports one or more pages of an existing PDF and adds individual text and images to create individual letters. The programmer adds code for retrieving text or graphics from a text file or database. A single large PDF containing all letters can be produced for printing, or many personalized small PDFs for e-mailing to the recipients.

If you need more flexibility because slightly different mailings must be produced or changes in the page design occur frequently, you can use the PDFlib Personalization Server (PPS). This facilitates both the designer's and the programmer's job when it comes to variable data processing.

Invoices and Reports from Office Applications

I'm unsatisfied with the look of invoices and reports created by our office applications. How can I create nice PDF documents?

PDFlib can be attached to common office applications. You can add PDF capability to MS Office and other applications with the popular Visual Basic scripting language. Use PDFlib to create invoices from a database in order to print or e-mail them to customers. Use PDFlib+PDI to incorporate PDF company stationery. Make PDF processing even more efficient by deploying the PDFlib Personalization Server (PPS).

Commercial Printing

Can I use PDFlib to prepare prepress data for commercial printing?

Customers use PDFlib to build systems for creating, assembling, or personalizing PDF documents for commercial printing. In many cases these production systems are accessible via a Web browser. The PDFlib product family supports a variety of features for the graphic arts industry, including color management with ICC profiles, CMYK color, spot colors with built-in PANTONE® and HKS® tables, and PDF/X-3, PDF/X-4 or PDF/X-5 conforming output.

Mass Generation of Phone Bills

I am responsible for creating the monthly phone bills at a major telecommunications provider. We plan to migrate from paper-based bills to online PDFs and distribute them via e-mail or Web.

PDFlib has a proven track record in mission-critical environments. Even with several millions of bills in each run you won't experience performance or reliability problems. PDFlib works on any kind of server, including midrange and mainframe systems.

Spice up existing PDFs

Can I add stamps and page numbers to existing PDF pages?

First, you import the pages from your PDF with PDFlib+PDI. Then you add a transparent stamp, running page numbers, barcodes, company logo, or any other content. You can even add interactive elements including links, form fields, bookmarks, etc. With these features you can approach PDF problems with a postprocessing solution.

Other PDFlib GmbH Products

PDFlib TET. Extract text and images from any PDF and normalize it to Unicode. TET includes patented content analysis algorithms for identifying word boundaries or dehyphenating text, and much more.

PDFlib TET PDF IFilter. Extracts the text and metadata of PDF documents and makes it available to search and retrieval softwar on Windows.

PDFlib PLOP. Linearize, optimize, and protect PDF documents, or add XMP metadata.

PDFlib PLOP DS. Apply digital signatures to PDF occuments.

PDFlib pCOS. Query any kind of information from PD

Supported Development Environments

PDFlib is everywhere – it runs on practically all computing platforms. We offer 32- and 64-bit variants for all common flavors of Windows, Mac OS X, Linux and Unix as well as for IBM eServer iSeries and zSeries mainframes

The PDFlib core is written in tighty optimized C code for maximum performance and small overhead. Via a simple API (Application Programming Interface) the PDFlip functionality is accessible from a variety of development environments:

- ► COM for use with **XB**, ASP, Borland Delphi, etc.
- ► C and C++
- ► CoboNIBM eServer zServes)
- ► Java, including servlets and Java Application Server

- ► .NET for use with C#, VB.NET, ASP.NET, etc.
- PHP hypertext processor
- ► Perl
- ► Python
- REALbasic
 REC (IRM of on vor iSc
- RPG (IBM eServer iSeries)
 Rubv
- ► Tcl

Benefits of using PDFIth Software

Rock-solid Products

Tens of thousands of programmers workwide are working with our software. PDFlib meets all quality and performance requirements for server deployment. All PDFlib products are suitable for robust 24x7 server deployment and nattended batch processing.

Speed and Simplicity

PDFlib products are increably fast – up to thousands of pages per second. The programming interface is straightforward and easy to learn.

PDFlib Products all over the World

Our products support all international languages as well as Unicode. They are used by customers in all parts of the world.

Professional Support

If there is a problem, we will try to help. We offer commercial support to meet the requirements of your business-critical applications. By adding support you will have access to the latest versions, and have guaranteed response times should any problems arise.

Licensing

We offer various licensing programs for server licenses, integration and site licenses, and source code licenses. Support contracts for extended technical support with short response times and free updates are also available.

About PDFlib GmbH

PDFlib GmbH is completely focused on PDF technology. Customers worldwide use PDFlib products since 1997. The company closely follows development and market trends, such as ISO standards for PDF. PDFlib GmbH products are distributed all over the world with major markets in North America, Europe, and Japan.

Contact

Fully functional evaluation versions including documentation and samples are available on our Web site. For more information please contact:

PDFlib

PDFlib GmbH





What is PDFlib?

PDFlib is the leading developer toolbox for generating and manipulating files in the Portable Document Format (PDF). PDFlib's mair targets are dynamic PDF creation on a Web server or any other server system, and to implement »Save as PDF« in existing applic tions. You can use PDFlib to dynamically create PDF documents from database contents, similar to dynamic Web pages. PDFlib has proven itself in a wide range of other use cases as well. Application programmers need only decent graphics or print output experience to be able to use PDFlib quickly. Since PDFlib frees you from the technicalities of the PDF file format, you can focus on acquiri ng the data and arranging text, graphics, and images on the page.

The PDFlib product family is available in three different flavors PDFlib, PDFlib+PDI (PDF Import), and the PDFlib Personalization Server (PPS).

