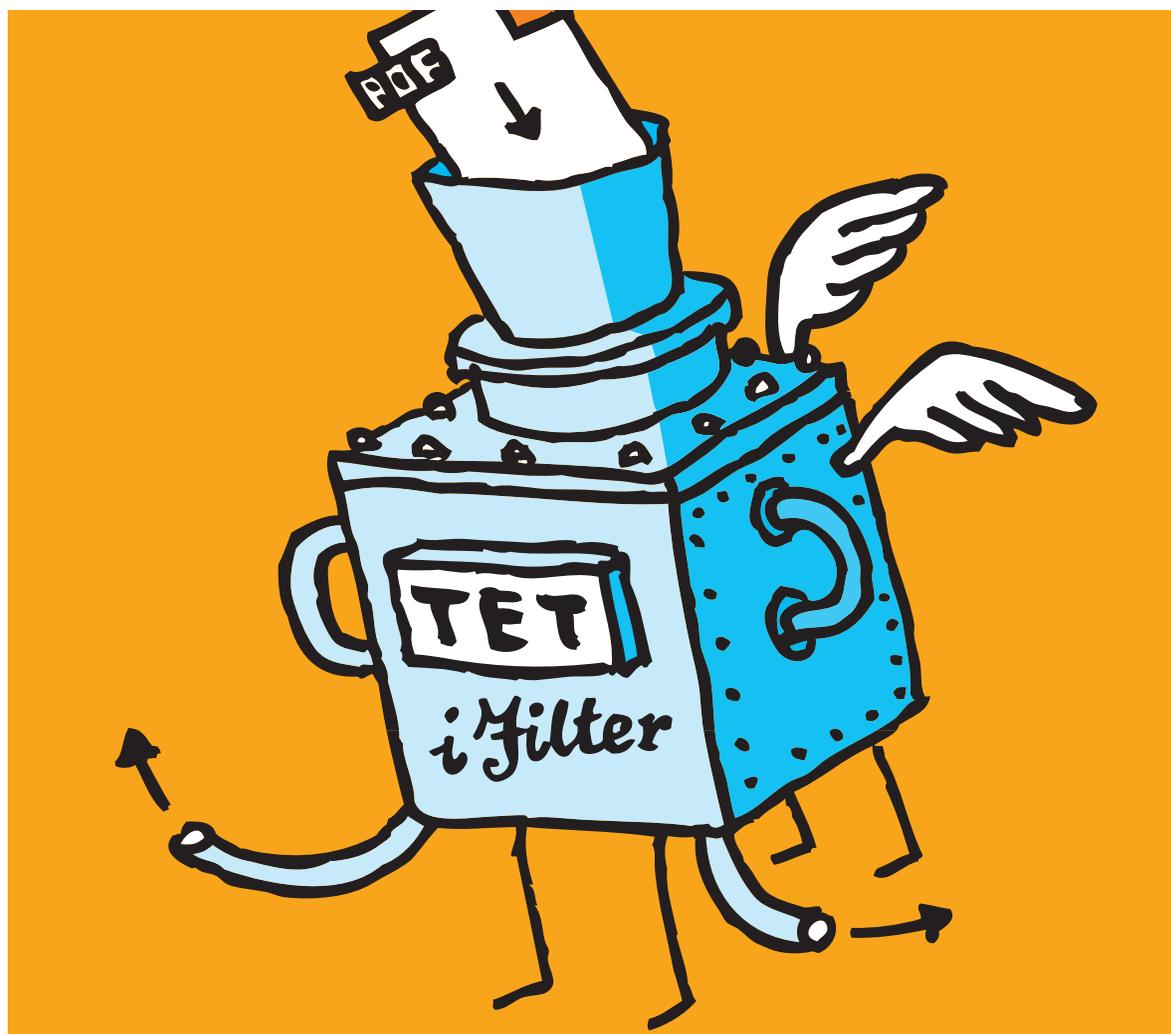


TET PDF IFilter

Version 4.1

Windows 用ビジネス向け PDF 検索



Copyright © 1997–2012 PDFlib GmbH and Thomas Merz. All rights reserved.
Protected by European and U.S. patents.

PDFlib GmbH
Franziska-Bilek-Weg 9, 80339 München, Germany
www.pdflib.com
電話 +49・89・452 33 84-0
fax +49・89・452 33 84-99

疑問がおありの際は、PDFlib メーリングリストと、tech.groups.yahoo.com/group/pdflib にあるアーカイブを
チェックしてください。

ライセンスご希望の際の連絡先 : jp.sales@pdflib.com
商用 PDFlib ライセンス保持者向けサポート : jp.support@pdflib.com (お使いのライセンス番号をお書きく
ださい)

この出版物およびここに含まれた情報はありのままに供給されるものであり、通知なく変更されることが
あり、また、PDFlib GmbH による誓約として解釈されるべきものではありません。PDFlib GmbH はいかな
る誤りや不正確に対しても責任や負担を全く負わず、この出版物に関するいかなる類の (明示的・暗示的
または法定に関わらず) 保証も行わず、そして、いかなるそしてすべての売買可能性の保証と、特定の
目的に対する適合性と、サードパーティの権利の侵害とを明白に否認します。

PDFlib と PDFlib ロゴは PDFlib GmbH の登録商標です。PDFlib ライセンス保持者は PDFlib の名称とロゴを彼
らの製品の文書内で用いる権利を与えられます。ただし、これは必須ではありません。

Adobe・Acrobat・PostScript・XMP は Adobe Systems Inc. の商標です。AIX・IBM・OS/390・WebSphere・
iSeries・zSeries は International Business Machines Corporation の商標です。ActiveX・Microsoft・OpenType・
Windows は Microsoft Corporation の商標です。Apple・Macintosh・TrueType は Apple Computer, Inc. の商標
です。Unicode・Unicode ロゴは Unicode, Inc. の商標です。Unix は The Open Group の商標です。Java・
Solaris は Sun Microsystems, Inc. の商標です。HKS は the HKS brand association: Hostmann-Steinberg, K+E
Printing Inks, Schmincke の登録商標です。他の企業の製品とサービス名は他の商標やサービスマークである
場合があります。

TET PDF IFilter は以下のサードパーティソフトウェアの変更された部分を含んでいます :
Zlib 圧縮ライブラリ、Copyright © 1995-2002 Jean-loup Gailly and Mark Adler
TIFFlib 画像ライブラリ、Copyright © 1988-1997 Sam Leffler, Copyright © 1991-1997 Silicon Graphics, Inc.
Eric Young の書いた Cryptographic ソフトウェア、Copyright © 1995-1998 Eric Young (ey@cryptsoft.com)
Independent JPEG Group の JPEG ソフトウェア、Copyright © 1991-1998, Thomas G. Lane
Cryptographic ソフトウェア、Copyright © 1998-2002 The OpenSSL Project (www.openssl.org)
Expat XML パーサ、Copyright © 1998, 1999, 2000 Thai Open Source Software Center Ltd
ICU International Components for Unicode、Copyright © 1995-2009 International Business Machines
Corporation and others

TET PDF IFilter は RSA Security, Inc. の MD5 メッセージダイジェストアルゴリズムを含んでいます。



目次

○ TET PDF IFilter をインストール 5

1 動作開始 7

- 1.1 Windows Search 7
- 1.2 SharePoint と FAST Search 10
- 1.3 Search Server 12
- 1.4 SQL Server 13
- 1.5 インデックスサービス 16

2 PDF の内容をインデックス 19

- 2.1 PDF 文書のさまざまな領域 19
- 2.2 自動言語検知 24
- 2.3 PDF の各種バージョンと暗号化文書 27
- 2.4 Unicode 後処理 28
 - 2.4.1 Unicode 字形統合 28
 - 2.4.2 Unicode 分解 30
 - 2.4.3 Unicode 正規化 33
- 2.5 カスタムグリフマッピングテーブル 35

3 メタデータをインデックス 37

- 3.1 PDF 内のメタデータ源 37
- 3.2 メタデータの編成 40
- 3.3 定義済みメタデータプロパティ 41
- 3.4 カスタムメタデータプロパティ 42
- 3.5 多値プロパティ 44
- 3.6 メタデータプロパティをテキストとしてインデックス 45
- 3.7 メタデータを採りページ内容を無視 47

4 各種 IFilter クライアントでメタデータ処理 49

- 4.1 Windows Search でメタデータ 49
- 4.2 SharePoint ・ Search Server でメタデータ 55
- 4.3 SQL Server でメタデータ 60

4.4 インデックスサービスでメタデータ 61

5 トラブルシューティング 65

5.1 TET PDF IFilter が全く動かない 65

5.2 TET PDF IFilter の動作上の問題 67

5.3 PDF 文書がインデックスされないか不完全 68

5.4 デバッグ機能 70

6 XML 構成ファイル 73

6.1 構成ファイルを使用 73

6.2 XML 要素・属性一覧 75

6.3 サンプル構成ファイル 80

A 定義済みメタデータプロパティ 81

B 更新履歴 87

索引 89

○ TET PDF IFilter をインストール

TET PDF IFilter は、Windows 各種システム用の MSI インストーラとして頒布されています。すべての TET PDF IFilter パッケージは、署名済みの IFilter DLL のほかに、サポートファイル群・文書・サンプル群を含んでいます。この MSI インストーラを走らせるには管理者権限が必要です。このインストーラは TET PDF IFilter をインストールして登録します。これに加えて個別の検索環境（Windows デスクトップサーチ・SharePoint など）で必要になる手順や、カスタム構成については、本マニュアル内で別途解説します。

32 ビット版と 64 ビット版 TET PDF IFilter には、32 ビットプラットフォーム用と 64 ビットプラットフォーム用があります。この 2 種類はインストーラが別々であり、同一システム上に共存してインストールする必要があるかもしれませんが、64 ビット版は、64 ビット実行形式とともにのみ動作するネイティブ 64 ビット実装です。64 ビットインストーラは 32 ビットシステムへのインストールを拒みますが、32 ビット版は 32 ビットシステムと 64 ビットシステムの両方で動作します。

TET PDF IFilter ライセンスキーを適用 TET PDF IFilter をサーバシステム上で業務目的に使用するには、有効なライセンスキーが必要です。TET PDF IFilter ライセンスを購入いただいたら、無制限に大きな文書の処理を許可するために、ご自分のライセンスキーを適用する必要があります。通常は、TET PDF IFilter をその MSI インストーラでインストールする際にそのライセンスキーを入力することになります。ただし、インストールの後でライセンスキーをレジストリ内で手動で適用することも可能です（6 ページの「手動インストール」を参照）。32 ビット版と 64 ビット版は同一のライセンスキーを受け入れます。