PDFlib

PDFlib offers all functions required to generate PDF documents with text, graphics, images, and interactive dements such notations or bookmarks. Use PDFlib for the following tasks

- Add »Save as PDF« capability to your application
- Create PDF documents on a Web server in rest time
- ► Create database reports in PDF
- ► Create PDF/X-1/3/4/5 documents for commercial printing
- ► Convert TIFF, JPEG, or other image formats to PDF
- Create PDF/A for archiving

PDFlib+PDI (PDF Import)

PDFlib+PDI includes all PDFlib functions, plus the PDF Import Library (PDI). With PDI you can open existing PDF documents and incorporate some pages nto the PDFH output. Use PDFlib+PDI for all PDFlib tasks plus the following:

- ► Impose multiple PDF pages on a single sheet for printing
- ► Add text, such as beaders, footers, stamps, or page numbers to existing PDF pages
- Place mages, e.g. company logo, on existing pages
- Add baccodes to existing PDF pages
 Assemble existing PDF pages
- dd content to POF/X or PDF/A documents

PDFIb Personalization Server (PPS)

The PDFlib Personalization Server (PPS) includes PDFlib+PDI plus dditional functions for variable data processing using PDFlib Blocks PPS makes applications independent from layout changes. The designer creates the page layout and converts it to PDF. She takes into account areas as placeholders for variable text and nages. In Acrobat she drags a rectangular Block for each area using the PDFbb Block Plugin. Each Block contains a variety of Block properties such as font size, color, image scaling. The PDFlib Block Plugin offers a Preview feature which shows the results of filling locks according to their properties.

he programmer writes code to fill PDFlib Blocks with text, images, or PDF pages. He doesn't need to know the formatting or position of a Block. Use PPS for all PDFlib+PDI tasks plus the following:

- Customize direct mailings with text and images
- Fill templates for transactional and statement processing
- ► Personalize promotional material with address data
- Generate individual parts catalogs from a database
- Produce customized documentation for multiple similar products

What's new in PDFlib 8?

New PDF Features for Acrobat 9

PDFlib supports various PDF features according to Acrobat 9 (technically: PDF 1.7 Adobe extension level 3):

- External graphical content (Reference XObjects)
- ► Layer variants (also called layer configurations)
- ► PDF Portfolios
- ► Georeferenced PDF
- ► AES-256 encryption and Unicode passwords
- ▶ PDFlib+PDI and PPS can import and process Acrobat 9 data.

Font Handling and Text Output

Quite a number of new typographical features can be found in PDFlib 8:

- Complex script shaping and bidirectional formatting for Arabic, Thai, Hindi, and many other writing systems
- ► Fallback fonts
- ► OpenType layout features, e.g. ligatures and swash characters
- ► Retain fonts across documents
- ► SING fonts for CJK Gaiji characters
- ► Redesigned font engine
- ► Wrap text around image clipping paths
- ► Text on a path

PDFlib Block Plugin and the PDFlib Personalization Server

The PDFlib Block Plugin is used to prepare PDF documents for Block filling (personalization) with the PDFlib Personalization Server (PPS). New features:

- ► Preview PPS Block processing in Acrobat
- ► Redesigned user interface
- ► Snap-to-grid for quickly layout out Blocks in a raster
- Additional Block properties, e.g. for transparency
- ► Clone PDF/A or PDF/X states of the Block container
- Leverage PDFlib 8 features with Blocks

Other important features

There are a number of other important new features, details can be found in the product documentation:

- ► Reusable path objects
- ► PDF/X-4 and PDF/X-5

ions

- Alpha channel in TIFF and PNG images
- ► JBIG2-compressed images
 - Compressed object streams and cross-reference streams Built in PANTONE® Goe™ color libraries

PEFlib 8 also introduces a variety of improvements in existing func-

Common Features in PDFlib, PDFlib+PDI, and the PDFlib Personalization Server

PDFoutput	Generate PDE documents on disk file or directly in memory (for Web servers)
- Di output	High-volume output and arbitrary PDE file size (even beyond to GB)
	Suspend/resume and insert page features to create pages out of order
PDE flavors	DDE 1.2. DDE 1.70yt2* (Acrobat 4. a) including ISO 22000 1 (-DDE 1.7)
	Linearized (web entimized) DDE for but coming over the Web
	Tagged DDE for accessibility and reflew
	Marked Content for adding application specific data or alternational to the without Tagging*
ISO standards	Marked Content for adding application-specific data of alternate text without agging
ISO standards	ISO 15930: PDF/A for the graphic arts industry
Constant and	ISO 32000: standardized version of PDF 1.7
Graphics	Common vector graphics primitives: lines, curves, arcs empses, rectingles, etc.
	Smooth shadings (color blends), pattern fus and strokes
	Iransparency (opacity) and blend modes
	External graphical content (Reference (Objects) for variable data printing*
	Reusable path objects and clipping paths imported from mages*
Layers	Optional page content which car selectively be displayed
	Annotations and form fields can be placed on layers
	Layers can be locked, automatically activated depending on zoom factor, etc.
	Layer variants* (production-safe groups of layers) for PDF/X-4 and PDF/X-5
Fonts	TrueType (TTF and TTC) and PostScript Type 1 fonts (PFB and PFA, plus LWFN on the Mac)
	OpenType fonts with PostScript or TrueType outlines (TTF, OTF)
	Support for dozens of Oper Type layout features for Western and CJK text output, e.g. ligatures, small caps, old-style numerals, swach characters, simplified/traditional forms, vertical alternates*
	Directly use fonts which are installed on the Windows or Mac system (»host fonts«)
	Font embedding for all font types; subsetting for TrueType, OpenType, and Type 3 fonts
	User-defined (Type 3) fonts for bitmap fonts or custom logos
	EUDC and SING* for ts (graphlets) for CJK Gaiji characters
	Fallback fonts (pull missing glyphs from an auxiliary font)*
	Retain fonts across documents to increase performance*
Text output	Text output in different fonts; underlined, overlined, and strikeout text
· ·	Glyphan a font car be addressed by numerical value, Unicode value, or glyph name*
	Kerning for improved character spacing
	Artificial old, talic, and shadow* text
	Create text on a path*
_	Proportional widths for standard CJK fonts*
	Configurate replacement of missing glyphs
Internationalization	Unicode styings for page content, interactive elements, and file names*; UTF-8, UTF-16, and UTF-32
$\langle \rangle$	formats
\sim	pport for a variety of 8-bit and legacy multi-byte CJK encodings (e.g. Shift-JIS; Big5)
	Feton code pages from the system (Windows, IBM eServer iSeries and zSeries)
	standard and custom CJK fonts and CMaps for Chinese, Japanese, and Korean text
$\langle \gamma \rangle$	Vertical writing mode for Chinese, Japanese, and Korean text
\wedge	Character shaping for complex scripts, e.g. Arabic, Thai, Devanagari*
	Bidirectional text formatting for right-to-left scripts, e.g. Arabic and Hebrew*
\sim	Embed Unicode information in PDF for proper text extraction in Acrobat
$\langle \rangle$	
\sim	
$\langle \neg \rangle$	

I

I	4
L	

l

4	PDFlib, PDFlib+PDI, PPS, 2009-12 PDFlib GmbH www.pdflib.com
	(\frown)
Imagos	
inages	Automatic dataction of image file formate
	Automatic detection of image file formats
	Query image information (pixel size, resolution, ICC profile, clipping path, etc.)
	Interpret clipping paths in TIFF and JPEG images
	Interpret alpha channel (transparency) in TIFF and PNG images*
	Image masks (transparent images with a color applied), colorize images with a spot color
Color	Grayscale, RGB (numerical, hexadecimal strings, HTML color names), CMVK, QE Lab color
	Integrated support for PANTONE [®] colors (incl. PANTONE [®] Goe [™])* and HKS [®] colors
	User-defined spot colors
Color management	ICC-based color with ICC profiles; support for ICC 4 profiles*
	Rendering intent for text, graphics, and raster images
	Default gray, RGB, and CMYK color spaces to remap device-dependent colors
	ICC profiles as output intent for PDF/A and PDF/X
Archiving	PDF/A-1a and PDF/A-1b (ISO 19005-1)
	XMP extension schemas for PDF/A-1
Graphic arts	PDF/X-1a, PDF/X-3, PDF/X-4*, PDF/X-4p*, PDF/X-55*, PDF/X-5pg* (SQ 15930)
	Embedded or externally referenced* output intent IC profile
	External graphical content (referenced pages) for PDE/X-5p and PDF/X-5pg*
	Create OPI 1.3 and OPI 2.0 information for imported images
	Separation information (Plate Color)
	Settings for text knockout overning etc
Textflow Formatting	Format text into one or more restangular or arbitrativ shaped areas with hyphenation (user supplied
lextnow ronnatting	hyphenation points required front and color changes justification methods tabs leaders control com-
	mands; wrap text around images
	Advanced line-breaking with language-specific processing
	Flexible image placement and formatting
	Wrap text around images or image clipping paths*
Table formatting	Table formatter places rows and columns, and automatically calculates their sizes according to a vari-
0	ety of user preferences. Tables can be split across multiple pages.
	Table cells can hold single or multi-line text, images, PDF pages, path objects, annotations, and form
	fields
	Table cells can be formatted with ruling and shading options
	Flexible stamping function
	Matchoox concept for referencing the coordinates of placed images or other objects
Security	Encrypt PD output with RC4 (40/128 bit) or AES encryption algorithms (128/256* bit)
	Unicode passwords*
	specify permission settings (e.g. printing or copying not allowed)
Interactive elements	Create form fields with all field options and JavaScript
	Create actions for bookmarks, annotations, page open/close and other events
\wedge	Cheate bookmarks with a variety of options and controls
\langle	Page transition effects such as shades and mosaic
	Reate all PDE annotation types such as PDE links Jaunch links (other document types) Web links
	Named destinations for links bookmarks and document open action
\mathcal{I}	Veste page labels (symbolic pames for pages)
Multimedia	Embed 2D animations in PDE
	Croate DDE with geographial reference information*
	Create PDF with geospatial reference information Create Tagged PDF and structure information for accessibility or any finance information
	Create Tagged PDF and structure information for accessibility, page reflow, and improved content renurnosing links and other apportations can be integrated in the document structure.
	וכירי איז איז איז איז איז איז איז איז איז אי
\frown	
\sim	

Metadata	Document information: common fields (Title, Subject, Author, Keywords) and user-defined fields
	Create XMP metadata from document info fields or from client-supplied XMP stream
	Process XMP image metadata in TIFF, JPEG, and JPEG 2000 images* $\langle \cdot \rangle$
Programming	Language bindings for Cobol, COM, C, C++*, Java, .NET, Perl, PHP, Python, REALbasic, RPG, Ruby, Vcl
	Virtual file system for supplying data in memory, e.g., images from a database
	$\langle \frown \rangle$

* New or considerably improved in PDFlib/PDFlib+PDI/PPS 8

Additional Features in PDFlib+PDI and the PDFlib Personalization Server

	\sim / \sim
PDF input (PDI)	Import pages from existing PDF documents
	Import all PDF versions up to PDF 1.7 extension level 3 (Acrobat 9)
	Import documents which are encrypted with any of PDF's standard encryption algorithms (master password required)*
	Query information about imported pages*
	Clone page geometry of imported pages (e.g. BleedBox Tri nBox, CropBox)*
	Delete redundant objects (e.g. identical fonts) across multiple imported PDF documents
	Repair malformed input PDF ocuments*
	Copy PDF/A or PDF/X or tput intent from imported PDF documents
pCOS interface	pCOS interface for querying details about imported PDF documents*
	\rangle \rangle \vee

()

*New or considerably improved in PDFlib+PDI and PPS 8

Additional Features in the PDFlib Personalization Server

Variable Data Processing (PPS)	PDF personalization with PDFite Blocks for text, image, and PDF data
PDFlib Block Plugin	PDFlib Block plugin for reating PDFlib Blocks interactively in Acrobat on Windows and Mac
	Redesigned user interface*
	Preview PPS Block filling in Acrobat*
~	Snap-to grid for interactively creating or editing Blocks in Acrobat*
4	Clone DFX or PDF/A properties of the Block container*
\land	Convert PDF form fields to PDFlib Blocks for automated filling
$\langle \rangle$	Textflow Blocks can be linked so that one Block holds the overflow text of a previous Block
	List of PANTONE® and HKS® spot color names integrated in the Block plugin*
*New or considerably improved in	PPS 8
	>

Supported Development Environments

PDFlib is everywhere – it runs on practically all computing platforms. We offer 32- and 64-bit variants for all common flavors of Windows, Mac OS X, Linux and Unix, as well as for IBM eServer iSeries and zSeries mainframes.