ライセンスキー 0（ゼロ）を使用すると、評価版をサーバシステム上にインストールすること、および無償デスクトップ版を非商用目的でインストールすることが可能です。

評価版の制限 TET PDF IFilter は、商用ライセンスがなくても完全機能の評価版として使用できます。非ライセンス版はすべての機能をサポートしていますが、ただし 10 ページかつ 1 MB サイズ以下の PDF 文書を処理することのみ可能になっています。TET PDF IFilter の評価版は、業務目的に使用してはならず、同製品の評価のためにのみ使用できます。TET PDF IFilter を業務目的に使用するには有効なライセンスが必要です。

無償デスクトップ版を非商用目的で使用 各種デスクトップシステム、すなわち Windows 2000/XP/Vista/7/8 用の TET PDF IFilter は、個人的利用のために、すなわち非商用目的であれば自由に使用できます。ただし、商用利用と見なされうる状況であれば必ず、このデスクトップ版を運用するにも商用ライセンスが必要です。

各種サーバオペレーティングシステム、すなわち Windows Server 2003/2008 用の TET PDF IFilter は、必ず商用ライセンスが必要です。

対応する IFilter クライアント TET PDF IFilter は、Microsoft の IFilter インタフェースを実装しています。さまざまなインデックス作成製品がこの IFilter インタフェースに対応しています。本文書ではそのような製品を IFilter クライアントと称します。TET PDF IFilter は下記の製品で試験済みですが、これ以外の、IFilter インタフェースに対応する Microsoft やサードパーティ製品でも動作する可能性があります：

- ▶ Windows Search。Windows Vista/7/8 のエクスプローラに内蔵されています。Windows XP 用の無料アドオンとしても入手可能です。
- ▶ インデックスサービス

- ▶ Windows SharePoint Services 3.0・SharePoint Portal Services 2003・Office SharePoint Server 2007・SharePoint Server 2010・FAST Search Server 2010 for SharePoint
- ▶ SharePoint Foundation 2010 を、Search Server か Search Server Express とともに使用の場合。なお、SharePoint Foundation 2010 はスタンドアロンでは、インデックス作成したいファイル種別の追加に対応しておらず、もともと入っているファイル種別群の中に PDF は含まれていません。
- ▶ SQL Server 2005・SQL Server 2008
- ▶ Search Server 2008・Search Server Express 2008、Search Server 2010・Search Server Express 2010
- ▶ Exchange Server 2007
- ▶ Site Server

TET PDF IFilter には、32 ビット版と 64 ビット版があります。この IFilter の 64 ビット版は、上記製品の 64 ビット版とともに動作します。

インストール後の手順 IFilter クライアントによっては、TET PDF IFilter のインストールが済んだ後、さらに手順をふむ必要がある場合があります。この手順については、7 ページの 1 章「動作開始」のそれぞれの節で解説しています。

手動インストール インストーラは、TET PDF IFilter を使用するために必要な手順をすべて適用しますが、場合によっては、特定の手順を手動で適用する必要が生じるかもしれません。その場合は下記の情報を参照してください。

ライセンスキーを手動で追加するには、下記のレジストリ値にそれを入力します：

```
HKEY_LOCAL_MACHINE\SOFTWARE\PDFlib\TET PDF IFilter4\license
```

TET PDF IFilter DLL をシステムに登録する (IFilter クライアントから確実に見つかるようにするために) には、下記のコマンド (インストレーションディレクトリが違う場合にはその部分を変更) をコンソールウィンドウで実行します：

```
regsvr32 "C:\Program Files\PDFlib\TET PDF IFilter 4.1\bin\TETPDFIFilter.dll"
```

注記 6 ページの「Windows Vista/7/8 で特権コマンドを実行」の項 (下記) を、それらのシステムで作業する場合には読んでください。

すでに TET PDF IFilter を使用している IFilter クライアントがある場合には、その DLL を再度登録する前に必ず、TET PDF IFilter にアクセスする関連サービスをすべて停止させてください。また、Windows イベントログも必ず閉じておいてください。

Windows Vista/7/8 で特権コマンドを実行 Windows Vista/7/8 の場合、レジストリへの書き込みアクセス (*regsvr32* プログラムや *registerpropdesc* プログラムなどによる) には管理者権限が必要です。コマンドプロンプトを管理者権限で起動させる方法は右記のとおりです：スタートをクリック→「コマンドプロンプト」項目を右クリック→「管理者として実行」を選択。すると、UAC プロンプトが表示されるので、確認の後、管理者権限で動作するコマンドプロンプトが開きます。

1 動作開始

この章では、TET PDF IFilter が対応しているいくつかの検索製品 (IFilter クライアント) を構成して使用するために必要な最初の手順を解説します。この解説は、まずは TET PDF IFilter が動作開始した状態まで導くことを目的としています。高度な構成事項については 37 ページの 3 章「メタデータをインデックス」で解説しています。完全な TET PDF IFilter インストールには、サンプル PDF 文書がひと揃い含まれていますので、それを利用すれば、インストールを試験したり、その強力なメタデータ諸機能に親しんだりすることができます。

1.1 Windows Search

システム要件 Windows デスクトップサーチ (WDS) 3 と、その後継である Windows Search は、レガシなインデックスサービスを置き換えるものです。それは、パフォーマンスと拡張性という特長をもたらすものであり、メタデータのための新しいプロパティシステムをサポートしており、TET PDF IFilter はこのプロパティシステムに完全対応しています。WDS の利用可能状況は、オペレーティングシステムのバージョンによって異なります：

- ▶ WDS 2 は Windows XP 用があります。WDS 3 の一部機能には、それ以前のバージョンである WDS 2.6.5 でも対応しています。TET PDF IFilter は、これらの種類すべてで動作します。
- ▶ WDS 3 は Windows Vista に内蔵されています。WDS 3.01 は、Windows XP SP2 と Windows Server 2003 SP1 の 32 ビット版用と 64 ビット版用の無償アドオンがあります。
- ▶ Windows Search 4 は、Windows Vista/7/8 用と、Windows XP SP2・Windows Server 2003 SP2 の 32 ビット版用・64 ビット版用と、Windows Server 2008 用があります。

セットアップと構成 デフォルトでは、Windows デスクトップサーチは「マイドキュメント」フォルダ内の文書群をインデックスします。それ以外の場所 (ネットワークドライブも含め) にある文書群をインデックスするよう WDS に指示することも可能であり、その方法は以下のとおりです：

- ▶ Windows タスクバーの通知領域内で、インデックス作成アイコン  を右クリック。
- ▶ 「Windows Search のオプション ...」をクリックし、そして「変更」をクリック。
- ▶ 「選択された場所の変更」セクション内で、望みの場所群を選び、そして OK をクリック。
- ▶ インデックス作成アイコンを再度右クリックし、「今すぐインデックスを作成」をクリックすると、その文書群がただちにインデックスされます。

Windows デスクトップサーチを開始/停止 検索サービスを (より正確には、TET PDF IFilter を呼び出すインデックス作成プロセスを) 手動で開始させたり停止させたりすることも可能であり、その方法は以下のとおりです：

- ▶ このサービスを開始させるには、コンソールウィンドウで下記をタイプします：

```
net start wsearch
```

このサービスを停止させるには下記をタイプします：

```
net stop wsearch
```

- ▶ このサービスをコントロールパネルで制御するには：
「スタート」→「設定」→「コントロールパネル」→「管理ツール」→「コンピュータ

の管理」をクリックします。「サービスとアプリケーション」→「サービス」へ進み
ます。利用可能なサービスの一覧の中で「Windows Search」を見つけます。そのサー
ビス名をダブルクリックすると、ダイアログが現れ、そこに「開始」／「停止」やそ
他のコントロールがあります。

- ▶ カタログを再構築するには：

インデックス作成アイコンを右クリックし、「Windows Search のオプション ...」→
「詳細設定」→「再構築」をクリックします。すると、すべての文書が再インデックス
されます。

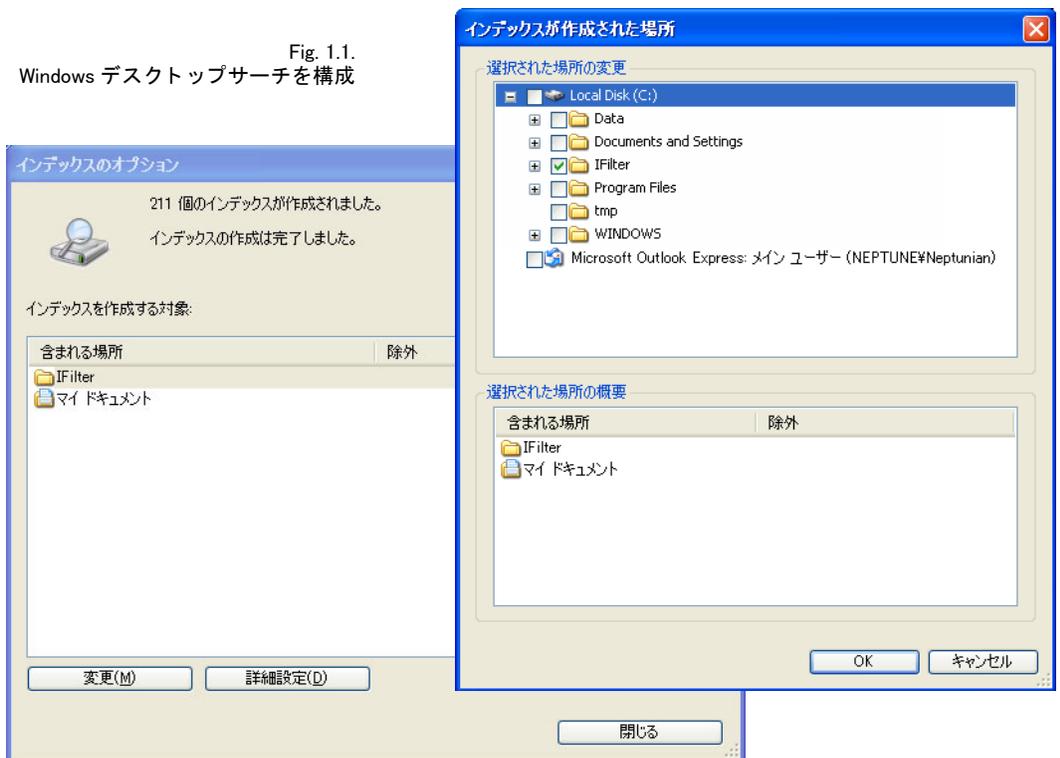
なお、WDS は、状況によっては自動的にインデックス作成サービスを開始させます。