The PDFlib core is written in highly optimized C code for maximum performance and small overhead. Via a simple API (Application Programming Interface) the PDFlib functionality is accessible from a variety of development environments:

- ► COM for use with VB, ASP, Borland Delphi, etc.
- ► C and C++
- ► Cobol (IBM eServer zSeries)
- ► Java, including servlets and Java Application Server
- ► .NET for use with C#, VB.NET, ASP.NET, etc.
- ► PHP
- ► Perl
- Python
- ► REALbasic
- RPG (IBM eServer iSeries)
- ► Ruby
- ► Tcl

Benefits of using PDFlib Software

Rock-solid Products

Tens of thousands of programmers worldwide are working with our software. PDFlib products meet all quality and performance requirements for server deployment. All products are suitable for robust 24x7 server deployment and unactended batch processing.

Speed and Simplicity

PDFlib products are incredibly fast – up to thousands of pages per second. The programming interface is straightforward and easy to learn.

PDFlib Products all over the World

Our products support all international languages as well as Unicode. They are used by custome is in all parts of the world.

Professional Support

icensing

If there so problem, we will try to help. We offer commercial support to meet the requirements of your business-critical applications by adding support you will have access to the latest versions, and have guaranteed response times should any problems arise.

We offer various licensing programs for server licenses, integration and site licenses, and source code licenses. Support contracts for extended technical support with short response times and free brodates are also available.

About PDFlib GmbH

PDFlib GmbH is completely focused on PDF technology. Customers worldwide use PDFlib products since 1997. The company closely follows development and market trends, such as ISO standards for PDF. PDFlib GmbH products are distributed all over the world with major markets in North America, Europe, and Japan.

Contact

Fully functional evaluation versions including documentation and samples are available on our Web site. For more information please contact:



PDFlib GmbH

PDFIib

datasheet **PDFlib TET 4** Text Extraction Tool

What is PDFlib TET?

TET

PDFlib TET (Text Extraction Toolkit) reliably extracts text, images and metadata from PDF documents. TET makes available the text contents of a PDF as Unicode strings, plus detailed glyph and fort information as well as the position on the page. Raster images are extracted in common raster formats. TET optionally converts PDF documents to an XML-based format called TETML which contains text and metadata as well as resource information.

TET contains advanced content analysis algorithms for determining word boundaries, grouping text into columns and removing redundant text. Using the integrated pCOS interface you can retrieve arbitrary objects from the PDF, such as metadata, interactive elements, etc.

With PDFlib TET you can:

- ► Implement the PDF indexer for a search engine
- Repurpose the text and images in PDFs
- ► Convert the contents of PDFs to other formats
- Process PDFs based on their contents, e.g. splitting based on headings (requires PDFlib+PDI in addition to TET)

PDFlib TET Features

Accepted PDF Input

- TET supports all relevant flavers of PDF input
- ► All PDF versions up to Actobal 9, including 190 32000-1
- Protected PDFs which do not require a password for opening the document
- Damaged PDF documents will be repaired

Unicode

Since text in PDF is usually not encoded in Unicode, PDFlib TET normalizes the text in a PDF document to Unicode:

- TET converts all text contents to Unicode. In C and other non-Unicode aware languages the text is returned in the UTF-8 or UTF-16 formats, and as native strings in Unicode-capable programming languages
- Ligatures and other multi-character glyphs are decomposed into a sequence of the corresponding Unicode characters.

- Gipples without appropriate Unicode mappings are identified as cuch, and are mapped to a configurable replacement character in order to avoid misin erpretation.
- NTT implements various workarounds for problems with specific occurrent creation packages, such as InDesign and TeX documents or PDEs generated on mainframe systems.

Content Analysis and Word Detection

- TET includes advanced content analysis algorithms:
- Patented algorithm for determining word boundaries which is required to retrieve proper words
- Recombine the parts of hyphenated words (dehyphenation)
- Remove duplicate instances of text, e.g. shadow and artificially bolded text
- Recombine paragraphs in reading order
- ► Correctly order text which is scattered over the page

Page Layout and Table Detection

The page content is analyzed to determine text columns. Tables are detected, including cells which span multiple columns. This improves the ordering of the extracted text. Table rows and the contents of each table cell can be identified.

Geometry

TET provides precise metrics for the text, such as the position on the page, glyph widths, and text direction. Specific areas on the page can be excluded or included in the text extraction, e.g. to ignore headers and footers or margins.

Image Extract

Images on PDF pages can be extracted as TIFF, JPEG, or JPEG 2000 files. Precise geometric information (position, size, and angles) are reported for each image. Fragmented images will be combined to larger images to facilitate repurposing. Image fidelity is guaranteed since no downsampling or color space conversion occurs. This ensures the highest possible image quality.

PDF Analysis

The TET library includes the pCOS interface for querying details about a PDF document, such as document info and XMP metadata, font lists, page size, and many more (see separate datasheet for the pCOS product).

Configuration Options for problematic PDF

TET contains special handling and workarounds for various kinds of PDF where the text cannot be extracted correctly with other products. In addition, it includes various configuration features to improve processing of problem documents:

- Unicode mapping can be customized via user-supplied tables for mapping character codes or glyph names to Unicode.
- PDFlib FontReporter is an auxiliary tool for analyzing fonts, encodings, and glyphs in PDF. It works as a plugin for Adobe Acrobat. This plugin is freely available for Mac and Windows.
- Embedded fonts are analyzed to find additional hints which are useful for Unicode mapping. External font files or system fonts are used to improve text extraction results if a font is not embedded.

Unicode Postprocessing

TET supports various Unicode postprocessing steps which can be used to improve the extracted text:

- ► Foldings preserve, remove or replace characters, e.g. remove punctuation or characters from irrelevant scripts.
- Decompositions replace a character with an equivalent sequence of one or more other characters, e.g. replace narrow, wide or vertical Japanese characters or Latin superscript (e.g. ^a) variants with their respective standard counterparts.
- Text can be converted to all four Unicode normalization forms e.g. emit NFC form to meet the requirements for Web text or a database.

Document Domains

PDF documents may contain text in other places than the page contents. While most applications will deal with the page contents only, in many situations other document domains may be relevant as well. TET extracts the text from all of the following document domains:

- ► page contents
- predefined and custom document info entries
- ► XMP metadata on document and image level
- bookmarks
- ► file attachments and PDF portfolios can be processed recursively
- ► form fields
- comments (annotations)
- general PDF properties can be queried, such as page count, conformance to standards like PDF/A or PDF/X etc.

XMP Metadata

TET supports XMP metadata in several ways

- Using the integrated pCOS noterface, XMP metadata for the document, individual pages, intages, or other parts of the document can be extracted programmatically.
- TETML output contains XMP occument and image metadata if present in the PDF.
- Images extracted in the TIPE or JPEC formats contain image metadata if present in the PDE

TETML represents PDF Contents as XML

TET optionally represents the PDF contents in an AML flavor called TETML. It contains a variety of PDF information in a form which can easily be processed with common XML tools. TETML contains the actual text plus optionally font and position information, resource details (fonts, images, colorspaces), and metadata

TETML is governed by a corresponding KML schema to make sure that TET always creates consistent and reliable XML output. TETML can be processed with XSLT stylesheets, e.g. to apply certain filters or to convert TETML to other formats. Sample XSLT stylesheets for processing TETML are included in the TET distribution.

The following fragment shows TETML output with glyph details:

<Word> <Text>PDFlib</Text>

<Box llx="111.48" 11y="636-33" urx="161)14" ury="654.33">

y="636.33" width="9.65">P</Glyph> y="636.33" width="11.88">D</Glyph> <Glyph 'F1' 18' 11.48 size X 1.12" <Glyp font 18 "13.00" y="636.33" width="1.833">F</Glyph> "141.33" y="636.33" width="4.88">l</Glyph> "146.21" y="636.33" width="4.88">l</Glyph> "151.08" y="636.33" width="10.06">b</Glyph> <Glyph ont="F1 e="18' F1" <Glyph font siz "18' х **ɗ**lyph font="F size= < lyph 'F1' "18 </Box> /Word>

TET connectors provide the necessary glue code to interface TET with other software. The following TET connectors make PDF text extraction functionality available for various software environments:

- TET onnector for the Lucene Search Engine
- TET connector for the Solr Search Server
- TET connector for Oracle Text
- TNJ connector for MediaWiki

THT PDF IFilter for Microsoft products is available as a separate product. It extracts text and metadata from PDF documents and makes it available to search and retrieval software on Windows (see separate datasheet for details).

TET Cookbook

TET Connectors

The TET Cookbook is a collection of programming examples which demonstrate the use of TET for various text and image extraction tasks. Several Cookbook samples show how to combine the TET and PDFlib+PDI products in order to process and enhance PDF documents, e.g. add bookmarks or links based on the text on the page.



trategische Grundsätze – der	
r der Nutzung von Synergie-	
in Branchen sowie in Unter-	
lukterstellung. So verringert	
t bei der Produkterstellung –	
g – seit längerem nicht nur	

TET correctly removes the hyphen, but keeps the dash.

Introduction

Other products extract »Inttrroduccttiion«. TET correctly extracts »Introduction«.

> Canadian Institute for Theoretical Astrop Observatoire de Paris, LERMA, 61 avenu Observatoire Midi-Pyrénées, UMR 5572, Oepartment of Astronomy, University of Observatorio Astronomico di Bologna, vi

Other products extract »Midi-Pyr´en´ees«. TET correctly extracts »Midi-Pyrénées«.

> is permanently hidden from Earth. The first photographs of the hid

cial satellite; modern satellites prov

Other products extract » e rst photographs. TET correctly extracts »The first photographs«.

tellen Sie sich vor, Sie stehen an einem Kinder ins Wasser springen und schwim vor, Sie graben am Sandstrand zwei klei Schritte landeinwärts, jeder eine Hand breit, so Kanäle fließen kant. Stellen Sie sich jetzt noc mittels einer Streichholzes und kleiner weißer

Other products extract two words: the drop cap »S« and »tellen«. TET correctly extracts the single word »Stellen«.

Challenges with PDF Text Extraction

Dehyphenation

TET detects hyphenated words which span inultiple lines, removes the hyphen, and combines the individual parts to form a complete word. This is important to make sure that searches for the full word will be successful although only hyphenated parts are present in the document. Dashes (different from hyphens) will be treated separately since they must not be removed.

Shadow and artifical bold Text Detection

Digital documents often contain shadowed text where the shadow effect is achieved by placing the text multiply on the page, using a small offset between the instances of text. Similarly, bold text is often simulated by overprinting the same text multiply. As a result, the document contains the characters in the shadowed or bold word more than once. TET patented shadow detection algorithm identifies and removes redundant instances of text to avoid excess text extraction. While other software extracts the shadowed or bold text multiply, TEF correctly removes the redundant copies. While extra instances of a word will still result in a search engine bit, no more hits would be found if the text is duplicated character by character as in the example.

ccented Characters

In many languages accents and other diacritical marks are placed close to other characters to form combined characters. Some typesetting programs, most notably TeX, emit two characters (base character and accent) separately to create a combined character. For example, to create the character *ä* first the letter *a* is placed on the page, and then the dieresis character [¨] is placed on top of it. IET detects this situation and recombines both characters to form the appropriate combined character.

Ligatures

Ligatures combine two or more characters in a single glyph. The most common ligatures are in use for the combinations *fi*, *fl*, and *ffi*; less common ligatures are used for the combinations *Th*, *sp*, *ct*, *st*, and many others. When extracting text from digital documents, ligatures must be analyzed and separated to the constituent characters to allow proper text processing. TET detects ligatures and delivers two or more characters as appropriate.

Drop Caps

Drop caps are large initial characters at the beginning of a paragraph where the top of the initial aligns with the top of the line, and the remainder of the character drops down several lines. Drop caps are used to emphasize the start of a paragraph. If they are not treated properly the initial word is extracted in two parts: the single initial character and the remainder of the word.

3



Other products extract unusable garbage, while TET delivers text.