シンプルなテキスト検索 テキストを検索するには、何通りかの方法が利用できます：

- ▶ WDSは、Windows タスクバー内の通知領域のとなりにデスクトップ検索ボックスを表示
しています。この検索ボックスに検索語群をタイプし、**Enter** を押します。
- ▶ 独立した検索ウィンドウを表示させることも可能であり、そのためには「スタート」→
「検索」→「ファイルとフォルダすべて」をクリックします。
- ▶ Windows エクスプローラのウィンドウ内で、上端にある「検索」ボタンを押します。

その方法を選んだとしても、それぞれのウィンドウ内で検索語群をタイプすることが
でき、WDS は、その検索語を含んでいる文書群（あれば）のヒットリストを表示します。何
らかの PDF ビューア（Adobe Acrobat や Adobe Reader など）がシステムにインストールさ
れていれば、選んだ文書の内容が結果ウィンドウ内に表示されます。

Fig. 1.1.
Windows デスクトップサーチを構成



高度なテキスト検索 シンプルに検索語群をタイプするだけでなく、高度なクエリ演算子を用いて検索を絞り込むことも可能です。表 1.1 に例をいくつか示します。高度なクエリのための文法要素の全一覧を見るには、Windows タスクバーの通知領域でインデックス作成アイコン  を右クリックし、「ヘルプ」→「高度な検索テクニック」をクリックします。メタデータのクエリについては、49 ページの 4.1 「Windows Search でメタデータ」を参照してください。

表 1.1 Windows Search 3.0 以上におけるクエリ文法の例

検索語の例	説明
Hol	文書の内容もしくは何らかのプロパティが、Hol で始まる語を含んでいる
Holmes AND Watson Holmes + Watson Holmes Watson (Holmes Watson)	文書内容が、Holmes と Watson を両方とも含んでいる
"Sherlock Holmes"	文書内容が、完全一致フレーズ Sherlock Holmes を含んでいる
Holmes OR Watson	文書内容が、Holmes か Watson のいずれかを含んでいる
Holmes NOT Mowgli Holmes -Mowgli	文書内容が、語 Holmes を含んでいるが、語 Mowgli を含んでいない

プログラマ的にインデックスをクエリ 対話的なクエリだけではなく、WDS インデックスをプログラマ的にクエリすることも可能です。検索インデックスへのアクセスは、AQS (*Advanced Query Syntax*) と SQL 文法拡張を通じて実現され、これは検索インデックスをデータベースライクなプログラミングインタフェースを通じて呈示します。WDS の SQL によるメタデータのクエリについては、52 ページの「メタデータプロパティに対する SQL クエリ」を参照してください。

制約 試験の結果我々は、下記の WDS 固有の制約に遭遇しました：

- ▶ 文書 1 個に対して生成されるカタログ項目が、サイズ約 1 MB を超えない。文書 1 個がそれより大きなカタログ項目を生成する場合には、WDS は残りのテキストを無視します。ページあたりテキストおよそ 3KB という典型的なサイズであれば、これは 300 ページ超のテキストに相当する量です。ただしこれは厳密なしきい値ではないので結果は変動する可能性があります。WDS でこの制約を回避する方法として唯一知られているのは、大きな文書を複数の小さな部分に分割することです。
- ▶ PDF 文書が Acrobat で開かれている間はインデックスされない。検索インデックスを作成する前には Acrobat を閉じることを推奨します。
- ▶ カスタムメタデータプロパティの対話的な検索が動作しない。プログラマ的な (PowerShell などによる) 検索は動作します。

1.2 SharePoint と FAST Search

システム要件 TET PDF IFilter は、下記の SharePoint 構成とともに動作します：

- ▶ SharePoint Server 2010、SharePoint Foundation 2010 を Search Server か Search Server Express とともに使用の場合。なお、SharePoint Foundation 2010 はスタンドアロンでは、インデックス作成したいファイル種別の追加に対応しておらず、もともと入っているファイル種別群の中に PDF は含まれていません。
- ▶ SharePoint 用 FAST Search Server 2010
- ▶ Microsoft Office SharePoint Server 2007 (略称 MOSS)
- ▶ SharePoint Portal Server 2003
- ▶ Windows SharePoint Services 3.0

SharePoint に関する詳しい情報は下記にあります：

msdn.microsoft.com/en-us/sharepoint/default.aspx

セットアップと構成 TET PDF IFilter を SharePoint とともに使用するように構成するには下記の手順を踏みます (SharePoint Foundation 2010 を Search Server か Search Server Express とともに使用の場合には、PDF アイコンを後述のように構成し、かつ TET PDF IFilter 自体を 12 ページの 1.3 「Search Server」に示すように構成します)：

- ▶ TET PDF IFilter を MSI パッケージからインストール。
- ▶ 後述のインストール後手順を適用。
- ▶ SharePoint 管理ウェブページを開く：「スタート」→「すべてのプログラム」→「Microsoft Office Server」→「SharePoint 3.0 サーバーの全体管理」。
- ▶ すると、ウェブページ「共有サービス管理：SharedServices1」が開きます。「検索」の下で「検索の設定」をクリック。「ファイルの種類」をクリック。
- ▶ 開いたページで「新しいファイルの種類」をクリックし、pdf ファイル拡張子を作成。
- ▶ 後述のように、PDF ファイルアイコンを構成することも可能です。

インストール後手順 TET PDF IFilter がインストールされている間 SharePoint が動作している場合、SharePoint クローラは新たにインストールされた PDF IFilter のことを、SharePoint 検索サービスが再起動されるまで知りません。これを再起動しないと、クローラから下記エラーメッセージを受けます：

```
The filtering process could not load the item.  
This is possibly caused by an unrecognized item format or item corruption.
```

MOSS 2010 の場合、SharePoint 検索サービスを再起動するには下記コマンドを適用します：

```
net stop osearch14  
net start osearch14
```

MOSS 2007 以前の場合、SharePoint 検索サービスを再起動するには下記コマンドを適用します：

```
net stop osearch  
net start osearch
```

PDF ファイルアイコンを追加 SharePoint クエリの結果一覧の中で PDF 文書を視覚的に表現するために、PDF アイコンを構成することもできます。このアイコンは PDF 文書に対して表示されます。下記のように操作します：

- ▶ PDF アイコンを下記からダウンロードし、
www.adobe.com/misc/linking.html#pdficon
pdficon_small.gif という名前で保存。
- ▶ SharePoint Server 2010 と SharePoint Foundation 2010 の場合、このアイコンを下記ディレクトリへ複製：
C:\Program Files\Common Files\Microsoft Shared\Web Server Extensions\14\TEMPLATE\IMAGES
SharePoint Server 2007 と Windows SharePoint Services 3.0 の場合、このアイコンファイルを下記ディレクトリへ複製：
C:\Program Files\Common Files\Microsoft Shared\Web Server Extensions\12\Template\Images
SharePoint Portal Server 2003 の場合、このアイコンファイルを下記ディレクトリへ複製：
C:\Program Files\Common Files\Microsoft Shared\Web Server Extensions\60\Template\Images
- ▶ SharePoint Server 2010 と SharePoint Foundation 2010 の場合、下記ディレクトリにある *Docicon.xml* 構成ファイルをテキストエディタに読み込み：
C:\Program Files\Common Files\Microsoft Shared\Web Server Extensions\14\TEMPLATE\XML
SharePoint Server 2007 と Windows SharePoint Services 3.0 の場合、下記ディレクトリにある *Docicon.xml* 構成ファイルをテキストエディタに読み込み：
C:\Program Files\Common Files\Microsoft Shared\Web server extensions\12\Template\Xml
SharePoint Portal Server 2003 の場合、下記ディレクトリにある *Docicon.xml* 構成ファイルをテキストエディタに読み込み：
C:\Program Files\Common Files\Microsoft Shared\Web server extensions\60\Template\Xml
- ▶ <ByExtension> 要素の子要素として下記の行を追加：
<Mapping Key="pdf" Value="pdficon_small.gif"/>
- ▶ IIS Web server を再起動。

シンプルなテキスト検索と高度なテキスト検索 SharePoint で検索クエリを構築するにはいくつかの方法があります：

- ▶ クエリキーワード
- ▶ SQL クエリ
- ▶ URL 内に符号化したクエリ

詳しくは下記を参照してください：

msdn.microsoft.com/en-us/library/ms497338.aspx

メタデータクエリについては 55 ページの 4.2 「SharePoint・Search Server でメタデータ」を参照してください。

1.3 Search Server

システム要件 TET PDF IFilter は、下記のエディションの Search Server とともに動作します：

- ▶ Search Server 2008 と Search Server 2008 Express
- ▶ Search Server 2010 と Search Server 2010 Express

セットアップと構成 Search Server の構成の方法は、SharePoint の構成の方法と非常に似ています。TET PDF IFilter を Search Server とともに使用するよう構成するには下記の手順を踏みます：

- ▶ TET PDF IFilter を MSI パッケージからインストール。
- ▶ 後述のインストール後手順を適用。
- ▶ Search 管理ウェブページを開く：「スタート」→「すべてのプログラム」→「Microsoft Search Server」→「Search Server 2008 Administration」。
- ▶ すると、ウェブページ「Search Administration」が開きます。「Crawling」の下で「File Types」をクリック。
- ▶ 開いたページで「New File Type」をクリックし、pdf ファイル拡張子を作成。
- ▶ 10 ページの「PDF ファイルアイコンを追加」で SharePoint について説明したように PDF ファイルアイコンを構成することも可能です。

インストール後手順 TET PDF IFilter がインストールされている間 Search Server が動作している場合、Search Server クローラは新たにインストールされた PDF IFilter のことを、その検索サービスが再起動されるまで知りません。これを再起動しないと、クローラからエラーメッセージを受けます。

Search Server を再起動するには下記コマンドを適用します：

```
net stop spsearch  
net start spsearch
```