The page contents are not even displayed in Acrobat, but TET still correctly extracts the text.



TET reorders the visual mixture of right-to-left and left-to-right text to create proper logical text output.

Other products extract 133 tive little strips. TET extracts a single large image

Challenges with PDF Text Extraction

Unicode Mapping

Unicode mapping forms the foundation of PDF text extraction: every glyph on the page must be assigned the corresponding Unicode value. PDF complicates this tasks by supporting a variety of font and encoding variants which may or may not provide the information required to assign proper Unicode values. In the worst case the document does not provide enough information with the result that no usable text can be extracted from the document.

TET's patented Unicode mapping algorithm implements a cascaded algorithm which takes all available pieces of information in order to determine Unicode values. For many problematic documents TET extracts proper Unicode text where other products deliver only unusable garbage.

Damaged PDF Documents

ment.

PDF documents may get tamaged because of transmission errors or other problems. FET's repair mode recovers many kinds of damared PDFs. Sometimes PDF documents are damaged so heavily that the pages cannot ever be displayed in Acrobat. Even in such extreme cases TEN often delivers the page contents of the docu-

Bidirectional Text with Arabic and Hebrew

PIPE does not encode logical text, but is simply a container for glyphs on the page. Text in the Arabic and Hebrew script runs from right to left. Since it often contains left-to-right inserts such as numbers of names in Western languages, text must be interpreted in both directions – hence the term »bidirectional«. Arabic poses additional challenges since the characters can be used in up to four different contextual forms. These shaped forms of characters must be normalized to the corresponding standard (isolated) form.

Challenges with PDF Image Extraction

Color Spaces and Compression

Raster image data in PDF may be encoded in any combination of eleven color spaces and nine compression filters, but common image file formats such as JPEG and TIFF support only a subset of those. TET's image extractor carefully balances the characteristics of the PDF image with the capabilities of the image output format. Regardless of the internal structure of the PDF image, the pixel image will be extracted in one of the common image file formats.

Image Merging

The images in many PDF documents are broken into smaller pieces by the software producing the PDF. What appears as a single image on the page may actually consist of hundreds or thousands of small fragments. Among others, Microsoft Office applications and TeX are known to produce such documents. TET detects fragmented images and merges the pieces to form a usable larger image. Only with image merging such images can be repurposed in any way.

Many Ways to use TET

TET is available as a programming library for various development environments, and as a command-line tool for batch operations. Both offer similar features, but are suitable for different deployment scenarios. Both the TET library and the TET command-line tool can create TETML, TET's XML-based output format.

TET offers the following deployment options:

- The TET programming library (component) is used for integration into desktop or server applications. Examples for using the library are included in the TET package.
- The TET command-line tool is suited for batch processing PDF documents. It doesn't require any programming, but offers command-line options which can be used to integrate it into complex workflows.
- TETML output is suited for XML-based workflows and developers who are familiar with the wide range of XML processing tools and languages, e.g. XSLT.
- ► TET connectors are suited for integrating TET in various common software packages, e.g. databases and search engines.

The TET Family of Products

The TET family comprises the following products:

- ► The TET core product as described in this datasheet.
- TET PDF IFilter is available as a separate product. It is suitable for use with Microsoft search products, e.g. Windows Search, Share-Point and SQL Server (see separate datasheet for details).
- The TET Plugin for Adobe Acrobat is a free utility for extracting text and images from PDF. It can be used to evaluate TET interactively.

Supported Development Environments

PDFlib TET is everywhere – it runs on practically all computing platforms. We offer 32-bit and 64-bit packages for all common flavors of Windows, Mac OS, Linux and Unix, as well as for IBM i5/iSeries and zSeries systems.

The TET core is written in highly optimized C code for maximum performance and small overhead. Via a simple API (Application Programming Interface) the TET functionality is accessible from a variety of development environments:

- ► COM for use with VB, ASP, Borland Delphi, et
- ► C and C++
- ► Java, including servlets and Java Application Serv
- ► .NET for use with C#, VB.NET, ASP.NET, etc.
- ► Perl
- ► PHP
- ► Python
- ► REALbasic
- ► RPG (IBM i5/iSeries)

Benefits of using PDFlib Software

Rock-solid Products

Tens of thousands of programmers worldwide are working with our software. PDFlib meets all quality and performance requirements for server deployment. All PDFlib product, are suitable for robust 24x7 server deployment and unactended batch processing.

Speed and Simplicity

PDFlib products are incredibly fast – up to thousands of pages per second. The programming interface is straightforward and easy to learn.

PDFlib Products all over the World

Our products support all international languages as well as Unicode. They are used by custome is in all parts of the world.

Professional Support

licensing

If there so problem, we will try to help. We offer commercial support to meet the requirements of your business-critical applications by adding support you will have access to the latest versions, and have guaranteed response times should any problems arise.

We offer various licensing programs for server licenses, integration and site licenses, and source code licenses. Support contracts for extended technical support with short response times and free updates are also available.

About PDFlib GmbH

PDFlib GmbH is completely focused on PDF technology. Customers worldwide use PDFlib products since 1997. The company closely follows development and market trends, such as ISO standards for PDF. PDFlib GmbH products are distributed all over the world with major markets in North America, Europe, and Japan.

Contact

Fully functional evaluation versions including documentation and samples are available on our Web site. For more information please contact:



PDFlib GmbH



datasheet PDFlib PLOP 4.1

Linearization, Optimization, Protection

What is PDFlib PLOP?

PDFlib PLOP is a versatile tool for linearizing, optimizing, repairing, analyzing, encrypting and decrypting PDF documents. PLOP linearzation and optimization features create efficient and small PDF documents for fast Web delivery. PLOP protection features encrypt or decrypt PDF documents and apply or remove permission settings, such as »printing not allowed« or »content extraction not allowed«. PLOP's repair mode automatically detects damaged PDF documents and fixes the problems if possible. PLOP analysis features can be used to query arbitrary properties of a PDF document. Document info entries and XMP metadata can be retrieved and set in a PDF/A and PDF/X conforming manner.

PDFlib PLOP Features

Linearization

With PDFlib PLOP you can linearize a PDF document for fast delivery over the Web (byteserving). Byteserving increases the perceived download speed since the first page is already visible while the remainder of the document is downloaded in the background.

Optimization

PLOP can significantly reduce the file size of a POF ocument without affecting quality. It achieves this by removing unnecessary or redundant identical objects, such as epertedly embedded fonts, images, identical ICC color profiles, etc.

Protection

PLOP can apply user and master passwords, and set access permissions to prevent the document from being printed with Acrobat, disallow text extraction or modification, etc.

PLOP supports all relevant PDF encryption methods including strong AES-256 encryption and Unicode passwords for Acrobat X. With PLOP you can:

- encrypt a PD document with user or master password, or both;
- remove PDE encyption (if you know the master password);
- add or remove permission settings, e.g. »text extraction not allowed« if you know the master password);
- query information about the security status, encryption scheme, permission settings, and document info fields

Various kinds of damaged PDF documents are detected and automatically repaired, if possible.

DF Analysis

Repair Mode

The PLOP library includes the pCOS interface for querying details about a PDF document, such as document info and XMP metadata, font lists, page size, and many more (see separate datasheet for the pCOS product).

XMP Metadata

Metadata (»data about data«) is an important topic in many areas of opplication software. XMP (Extensible Metadata Platform) is an XML-based framework with many predefined metadata properties. XMP is integrated in Acrobat/PDF, and much more powerful than simple document info entries. XMP is required for PDF/A and other ISO standards. Many industry groups have published XMP-based recommendations for vertical applications.

With PLOP you can insert XMP metadata in PDF documents and extract XMP from PDF. Inserted XMP will be validated to make sure that valid output can be created. If the input document conforms to the PDF/A-1 standard, the user-supplied XMP must conform to the XMP rules set forth in PDF/A.

XMP insertion with PLOP can be used in the following and other situations (sample XMP is contained in the PLOP distribution):

- ► Add XMP metadata to PDF/A-1 documents, including support for XMP extension schemas as defined in the PDF/A-1 standard.
- ► Add XMP metadata describing the scanning process for digitized legacy documents, also according to PDF/A-1.
- Add XMP metadata according to the Ghent Workgroup (GWG) Ad Ticket scheme.
- ► Add company-specific XMP metadata.

Document Info Entries

With PLOP you can add new document information entries or replace the values of existing info entries. Both predefined and custom entries can be set. If the input document contains XMP document metadata, all predefined info entries will automatically be synchronized to the XMP metadata in order to keep the metadata consistent (this is a requirement of PDF/A-1).

PDF Versions and Standards

PLOP supports all PDF versions up to Acrobat X, including PDF 1.7 (ISO 32000). PLOP is PDF/A-aware: if the input document conforms to the PDF/A standard (ISO 19005-1), the output document is guaranteed to still comply with PDF/A. PLOP fully supports XMP extension schemas as required by PDF/A-1. Similarly, PLOP is aware of PDF/X-1a/3/4/5 (ISO 15930).

The ability to insert PDF/A-conforming XMP metadata in PDF documents is an important advantage of PLOP.

PLOP Library or Command-Line Tool?

PLOP is available as a programming library (component) for various development environments, and as a command-line tool for batch operations. The library and the command-line tool offer similar features, but are suitable for different deployment tasks.

The PLOP programming library is used...

...for integration into your desktop or server application. Examples for using the library with all supported language bindings are included in the PLOP package. Since the PLOP library accepts PDF input documents from a disk file or directly in memory, it can easily be combined with other products.

The PLOP command-line tool is suited...

...for batch processing PDF documents. It doesn't require any programming, but offers powerful command-line options which can be used to integrate it into complex workflows. The PLOP command-line tool can also be called from environments which not support the use of the PLOP library.

Supported Development Environments

PDFlib PLOP is everywhere – it runs on practically all computing platforms. We offer 32-bit and 64-bit packages for all common flavors of Windows, Mac OS, Linux and Unix, as well as for IBW is iSeries and zSeries systems. The PLOP core is written in highly opt mized C code for maximum performance and small overhead. Via a simple API (Application Programming Interface) the PLOP functionality is accessible from a variety of development environments.

- ► COM for use with VB, ASP, Borland Delphi, etc.
- ► C and C++
- ► Java, including servlets and Java Application Server
- ▶ .NET for use with C#, VB.NET, ASP.NET, etc.
- ► Perl
- ► PHP
- ► Python
- ► RPG on i5/iSeries

PLOP DS for digitally signing PDF

The extended version PLOP DS supports all features of PLOP, plus the ability to apply digital signatures to PDF documents. Please see the separate PLOP DS datasheet for more information.



Benefits of using PDFlib Software

Rock-solid Products

Tens of thousands of programmers worldwide successfully use our software. PDFlib products meet all quality and performance requirements for server deployment. All products are suitable for robust 24x7 server deployment and unattended batch processing.

Speed and Simplicity

PDFlib products are incredibly fast – up to thousands of pages per second. The programming interface is straightforward and easy to learn.

PDFlib Products all over the World

Our products support all international languages as well as Unicode. They are used by custome is in all parts of the world.

Professional Support

Licensing

If there so problem, we will try to help. We offer commercial support to meet the requirements of your business-critical applications by adding support you will have access to the latest versions, and have guaranteed response times should any problems arise.

We offer various licensing options for server licenses, integration and site licenses, and source code licenses. Support contracts for extended technical support with short response times and free updates are also available.

About PDFlib GmbH

PDFlib GmbH is completely focused on PDF technology. Customers worldwide use PDFlib products since 1997. The company closely follows development and market trends, such as ISO standards for PDF. PDFlib GmbH products are distributed all over the world with major markets in North America, Europe, and Japan.

Contact

Fully functional evaluation versions including documentation and samples are available on our Web site. For more information please contact:

PDFlib

PDFlib GmbH





PDFlib pCOS PDF Information Retrieval Tool

What is PDFlib pCOS?

PDFlib pCOS provides a simple and elegant facility for retrieving any information from a PDF document which is not part of the page contents. For example, PDF metadata, interactive elements (links, form fields, etc.), or page dimensions can easily be queried with pCOS.