シンプルなテキスト検索と高度なテキスト検索 Search Server でクエリを構成する方法は SharePoint と同様です。詳しくは下記を参照してください：

msdn.microsoft.com/en-us/library/aag81100.aspx

1.4 SQL Server

システム要件 TET PDF IFilter は、下記のエディションの SQL Server とともに動作しません：

- ▶ SQL Server 2005 Workgroup ・ Standard ・ Enterprise エディション
- ▶ SQL Server 2005 Express Edition
- ▶ SQL Server 2008

SQL Server に関する詳しい情報は下記にあります：

msdn.microsoft.com/en-us/library/bb418498.aspx

SQL Server 2005 でのフルテキスト検索に関する詳しい情報は下記にあります：

msdn.microsoft.com/en-us/library/ms142571.aspx

セットアップと構成 SQL Server 内でフィルタ群の使用を完全に制御できるよう、インストーラは、SQL Server のどのインスタンスにおいても、自動的に TET PDF IFilter を登録することはしません。ですので TET PDF IFilter を手動で、下記の手順に従って登録する必要があります。いずれの方式を選んだ場合も、SQL Server のすべてのインスタンスに対して個別に IFilter を構成する必要があります。下記に示す 2 種類の構成方式があります。

構成方式 1 の詳細 この方式は、システムワイドにインストールされている IFilters にアクセスするよう SQL Server に指示します。TET PDF IFilter のインストーラは IFilter DLL をシステム内に登録しますので、これは評価に便利な方式です。ただし、これはシステム設定への依存を作り出しますので、業務環境では注意して用いられねばなりません。下記のように操作します：

- ▶ TET PDF IFilter を MSI パッケージからインストール。
- ▶ SQL Server Management Studio を走らせ、下記ステートメントを実行して、システムワイド文書フィルタ群を SQL Server のこのインスタンスで利用可能にします：

```
exec sp_fulltext_service 'load_os_resources', 1;  
GO
```

構成方式 2 の詳細 この方式は、SQL Server の特定のインスタンスに対して、システム依存のまったくない状態で TET PDF IFilter を構成します。これは方式 1 よりも手順が多いですが、業務システムに対して制御がよりよく効きます。下記のように操作します：

- ▶ TET PDF IFilter を MSI パッケージからインストール。このステップは、IFilter を SQL Server とともに登録しませんが、IFilter ライセンスキーがレジストリで必ず得られるようにするために必要です。
- ▶ インスタンス名からインスタンス番号を特定するために、レジストリ内で下記キーへ移動：

```
HKLM\SOFTWARE\Microsoft\Microsoft SQL Server\Instance Names
```

たとえば、SQL という名前のインスタンスが下記エントリを持つ場合：

```
HKLM\SOFTWARE\Microsoft\Microsoft SQL Server\Instance Names\SQL  
"MSSQLSERVER"="MSSQL.1"
```

これは、SQL という名前のインスタンスが次のステップでディレクトリ *MSSQL.1* を使用していることを意味します。

- ▶ TET PDF IFilter DLL を見つけます。下記のようなパスで見つけられます：

```
C:\Program Files\PDFlib\TET PDF IFilter 4.1\bin\TETPDFIFilter.dll
```

この DLL を、前ステップで特定した、SQL Server のターゲットインスタンスのためのバイナリ群を内容として持つディレクトリへ複製。例：

```
C:\Program Files\Microsoft SQL Server\MSSQL.1\MSSQL\Bin
```

必ず、*MSSQL.1* ではなく正しいディレクトリを用いてください。

- ▶ レジストリ内に下記キーを作成：

```
HKLM\SOFTWARE\Microsoft\Microsoft SQL Server\MSSQL.1\MSSearch\Filters\.pdf
```

このキーのデフォルト値として TET PDF IFilter のクラス ID を追加：

```
{47A1AF35-C345-475d-AE68-EB07E948BD07}
```

- ▶ レジストリ内に下記キーを作成：

```
HKLM\SOFTWARE\Microsoft\Microsoft SQL Server\MSSQL.1\MSSearch\CLSID\{47A1AF35-C345-475d-AE68-EB07E948BD07}
```

このキーに下記文字列値を追加：

```
Default value: TETPDFIFilter.dll
```

```
ThreadingModel: Both
```

構成を試験 いずれの方式の場合も、構成結果をチェックして、SQL Server の一つのインスタンスで TET PDF IFilter が利用可能になったことを確認することが可能です。下記ステートメントを用います：

```
SELECT document_type, path FROM sys.fulltext_document_types WHERE document_type = '.pdf'
```

下記のような結果出力行は、そのインスタンス向けに TET PDF IFilter が正しく構成されていることを示します（正確なパスは、お使いのインストールパスと、選択した構成方式によって異なります）：

```
.pdf C:\Program Files\PDFlib\TET PDF IFilter 4.1\bin\TETPDFIFilter.dll
```

フルテキスト PDF インデックス作成用データベーステーブルを作成 TET PDF IFilter は SQL Server によって、型 *varbinary(max)* の列に格納された PDF 文書に対するフルテキストインデックスを生成するために使用されます。このままでは文書種別がわからないので、ファイル拡張子を、そのテーブル内の別の列、いわゆる **型列** に格納する必要があります。型列は、任意の文字ベースのデータ型にすることができます。*VARCHAR(4)* を使用し、ファイル拡張子 *pdf* を格納しましょう。

下記ステートメント群は、1 個のサンプル PDF 文書を *data* 列に、そのファイル名を *name* 列に、その型を *extension* 列に持つ *DocumentTable* を作成します：

```
CREATE DATABASE TestDatabase
```

```
GO
```

```
USE TestDatabase
```

```
GO
```

```
CREATE TABLE DocumentTable
```

```
(pk INT NOT NULL IDENTITY CONSTRAINT DocumentTablePK PRIMARY KEY,
```

```
data VARBINARY(MAX), name VARCHAR(100), extension VARCHAR(4))
```

```
GO
```

```
INSERT INTO TestTable(data, name, extension) SELECT *, 'The_Hound_of_the_Baskervilles.pdf', 'pdf' FROM OPENROWSET(BULK 'C:\Program Files\PDFlib\TET PDF IFilter
```

```
4.1\PDF samples\The_Hound_of_the_Baskervilles.pdf', SINGLE_BLOB) AS Document
GO
```

これで、フルテキストインデックスを生成できます：

```
sp_fulltext_database 'enable'
GO
CREATE FULLTEXT CATALOG TestCatalog AS DEFAULT
GO
CREATE FULLTEXT INDEX ON DocumentTable (data TYPE COLUMN extension)
KEY INDEX DocumentTablePK
GO
```

フルテキストインデックスを削除・再生成 フルテキストインデックスを削除するには下記ステートメント群を用います：

```
USE TestDatabase
GO
DROP FULLTEXT INDEX ON DocumentTable
GO
```

フルテキストインデックスを再生成：

```
USE TestDatabase
CREATE FULLTEXT INDEX ON DocumentTable (data TYPE COLUMN extension)
KEY INDEX DocumentTablePK
GO
```

シンプルなテキスト検索と高度なテキスト検索 フルテキストインデックス内の個々の語をクエリすることが可能です：

```
SELECT name FROM DocumentTable WHERE CONTAINS(*, 'Watson')
GO
```

複数の語で構成されたフレーズを検索するには、そのフレーズをダブルクォーテーションでくくります：

```
SELECT name FROM DocumentTable WHERE CONTAINS(*, "Arthur Conan Doyle")
GO
```

これらの手順を、与えられた PDF サンプル群で演示するサンプルスクリプトが一つ、TET PDF IFilter とともにインストールされます。Transact-SQL における **CONTAINS** 述語に関する詳しい情報は下記にあります：

[msdn.microsoft.com/en-us/library/ms187787\(SQL.100\).aspx](http://msdn.microsoft.com/en-us/library/ms187787(SQL.100).aspx)

メタデータクエリについては 60 ページの 4.3 「SQL Server でメタデータ」を参照してください。

1.5 インデックスサービス

システム要件 インデックスサービス 3.0 は、Windows Search の前身であり、Windows NT・Windows 2000・Windows XP に内蔵されています。これは Windows デスクトップサーチ (WDS) にとってかわられたものの、依然広く使われています。TET PDF IFilter はインデックスサービス 3 とともに動作します。インデックスサービスに関する詳しい情報は下記にあります：

msdn.microsoft.com/en-us/library/aa163263.aspx

セットアップと構成 インデックスサービスを構成するには、管理コンソール (MMC) で下記のように操作します (図 1.2 参照)：

- ▶ 「スタート」→「設定」→「コントロールパネル」→「管理ツール」→「コンピュータの管理」をクリックし、「サービスとアプリケーション」→「インデックス サービス」→「System」へ移動。
- ▶ 「ディレクトリ」を右クリックし、「新規作成」→「ディレクトリ」を選択し、ローカルパスか UNC ロケーションを追加。そのファイル群へのアクセスが制限されている場合には、適切なユーザ名とパスワードも追加。
- ▶ 「インデックス サービス」を右クリックし、このサービスが動作していることを確認。もし動作していないなら「開始」をクリック。
- ▶ 「インデックスを付けるドキュメント」列には、インデックスされるのを待っている文書の数が表示されています。これが 0 になったら、そのカタログは完全となり、すべての文書をクエリすることが可能になります。

図 1.2
管理コンソールでインデックス
サービスの構成とクエリフォーム



インデックスサービスを開始/停止 インデックスサービスを手動で開始/停止するには下記の方法があります：

- ▶ コンソールウィンドウで、サービスを開始させるには下記をタイプ：

```
net start CiSvc
```

サービスを停止させるには下記をタイプ：

```
net stop CiSvc
```

- ▶ コントロールパネルでサービスを制御するには：
「スタート」→「設定」→「コントロールパネル」→「管理ツール」→「コンピュータの管理」をクリックし、「サービスとアプリケーション」へ移動。「インデックス サービス」を右クリックし、「開始」か「停止」が望むほうを選択。
- ▶ カタログを強制的に再構築させるには：
上記のようにしてインデックスサービスを停止。「System」を右クリックし、「すべてのタスク」→「カタログを空にする」を選択。インデックスサービスを再び開始。

シンプルなテキスト検索 別個の検索ウィンドウを開くには、「スタート」→「検索」→「ファイルとフォルダすべて」をクリックします。ただし、MMC を使ってシンプルな検索や高度な検索を試すことも可能です：

- ▶ MMC を開き（上述）、「インデックス サービス」→「System」をクリック。
- ▶ 「カタログのクエリ」を選択し、メインウィンドウで、キャプション「フリー テキスト クエリを入力してください」の下に何らかの検索語群を入力（図 1.2 参照）。
- ▶ 「検索」をクリックしてヒットリストを生成。

高度なテキスト検索 シンプルな検索の場合（上述）と同じですが、ただしクエリを発する前に「高度なクエリ」を選択します。このクエリ内では、検索を絞り込むためにさまざまな演算子を用いることも可能です。表 1.2 にいくつかの例を示します。クエリ文法の解説は下記にあります：

[msdn.microsoft.com/en-us/library/ms690538\(VS.85\).aspx](https://msdn.microsoft.com/en-us/library/ms690538(VS.85).aspx)

プロパティに関するクエリについては「高度なクエリ」を用います。プロパティベースのクエリの概要は下記にあります：

[msdn.microsoft.com/en-us/library/ms690492\(VS.85\).aspx](https://msdn.microsoft.com/en-us/library/ms690492(VS.85).aspx)

メタデータクエリについては 61 ページの 4.4 「インデックスサービスでメタデータ」を参照してください。

表 1.2 インデックスサービスのためのクエリ文法の例

検索語の例	説明
Holmes AND Watson Holmes & Watson	文書内容が Holmes と Watson の両方を含んでいる
"Sherlock Holmes" {phrase} Sherlock Holmes [/phrase]	文書内容が厳密にフレーズ Sherlock Holmes を含んでいる
Holmes OR Watson Holmes Watson	文書内容が Holmes か Watson のいずれかを含んでいる

表 1.2 インデックスサービスのためのクエリ文法の例

検索語の例	説明
Holmes NEAR Watson Holmes ~ Watson	文書内容が語 Holmes と語 Watson を近傍（距離 50 語以内）に含んでいる
Holmes AND NOT Mowgli	文書内容が語 Holmes を含んでいるが、語 Mowgli を含まない

プログラムのにインデックスをクエリ Microsoft は、インデックスサービスでプログラムのにクエリを発するさまざまな方法を提供しています：

- ▶ ADO オブジェクト。Visual Basic・VBScript などから利用できます。
- ▶ ISAPI Extensions。IIS Web サイト内の *.idq*・*.ida*・*.htx* ファイルから利用できます。
- ▶ OLE DB プロバイダ・ヘルパ関数。Visual C++ から利用できます。

下記の、プログラムのにインデックスをクエリする VBscript サンプルが、TET PDF IFilter とともにインストールされています：

- ▶ *list_documents.vbs* は、*System* カタログをクエリして、PDF 文書の一覧を、いくつかの基本的情報（ファイル名・サイズ・ページ数）とともに出力します。コマンドウィンドウで下記のようにすればこれを動作させることができます：

```
cscript list_documents.vbs
```

- ▶ *query.vbs* を使うと、*System* カタログに対してシンプルなクエリや複雑なクエリをコマンドラインから発することができます。コマンドウィンドウで下記のようにすればこのスクリプトでクエリを動作させることができます：

```
cscript query.vbs "@DocAuthor Doy*"
```

プログラムのにインデックスをクエリすることについて詳しい情報は下記にあります：

[msdn.microsoft.com/en-us/library/ms692266\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms692266(VS.85).aspx) (*file-based indexing*)

[msdn.microsoft.com/en-us/library/ms693094\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms693094(VS.85).aspx) (*Web indexing*)

インデックスサービスのクエリ言語のほかに、SQL 拡張を利用してクエリ関連語を受け付けることもできます。詳しくは下記を参照してください：

[msdn.microsoft.com/en-us/library/ms691676\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms691676(VS.85).aspx)

2 PDF の内容をインデックス

2.1 PDF 文書のさまざまな領域

PDF 文書とは、単なるページの束をはるかに超えたものです。それ以外の、注釈やしおりといったさまざまな場所にもテキストを内容として持ちえますし、Adobe の XMP 形式か固有の文書情報項目の形をとったメタデータを活用することもできます。PDF 文書内でテキストを内容として持ちうる場所を、PDF 文書の領域と呼びます。以下、PDF 文書のすべての領域を列挙して解説するとともに、それぞれのテキストを Acrobat で表示させる方法もあわせて示します。また、文書のすべての領域に対する TET PDF IFilter のデフォルト動作もこの列挙中に記していきます。簡単にいえば TET PDF IFilter は、関連するすべての場所にあるテキストをインデックスします。その結果として、ひと目ただけではなげヒットが生成されたのかかわからないような文書についても検索ヒットを得る場合があります。通常、IFilter クライアントには検索語ハイライト機能がありませんので、結果文書の中で検索語がどこにあるか見つけだす方法を知っておくことが重要です。検索されたテキストがページ内容本体以外の場所にあるかもしれないことを念頭においたうえで、TET PDF IFilter が検索ヒットを報告した PDF 文書内で検索テキストがどこにあるか見つけにくいときには、以下の列挙を参照してください。

以下の説明に関する注記：

- ▶ Acrobat による「複数の PDF」の検索とは、右記の類の検索のことです：「編集」→「検索」の後、「以下の場所にあるすべての PDF 文書：」プルダウンで PDF 文書群のフォルダを選択。
- ▶ 説明の中で、プロパティセットコレクション *documentXMP*・*imageXMP*・*shell*・*pdf*・*internal* について記している箇所がいくつかあります。これらは XML 構成ファイル内で有効にすることができます（41 ページの 3.3 「定義済みメタデータプロパティ」を参照）。デフォルトでは、*shell*・*internal* プロパティセットコレクションは有効にされており、*pdf*・*documentXMP*・*imageXMP* プロパティセットコレクションは無効にされています。プロパティセットコレクションについて詳しくは 41 ページの 3.3 「定義済みメタデータプロパティ」を参照してください。
- ▶ 記法 *@indexNestedPdf* は、XML 構成ファイル内の属性を意味します（75 ページの 6.2 「XML 要素・属性一覧」を参照）。

ページ上のテキスト ページ内容は、PDF 内においてテキストの主要な供給源です。ページ上のテキストはフォントで視覚表現され、PDF 内で利用可能な多様なエンコーディング技法の一つを用いて符号化されています。

- ▶ Acrobat 8/9/X での表示方法：ページ内容はつねに表示されています。
- ▶ Acrobat 8/9/X で一個の PDF を検索する方法：「編集」→「検索」または「編集」→「高度な検索」。Acrobat がグリフを Unicode 値へ正しくマップできない文書内のテキストを、TET PDF IFilter は処理できる可能性があります。このような状況では、TET Plugin を利用することができます。TET Plugin は、その自前の検索ダイアログを「Plug-Ins」→「PDFlib TET Plugin...」→「TET Find」で提供します。ただしこれは、完全な検索機能を提供するようには意図されていません。
- ▶ Acrobat 8/9/X で複数の PDF を検索する方法：「編集」→「検索」または「編集」→「高度な検索」の後、「検索する場所を指定してください。」の中で「以下の場所にあるすべての PDF 文書」を選択し、PDF 文書群が入っているフォルダへブラウザ。



図 2.1
Acrobat の
高度な検索ダイアログ

- ▶ TET PDF IFilter : ページ内容はデフォルトでインデックスされます。ただし特殊な場合においては、`@indexPageContents=false` でページ内容のインデックス作成を無効にすることも可能です。

定義済み文書情報項目 固有の文書情報項目はキー／値対です。

- ▶ Acrobat 8/9/X での表示方法 : 「ファイル」 → 「プロパティ ...」
- ▶ Acrobat 8/9/X で一つの PDF を検索する方法 : なし
- ▶ Acrobat 8/9/X で複数の PDF を検索する方法 : 「編集」 → 「検索」または「編集」 → 「高度な検索」の後、ダイアログの下端近くの「高度な検索オプションを使用」をクリック。「検索する場所 : 」プルダウンで PDF 文書群のフォルダを選択し、プルダウンメニュー「その他の条件」で「作成日」「更新日」「作成者」「タイトル」「サブタイトル」「キーワード」のいずれかを選択。
- ▶ TET PDF IFilter : 定義済み文書情報項目群は、プロパティセットコレクション *shell* が有効にされていれば、インデックスされます。

カスタム文書情報項目 標準項目に加えて、カスタム文書情報項目を定義することもできます。

- ▶ Acrobat 8/9/X での表示方法 : 「ファイル」 → 「プロパティ ...」 → 「カスタム」(無償の Acrobat Reader では利用できません)
- ▶ Acrobat 8/9/X での検索方法 : なし

- ▶ TET PDF IFilter : カスタム文書情報項目群は、カスタムプロパティ群がこの文書情報項目に基づいて定義されていれば、インデックスされます。

文書レベルの XMP メタデータ XMP メタデータは、拡張されたメタデータを内容とする XML ストリームから成ります。

- ▶ Acrobat 8/9/X での表示方法 : 「ファイル」 → 「プロパティ ...」 → 「その他のメタデータ ...」 (無償の Acrobat Reader では利用できません)
- ▶ Acrobat 8/9/X で一つの PDF を検索する方法 : なし
- ▶ Acrobat 8/9/X で複数の PDF を検索する方法 : 「編集」 → 「検索」または「編集」 → 「高度な検索」の後、「高度な検索オプションを使用」をクリック。「検索する場所 :」プルダウンで PDF 文書群のフォルダを選択し、プルダウンメニュー「その他の条件」で「XMP メタデータ」を選択 (無償の Acrobat Reader では利用できません)。
- ▶ TET PDF IFilter : 文書 XMP は、プロパティセットコレクション *documentXMP* が有効にされていれば、あるいはカスタムプロパティ群が文書 XMP に基いて定義されていれば、インデックスされます。

画像レベルの XMP メタデータ XMP メタデータは、画像・ページ・フォントといった文書構成要素にも付けることが可能です。しかし、XMP が通常見つかるのは画像レベルにおいてのみです (文書レベルに加えて)。

- ▶ Acrobat 8/9 での表示方法 : 「ツール」 → 「高度な編集」 → 「TouchUp オブジェクトツール」 → 画像を選択 → 右クリック → 「メタデータを表示 ...」 (無償の Acrobat Reader では利用できません)
- ▶ Acrobat X での表示方法 : 「ツール」 → 「コンテンツ」 → 「オブジェクトを編集」 → 画像を選択 → 右クリック → 「メタデータを表示 ...」 (無償の Acrobat Reader では利用できません)
- ▶ Acrobat 8/9/X での検索方法 : なし
- ▶ TET PDF IFilter : XMP 画像メタデータは、プロパティセットコレクション *imageXMP* が有効にされていれば、インデックスされます。