With pCOS you can extract a variety of interesting items and create output for different purposes. By processing multiple PDF documents with a single call you can easily create summaries of document info entries, page formats, fonts, or any other property. Combined with tabular output this provides a powerful PDF administration tool.

There are many application scenarios for the PDF Information Retrieval Tool PDFlib pCOS within PDF workflows, but you can also use PDFlib pCOS as a tool for learning or debugging PDF. Here are some typical situations:

- Check incoming documents for predefined criteria
- Identify problem files in a large collection
- Create metadata summaries for document management
- quality assurance before publishing document
- document retrieval and repository workflow
- summarize the bookmarks
- ▶ extract components of PDF documents, g. ICC profiles
- ► Check PDFs for security problems (Java) cript etc.)

The pCOS retrieval interface sincluded in other PDFlib GmbH products: if you use PDFlib PDI, DFlib Personalization Server, TET or PLOP you also have access to the pCOS interface. If you need access to text or images on the page use our product PDFlib TET for PDF content extraction

pCOS Cookbook

The pCOS Cook took is a collection of programming examples which demonstrate the use of pCOS for various PDF retrieval tasks. The Cookbook is available on the PDFlib Web site and includes sample code, input documents and sample output.

PDFkb pCOS Features

Supported Input

PDF b pCOs supports all flavors of PDF input:

- All PDF versions up to Acrobat X, including ISO 32000
- Encrypted documents (password may be required)
- Sophisticated security model: even if you don't know the password, you can query certain pieces of information as long as this
- doesn't violate the document author's intentions ► Daw aged PDF input documents will be repaired if possible

Information Retrieval

PD lib pCOS offers a simple query interface. With PDFlib pCOS you can extract a variety of interesting items, such as:

- ► Document info entries and XMP metadata
- General information: linearization and tagged PDF status, encryption details and permission settings, number of pages and fonts
- ► Fonts with name, embedding status, etc.
- Image data, such as bit depth, color space, compression
 - ► Color space details
- ► Target URLs and coordinates of Web links
- Bookmarks and the corresponding page numbers, e.g. to create a table of contents
- ► Form field data: full field names, contents, position, etc.
- ► Page size, CropBox, page rotation
- ► Status of ISO standards: PDF/X, PDF/A, PDF/UA, PDF/E, and PDF/VT
- Geospatial reference information
- ► List or extract file attachments
- ► Layer names, page labels, article threads
- Annotation details
- ► List all comments along with the reviewer's name
- Digital signature details: name of signature field(s), signed/unsigned, name of signer, date and reason of signature
- Extract ICC output intent profiles from PDF/X or PDF/A documents
- ► Block properties for PDFlib Personalization Server
- ► JavaScript on document, page, annotation, or field level

Output Formats

PDFlib pCOS can create output for different purposes:

- ► Plain text output
- ► Unicode text output in UTF-8 or UTF-16 formats
- Tabular output for processing with a spreadsheet/database
- ▶ Binary data, e.g. ICC profiles or file attachments
- ► User-defined output formats for custom post-processing

pCOS Paths: Simple Syntax for PDF Objects

Instead of getting bogged down by complex tree structures, e.g. for bookmarks or form fields, you can easily access PDF objects by using the simple pCOS path syntax. It offers convenient shortcuts for accessing commonly used PDF objects, such as pages, fonts, bookmarks, form fields etc.

pCOS Library or Command-Line Tool?

pCOS is available as a programming library (component) for many development environments, and as a command-line tool for batch operations. Both offer similar features, but are suitable for different deployment tasks.

The pCOS programming library is used...

... for integration into desktop or server applications. Examples for using the library with all supported language bindings are includ in the pCOS package.

The pCOS command-line tool is suited...

...for batch processing PDF documents. It doesn't require any programming, but offers powerful command-line options which can be used to integrate it into complex workflows. The pCOS command-line tool extends the features of the library:

- Simple retrieval of common PDF elements, such as bookmarks, annotations, metadata, form fields, etc.
- Extended mode for querying more complex objects and customer and cu izing the output format
- Extract data items such as file attachments, ICC profi s, etc. Emit information as comma-separated values or a user-definition
- format for import into a spreadsheet or database ► Recursion feature for dumping composite PDF objects, such as dictionaries and arrays

Supported Development Environments

PDFlib pCOS is everywhere – it runs on practice v al computing platforms. We offer 32-bit and 64-bit packages for all common flavors of Windows, Mac OS X, Linux and Uni

The pCOS core is written in highly optimized C and C++ code for maximum performance and small overhead. Via a simple API (Application Programming Interface) the pCOS functionality is accessible from a variety of development environments:

- COM for use with VX ASP and many other languages
- ► C and C++
- ► Java, including services d Java Application Server
- ith C#_VB.NET, ASP.NET, etc. ► .NET for use
- ► Perl





Benefits of using PDFlib Software

Rock-solid Products

Tens of thousands of programmers worldwide are working w our software. PDFlib meets all quality and performance requirements for server deployment. All PDFI products are suitable for robust 24x7 server deployment and unattended batch processing.

Speed and Simplicity

PDFlib products are incredibly fast – up to thousands of pages per second. The programming interface is straightforward and easy to learn.

PDFlib all over the World

Our products support all international languages as well as Unicode. They are used by customers in all parts of the world.

Professional Support

Licensing

If there's , we will try to help. We offer commercial supproblem tto meet the requirements of your business-critical applications y adding support yoy will have access to the latest versions, nd have guaranteed response times should any problems arise.

offer various licensing programs for server licenses, integration and site licenses, and source code licenses. Support contracts for extended technical support with short response times and free opdates are also available.

About PDFlib GmbH

PDFlib GmbH is completely focused on PDF technology. Customers worldwide use PDFlib products since 1997. The company closely follows development and market trends, such as ISO standards for PDF. PDFlib GmbH products are distributed all over the world with major markets in North America, Europe, and Japan.

Contact

Fully functional evaluation versions including documentation and samples are available on our Web site. For more information please contact:

PDFlib

PDFlib GmbH