フォームフィールド内のテキスト フォームフィールドはページにかぶさって表示されます。しかし、技術的にはこれはページ内容の一部ではなく、別途のデータ構造によって表現されます。

- ▶ Acrobat 8 での表示方法 : 「表示」 → 「ナビゲーションパネル」 → 「フィールド」
- ▶ Acrobat 9 での表示方法 : 「フォーム」 → 「フィールドを追加または編集 ...」
- ▶ Acrobat X での表示方法 : 「ツール」 → 「フォーム」 → 「編集」 (無償の Acrobat Reader では利用できません)
- ▶ Acrobat 8/9/X での検索方法 : なし
- ▶ TET PDF IFilter : フォームフィールド群は、カスタムプロパティ群が pCOS 擬似オブジェクト *fields* に基づいて定義されていれば、インデックスされます。もしメインページ内容 (すなわちフォームフィールドキャプション群) が一定であるので必要でない場合は、ページ内容インデックス作成をページオプション *skipengines={text image}* で無効にすることも可能です。下記の抜粋に、XML 構成ファイル (73 ページの 6 章「XML 構成ファイル」を参照) のこの場合における該当箇所を示します :

```
<n:Tet>
  <n:TetOptions></n:TetOptions>
  <n:DocOptions></n:DocOptions>
  <n:PageOptions>skipengines={text image}</n:PageOptions>
```

```

</n:Tet>

<n:Filtering metadataHandling="propertyAndText"/>

<n:Metadata>
  <n:PropertySet guid="E9CDA960-D09A-43bc-AAAA-BBBBBBBBBBBB">
    <n:Property friendlyName="Formfield" identifier="2">
      <n:Source pdfObject="fields[*]/V"/>
    </n:Property>
  </n:PropertySet>
</n:Metadata>

```

注釈内のテキスト フォームフィールドと同様、注釈（ノート・テキスト注釈等）は別のデータ構造によって表現されます。注釈において関心対象となるテキスト内容は、その種類によって異なります。たとえば Web リンクの場合、そこで関心対象となる部分は URL でしょうし、それ以外の種類の注釈では、印字されるテキスト内容が意味を持つでしょう。

- ▶ Acrobat 8/9 での表示方法：「表示」→「ナビゲーションパネル」→「注釈」
- ▶ Acrobat X での表示方法：「注釈」→「注釈のリスト」
- ▶ Acrobat 8/9/X で一つの PDF を検索する方法：「編集」→「検索」→ボックス「注釈を含める」をチェック、または注釈のリストツールバーで「検索」ボタンをクリック
- ▶ Acrobat 8/9/X で複数の PDF を検索する方法：「編集」→「検索」または「編集」→「高度な検索」の後、「高度な検索オプションを使用」をクリック。「検索する場所：」プルダウンで PDF 文書群のフォルダを選択し、プルダウンメニュー「その他の条件」で「注釈」を選択。
- ▶ TET PDF IFilter: 注釈群は、プロパティセットコレクション *pdf* が有効にされていれば、インデックスされます。

しおり内のテキスト しおりは、直接にはページと結びついていませんが、特定のページへ飛ぶアクションを内容として持つ場合があります。しおりは、入れ子にして階層構造を形成させることもできます。

- ▶ Acrobat 8/9 での表示方法：「表示」→「ナビゲーションパネル」→「しおり」
- ▶ Acrobat X での表示方法：「表示」→「表示切り替え」→「ナビゲーションパネル」→「しおり」
- ▶ Acrobat 8/9/X で一つの PDF を検索する方法：「編集」→「検索」または「編集」→「高度な検索」の後、ボックス「しおりを含める」をチェック
- ▶ Acrobat 8/9/X で複数の PDF を検索する方法：「編集」→「検索」または「編集」→「高度な検索」の後、「高度な検索オプションを使用」をクリック。「検索する場所：」プルダウンで PDF 文書群のフォルダを選択し、プルダウンメニュー「その他の条件」で「しおり」を選択（無償の Acrobat Reader では利用できません）
- ▶ TET PDF IFilter: しおり群は、プロパティセットコレクション *pdf* が有効にされていれば、インデックスされます。

ファイル添付 PDF 文書はファイル添付を含むこともでき（文書レベルかページレベルで）、そのファイル添付自体が PDF 文書であることも可能です。

- ▶ Acrobat 8/9 での表示方法：「表示」→「ナビゲーションパネル」→「添付ファイル」
- ▶ Acrobat X での表示方法：「表示」→「表示切り替え」→「ナビゲーションパネル」→「添付ファイル」

- ▶ Acrobat 8/9/X での検索方法：「編集」→「検索」または「編集」→「高度な検索」の後、ボックス「添付ファイルを含める」をチェック（無償の Acrobat Reader では利用できません）。入れ子になった添付は再帰的に検索されません。
- ▶ TET PDF IFilter：PDF 添付群は、@indexNestedPdf=true であれば、再帰的にインデックスされます。

PDF パッケージ・ポートフォリオ Acrobat 8 (PDF 1.7) では、PDF パッケージという概念が導入されました。これはファイル添付に追加のプロパティを与えたものです。Acrobat 9 (PDF 1.7 拡張レベル3) ではこの概念を拡張して、PDF ポートフォリオが導入されました。

- ▶ Acrobat 8/9/X での表示方法：Acrobat は、パッケージ／ポートフォリオの表紙と中身の PDF 文書群を、PDF パッケージ専用のユーザインタフェース要素群を用いて表現します。
- ▶ Acrobat 8/9 で一つの PDF パッケージを検索する方法：「編集」→「検索」→「検索する場所：」プルダウンで「PDF パッケージ全体」を選択
- ▶ Acrobat X で一つの PDF パッケージを検索する方法：「編集」→「ポートフォリオ全体を検索」
- ▶ Acrobat 8/9/X で複数の PDF パッケージを検索する方法：なし
- ▶ TET PDF IFilter：パッケージ／ポートフォリオに含まれている PDF 文書は、@indexNestedPdf=true であれば、再帰的にインデックスされます。

PDF の各種規格とその他各種 PDF プロパティ この領域は明示的にテキストを含んではおらず、PDF/X・PDF/A ステータス、タグ付き PDF ステータスといった、PDF 文書のさまざまな固有プロパティを集めたコンテナとして用いられます。

- ▶ Acrobat 8 での表示方法：なし
- ▶ Acrobat 9：「表示」→「ナビゲーションパネル」→「規格」（規格準拠 PDF でのみ現れます）
- ▶ Acrobat X：「表示」→「表示切り替え」→「ナビゲーションパネル」→「規格」（規格準拠 PDF でのみ現れます）
- ▶ Acrobat 8/9/X での検索方法：なし
- ▶ TET PDF IFilter：PDF プロパティ群は、プロパティセットコレクション pdf が有効にされていれば、インデックスされます。

2.2 自動言語検知

正しい言語検知の重要性 内容インデックス作成のいくつかの側面は、インデックスされた内容の言語独自の後処理によって著しく向上します。細部はそれぞれの IFilter クライアントによって異なるものの、下記の側面は高度に言語依存です：

- ▶ 単語区切り機能という構成要素。テキスト内容はこれを用いて単語群に区切られます。
- ▶ 単語語幹化機能。インデックスされた単語の原形を抽出します。この方法で検索クエリは、インデックスされた文書の中でその検索語が多少変形（屈折）された形で用いられていたとしても、ヒットを生成します。
- ▶ シソーラスベースの検索
- ▶ ノイズ語リスト。検索上意味を持たないと見なされる、それゆえインデックスへ入れ込まれない、*the・a・and・at・on* といった語群を内容として持っています。

自動言語検知 上記のさまざまな機能は、テキストがどの自然言語で書かれているかがわかる場合にのみ実装できるものです。デフォルトでは、さまざまな IFilter クライアントは、あらゆる言語独自処理に際してシステムロケールを用います。その結果、日本語文書を英語システム上でインデックスすれば品質二流のインデックス内容が生成されるので、検索にあてはまる文書があっても見つけそこなうおそれがあります。自然言語を正しく割り当てることは、東アジア言語の内容を持つ文書について特に重要です。

インデックス作成処理における言語独自処理のあらゆる側面を向上させるために、TET PDF IFilter は自動的に、文書内容群と型 *string* のプロパティ群の中のテキストについて自然言語を特定します。TET PDF IFilter は、その言語を検出するために 2 種類の技法を運用しています：

- ▶ いくつかの場合においては、用字系（表記体系）を調べるだけで十分に言語を特定できます。たとえば日本語のひらがなは、日本語でテキストを書くためにしか使われません。
- ▶ 多くの言語は、他の言語と同じ表記体系を共有しています。たとえば西洋語の大半は、ラテン用字系で書かれます。このような場合には TET PDF IFilter は、統計分析を適用して正しい言語を特定します。

自動言語検知は各テキスト塊を個別に分析しますので、1 個の文書に複数言語がどのようになでも混在してありえます。各切片に対応する言語が下記のように割り当てられます：

- ▶ 自動用字系検知が、テキストを複数の、同一用字系から成る塊に分割します。
- ▶ その用字系が言語を一意に特定しない場合には、自動言語検知を用いてテキストが分析されます。この処理は、塊全体についてもっともありうる言語を割り当てますが、その塊を各言語ごとにもっと小さな塊へさらに分割することはしません。たとえば、英語・ドイツ語・フランス語の混在したラテン字テキストを持つ塊は、その塊全体に対してこれらの言語のうちの 1 つを割り当てられます。どれを割り当てられるかは言語の分布状況によって決まります。

手動言語割り当て TET PDF IFilter における自動言語検知は、特殊な応用のためにカスタマイズすることも可能です。このカスタマイズ機能は LCID（ロケール識別子）を利用しています。LCID は自然言語を指定し、さらに、国別または地域別の言語変種どうし（英国英語と米国英語など）を区別することも可能です。表 2.1 に、よく使われる LCID 値をいくつか挙げます。LCID の全一覧は下記にあります：

[msdn.microsoft.com/en-us/library/ms776294\(VS.85\).aspx](https://msdn.microsoft.com/en-us/library/ms776294(VS.85).aspx)

注記 TET PDF IFilter は、言語検知より後の言語独自処理は一切適用しません。その LCID 情報を使うかどうかは IFilter クライアントに委ねられます。いくつかの IFilter クライアント (SharePoint・SQL Server など) は、LCID を取り扱う洗練された機能を持っていますが、それ以外の IFilter クライアントのなかには、LCID 情報を全く無視するものもあります。

表 2.1 よくある LCID 値と、対応する第一・第二言語

LCID	第一言語	第二言語 (国)
0x0000	中立ロケール言語	中立副言語
0x0401	アラビア語 (ar)	サウジアラビア (SA)
0x0404	中国語 (zh)	繁体字 (Hant)
0x0407	ドイツ語 (de)	ドイツ (DE)
0x0409	英語 (en)	米国 (US)
0x040c	フランス語 (fr)	フランス (FR)
0x0410	イタリア語 (it)	イタリア (IT)
0x0411	日本語 (ja)	日本 (JP)
0x0413	オランダ語 (nl)	オランダ (NL)
0x0419	ロシア語 (ru)	ロシア (RU)
0x0804	中国語 (zh)	簡体字 (Hans)
0x0c0a	スペイン語 (es)	スペイン (ES)
0x0800	システムデフォルトロケール言語	
0x1000	不特定カスタムロケール言語	不特定カスタム副言語

LCID のための XML 構成 自動 LCID 検知を上書きまたは補足するための LCID は、XML 構成ファイルの *LocaleId* 要素内で指定することができます：

```
<LocaleId default="1031" detection="auto" latin="1031" cyrillic="1026" chinese="4100"/>
```

detection 属性は値 *auto*・*disabled*・*script* をとりえます。他の属性は *default* 以外すべて、*detection=disabled* の場合には無視されます。デフォルトは *auto* です。*script* 設定は、用字系分析を有効にしますが統計分析を無効にします。

default 属性を使うと、*detection=disabled* の場合にすべてのテキストについて用いられるグローバルな LCID 設定を指定することが可能です。この属性がないときにはシステムロケールが用いられます。

detection 以外の属性ではすべて、10 進か 16 進文法の数値を指定できます。16 進値は先頭が *0x* でなければなりません。表 2.2 に、対応している用字系属性と、そのデフォルト値を挙げます。これ以外のあらゆる用字系のテキストに対する LCID は自動的に割り当てられます。

表 2.2 Attributes for specifying script-specific locale IDs (LCIDs)

属性	デフォルト値
default	0x0800 カレントシステムロケール (detection=disabled の場合にすべてのテキスト塊に対して用いられます)
latin	0x0409 英語 (US)
cyrillic	0x0419 ロシア語 (RU)
arabic	0x0401 アラビア語 (SA)
chinese	0x0804 中国語 (中華人民共和国)

2.3 PDF の各種バージョンと暗号化文書

PDF の各種バージョン TET PDF IFilter は、Acrobat X のファイル形式である PDF 1.7 拡張レベル 8 までのすべての PDF バージョンを受け付けます。これは、PDF/A・PDF/E・PDF/X といった PDF ベースのさまざまな ISO 規格を含んでいます。なお、ISO 32 000-1 は技術的に PDF 1.7 と同等ですので、これにも対応しています。ISO 32 000-2 で規定されている暗号化方式についても同様です。

暗号化 PDF 文書 TET PDF IFilter は、それが開くことができる限りにおいてあらゆる文書から、テキストとメタデータをインデックスします。これは下記の種類の PDF 文書を含みます：

- ▶ 非暗号化文書。
- ▶ マスターパスワードで暗号化されている、しかしユーザーパスワードを一切必要としない文書。Acrobat のセキュリティ設定「内容のコピー：許可／許可しない」の状態は、このグループに属する文書に影響を与えません。

一見、この 2 番目のカテゴリは文書作成者の文書保護意図の侵害に見えるかもしれませんが、しかしそうではないのです。なぜなら TET PDF IFilter はテキスト自体を複製する手段を一切提供していないからです。それは単に検索エンジンがその文書をインデックスすることを、そして後に検索でその文書を見つけだすことを可能にしているにすぎません。その文書が検索で特定されて Acrobat で開かれた後は、やはりそれは、もしその文書について内容複製に関する制限が指定されていたならば、それが何であれ従うのです。

開くことができない暗号化 PDF 文書はログ記録されます。このカテゴリは下記の場合を含みます：

- ▶ ユーザーパスワードを必要とする暗号化文書、すなわち、対応するパスワードを与えなければ Acrobat で開くことができない文書。
- ▶ 非暗号化文書内の暗号化 PDF 添付。
- ▶ ユーザー独自のセキュリティ証明書で暗号化されている文書。

破損 PDF 文書 PDF 文書の中に、破損したデータ構造が含まれてしまっている場合があります。これは PDF 生成ソフトウェアのバグや、あるいは何らかのアクシデントで変更が加わってしまったこと（ネットワーク転送が失敗したなどにより）が原因です。TET PDF IFilter は自動的に破損 PDF 文書を検知し、そのような文書を修復しようと試みます。これはテキストとメタデータをうまく抽出できるようにするためです。この修復モードは、インデックス作成処理の中で自動的に動作します。場合によっては、このモードでは充分ではない場合もあり、その場合には TET PDF IFilter はその文書をもっと徹底的な修復モードで処理する必要が出てきます。これはもっと時間がかかりますので、この強制修復モードは、自動修復モードではうまく処理できない、深刻に破損した PDF にのみ適用されます。

文書が、うまく開くことができたけれども、破損したページを 1 個ないし複数含んでいる場合には、それらのページは無視され、それに続くページ群について処理は続行されます。無視されたページそれぞれについて、アプリケーションイベントログに書き込みが行われます。

2.4 Unicode 後処理

TET PDF IFilter のために PDF テキスト抽出エンジンを提供している基盤である TET カーネルは、Unicode 後処理のためのさまざまな制御手段を提供しています。これらについては、TET PDF IFilter パッケージに含まれている TET マニュアルで詳しく解説しています。以下、もっとも重要な機能にしばって見ていきます。

さまざまな Unicode 後処理機能は、TET 文書オプション群によって制御されます。TET PDF IFilter では、これらのオプションは XML 構成ファイルの *DocOptions* 要素内で提供することが可能です。例：

```
<Tet>
  <DocOptions>decompose={canonical=_all}</DocOptions>
  <PageOptions/>
  <TetOptions/>
</Tet>
```

2.4.1 Unicode 字形統合

字形統合は、1 個ないし複数の Unicode キャラクタを処理して、各キャラクタに対して特定の操作を適用します。以下の操作が利用可能です：

- ▶ キャラクタを温存。
- ▶ キャラクタを除去。
- ▶ 別の（固定の）キャラクタへ置き換え。

表 2.3 に、字形統合のさまざまな応用を演示するさまざまな *fold* オプションの例を挙げます。これらのサンプルのオプションは、文書オプションリストの中で与える必要があります。TET は、全 Unicode キャラクタから選択された部分集合に対して字形統合を適用することもできます。これを Unicode 集合といいます。その文法は TET マニュアルで解説しています。

表 2.3 さまざまな fold オプションの例

説明とオプションリスト	字形統合前	字形統合後
1 個の Unicode 集合内のすべてのキャラクタを除去		
ISO 8859-1 (Latin-1) 内のキャラクタのみを出力内に保持、すなわち、基本欧文ブロック外のすべてのキャラクタを除去： fold={{[^U+0020-U+00FF] remove}}	A U+0104	なし
すべての非アルファベットキャラクタ（句読点・数字など）を除去： fold={{[:Alphabetic=No:] remove}}	7 U+0037	なし
	A U+0041	A U+0041
数字以外のすべてのキャラクタを除去： fold={{[^[:General_Category=Decimal_Number:]] remove}}	7 U+0037	7 U+0037
	A U+0041	なし
すべての未知キャラクタを、すなわち、PUA キャラクタ群と、有用な Unicode 値が決定できなかったキャラクタ群とを除去（残りのデフォルト字形統合は再有効化されます）： fold={{[:Private_Use:] remove} {[U+FFFD] remove} default}	U+FFFF	なし
すべてのダッシュ約物キャラクタを除去します： fold={{[:General_Category=Dash_Punctuation:] remove}}	- U+002D	なし

表 2.3 さまざまな fold オプションの例

説明とオプションリスト	字形統合前	字形統合後
<p>すべての双方向制御キャラクタを除去します： fold={{[:Bidi_Control:] remove}}</p>	U+200E	なし
1 個の Unicode 集合内のすべてのキャラクタを 1 個の特定のキャラクタへ置き換え		
<p>空白字形統合：Unicode 空白キャラクタのすべての種類を U+0020 へマップ： fold={{[:blank:] U+0020}}</p>	U+00A0	U+0020
<p>ダッシュ字形統合：Unicode ダッシュキャラクタのすべての種類を U+002D へマップ： fold={{[:Dash:] U+002D}}</p>	- U+2011	- U+002D
<p>すべての未割り当てキャラクタ（すなわち、キャラクタが割り当てられていない Unicode コード点）を U+FFFD へマップ： fold={{[:Unassigned:] U+FFFD}}</p>	U+03A2	◆ U+FFFD
個別のキャラクタに対して特別処理		
<p>改行位置のすべてのハイフンキャラクタを温存しつつ、残りのデフォルト字形統合を保持。これらのキャラクタは TET によって内部的に識別されていますので（固定の Unicode プロパティを持つのではなく）、この字形統合のドメインを指定するにはキーワード <code>_dehyphenation</code> を用います： fold={{_dehyphenation preserve}}</p>	- U+002D	- U+002D
<p>アラビア文字のタトゥールキャラクタ（デフォルトでは除去されます）を温存： fold={{[U+0640] preserve}}</p>	- U+0640	- U+0640
<p>さまざまな約物キャラクタをおのおの対応する ASCII キャラクタへ置き換え： fold={{[U+2018] U+0027} {[U+2019] U+0027} {[U+201C] U+0022} {[U+201D] U+0022}}</p>	" U+201C	" U+0022

2.4.2 Unicode 分解

正準分解 互いに正準等価であるキャラクタまたはキャラクタ列は、同一の抽象キャラクタを表現しており、したがって必ず同一の体裁と動作を持つべきものです。よくある例としては、合成済みキャラクタ（例： Ä ）と連結キャラクタ列（例： $\text{A} \text{ ¨}$ ）が挙げられます：2つの表現は互いに正準等価です。1つの表現から別の表現へ切り替えても情報は失われません。正準分解は、1つの表現を、正準表現と見なされる別の表現へ置き換えます。

Unicode コードチャート¹では、正準マッピングには記号 IDENTICAL TO \equiv U+2261 が付けられています（キャラクタテーブルにはなし）。分解名 `<canonical>` と暗黙的に見なされません。表 2.4 にいくつかの例を挙げます。

下記の文書オプションは、すべての正準等価キャラクタをおのおの対応キャラクタへマップします：

```
decompose={canonical=_all}
```

表 2.4 正準分解：decompose オプションのサブオプション（正準等価なキャラクタには、Unicode コードチャートで記号 IDENTICAL TO \equiv U+2261 が付けられています）

分解名	説明	分解前	分解後
<code>canonical¹</code>	正準分解	Ä U+00C0	$\text{A} \text{ ¨}$ U+0041 U+0308
		林 U+F9F4	林 U+6797
		Ω U+2126	Ω U+03A9
		ば U+3070	$\text{は} \text{ ¨}$ U+306F U+3099
		𐀀 U+FB2F	𐀀 U+05D0 U+05B8

1. デフォルトではこの分解は、特定のキャラクタ群を温存するために、すべてのキャラクタに対しては適用されません。詳しくは TET マニュアルを参照してください。

互換分解 互換等価なキャラクタどうしは、同一の抽象キャラクタを表しますが、その体裁や動作は互いに異なっている可能性があります。例としては、アラビア文字キャラクタの単独形（ س 等）と位置依存表示形（ $\text{س} \cdot \text{س} \cdot \text{س}$ 等）が挙げられます。互換等価キャラクタは組版上互いに異なるものです。この組版情報を除去すると、情報が失われる可能性があります。ある種の応用の場合（検索等）には処理を単純化できそうです。互換分解は組版情報を除去します。

Unicode コードチャートでは、互換マッピングには記号 ALMOST EQUAL TO \approx U+2248 が付いており、その後、分解名（「タグ」ともいいます）が `<noBreak>` のように山括弧にくくって書かれています。タグ名が書かれていないものについては、`<compat>` と見なさ

1. www.unicode.org/Public/5.2.0/charts/ を参照。

表 2.5 互換分解：decompose オプションのサブオプション一覧（互換等価なキャラクタには、Unicode コードチャートで記号 ALMOST EQUAL TO \approx U+2248 が付けられています）

分解名	説明	分解前	分解後（論理順で）
<i>circle</i>	丸囲みキャラクタ	② U+3251	2 1 U+0032 U+0031
<i>compat</i> ¹	その他の互換分解。例：広く用いられている合字	fi U+FB01	f i U+0066 U+0069
<i>final</i>	語尾形。特にアラビア文字の	س U+FEB2	س U+0633
<i>font</i>	フォントバリエント。例：数学の集合の文字、ヘブライ文字の合字	Ⓒ U+2102	Ⓒ U+0043
<i>fraction</i> ¹	俗用分数字体	¼ U+00BC	1 / 4 U+0031 U+2044 U+0034
<i>initial</i>	語頭形。特にアラビア文字の	س U+FEB3	س U+0633
<i>isolated</i>	単独形。特にアラビア文字の	سر U+FD0E	س ر U+0633 U+0631
<i>medial</i>	語中形。特にアラビア文字の	س U+FEB4	س U+0633
<i>narrow</i>	半角キャラクタ	ㄗ U+FF66	ㄗ U+30F2
<i>nobreak</i>	改行禁止キャラクタ	␣ U+00A0	␣ U+0020
<i>none</i>	明示的に decompose オプションリストで指定されていないすべての分解を無効化	(すべてのキャラクタを変更せず温存)	
<i>small</i>	CNS 11643 互換のための小型字体	， U+FE50	， U+002C
<i>square</i>	日中韓の組文字	キ □ U+3314	キ □ U+30AD U+30ED
<i>sub</i> ¹	下付き字体	₁ U+2081	₁ U+0031
<i>super</i> ¹	上付き字体	ₐ U+00AA TM U+2122	ₐ U+0061 T M U+0054 U+004D
<i>vertical</i>	縦書き字体	⎵ U+FE37	{ U+007B
<i>wide</i>	全角字体	£ U+FFE1	£ U+00A3

1. デフォルトではこの分解は、特定のキャラクタ群を温存するために、すべてのキャラクタに対しては適用されません。詳しくは TET マニュアルを参照してください。

れます。このタグ名は、表 2.5 のオプション名と同一です。いくつかの例に見られるように、分解の結果、1 個のキャラクタが複数キャラクタ列へ変換されることもあります。

下記の文書オプションは、全角（ダブルバイト）・半角キャラクタを温存します：

```
decompose={wide=_none narrow=_none}
```

2.4.3 Unicode 正規化

Unicode 規格では、正準等価と互換等価の概念に基づき、4 種類の正規形を定義しています。すべての正規形は、結合記号を特定の順序に置き、分解と合成を異なる方式で適用します：

- ▶ 正規形 C (NFC) は、正準分解の後に正準合成を適用します。
- ▶ 正規形 D (NFD) は、正準分解を適用します。
- ▶ 正規形 KC (NFKC) は、互換分解の後に正準合成を適用します。
- ▶ 正規形 KD (NFKD) は、互換分解を適用します。

正規形は、Unicode 規格付録 #15 「Unicode Normalization Forms」で定義されています (www.unicode.org/versions/Unicode5.2.0/cho3.pdf#G21796・www.unicode.org/reports/tr15/を参照)。

TET PDF IFilter は 4 種類の Unicode 正規形すべてに対応しています。Unicode 正規化は *normalize* 文書オプションで制御できます。例：

```
normalize=nfc
```

TET PDF IFilter はデフォルトでは正規化を適用しません。*decompose* オプションと *normalize* オプションはかち合うおそれがあることから、*normalize* オプションを *none* 以外の値に設定するとデフォルト分解は無効化されます。

正規形の選択はアプリケーションの要請に依存します。たとえば、いくつかのデータベースはテキストを NFC で受け付けます。NFC は、Web 上の Unicode テキストについて推奨される形式でもあります。表 2.6 に、さまざまなキャラクタに対する正規化の効果を示します。

表 2.6 Unicode 正規形：さまざまな例

正規化前	NFC	NFD	NFKC	NFKD
Ä U+00C4	Ä U+00C4	A ¨ U+0041 U+0308	Ä U+00C4	A ¨ U+0041 U+0308
A ¨ U+0041 U+0308	Ä U+00C4	A ¨ U+0041 U+0308	Ä U+00C4	A ¨ U+0041 U+0308
¨ A U+0308 U+0041				
fi U+FB01	fi U+FB01	fi U+FB01	f i U+0066 U+0069	f i U+0066 U+0069
3 5 U+0033 U+2075	3 5 U+0033 U+2075	3 5 U+0033 U+2075	3 5 U+0033 U+0035	3 5 U+0033 U+0035
Å U+212B	Å U+00C5	A ° U+0041 U+030A	Å U+00C5	A ° U+0041 U+030A
TM U+2122	TM U+2122	TM U+2122	T M U+0054 U+004D	T M U+0054 U+004D
IV U+2163	IV U+2163	IV U+2163	I V U+0049 U+0056	I V U+0049 U+0056

表 2.6 Unicode 正規形 : ささまざまな例

正規化前	NFC	NFD	NFKC	NFKD
ㄱ U+FB48	ㄱ U+05E8 U+05BC	ㄱ U+05E8 U+05BC	ㄱ U+05E8 U+05BC	ㄱ U+05E8 U+05BC
가 U+AC00	가 U+AC00	ㄱ ㅏ U+1100 U+1161	가 U+AC00	ㄱ ㅏ U+1100 U+1161
ぢ U+3062	ぢ U+3062	ぢ ㄷ U+3061 U+3099	ぢ U+3062	ぢ ㄷ U+3061 U+3099
10月 U+32C9	10月 U+32C9	10月 U+32C9	1 0 月 U+0031 U+0030 U+6708	1 0 月 U+0031 U+0030 U+6708

2.5 カスタムグリフマッピングテーブル

TET PDF IFilter にはさまざまな回避策が実装されてはいるものの、ある種のいくつかのまれな場合においては、PDF 文書が Unicode マッピングのための決定的な情報を持っていないければ、その文書からテキストを正しく抽出することが不可能なこともあります。

TET PDF IFilter は、PDFlib TET 製品のマニュアルで解説されているさまざまなテーブル形式に対応しています。無償の PDFlib FontReporter や、Adobe Acrobat 用 PDFlib TET Plugin を使って、そのようなマッピングテーブルを試用することも可能です。Unicode マッピングテーブルは、構成ファイル内で正しい文書オプション群を用いて構成する必要があり、TET PDF IFilter インストールディレクトリの *resource* ディレクトリ内に置くことができます。

TET オプション群のための XML 構成 TET コアの操作を制御するためのオプションリスト群 () は、TET 製品のマニュアルに解説されているオプションリスト文法に従って構築する必要があり、TET PDF IFilter に対して、XML 構成ファイルの *Tet* 要素の子要素である *DocOptions* ・ *PageOptions* ・ *TetOptions* 要素の中で与えることができます：

```
<Tet>
  <DocOptions>glyphmapping {{fontname=T* glyphlist={tex}}}</DocOptions>
  <PageOptions/>
  <TetOptions>searchpath={C:/glyphlists}</TetOptions>
</Tet>
```


3 メタデータをインデックス

3.1 PDF 内のメタデータ源

多くの PDF 文書は、「作成者」や「タイトル」などの文書情報項目群を内容として持っています。文書情報項目群だけではなく、PDF 文書は XMP メタデータを内容として持っている場合もあります。TET PDF IFilter は、PDF 文書内のいくつかの種類メタデータのインデックス作成に対応しています。

定義済みとカスタムの文書情報項目 文書情報項目群は、古くてシンプルな種類の PDF メタデータであるにとらえることができます。これらは Acrobat で「ファイル」→「プロパティ ...」で表示できます。下記の文書情報項目が PDF 内で定義済みです：

Title Author Subject Keywords Creator Producer CreationDate ModDate Trapped

定義済みの文書情報項目群に加えて、カスタムの項目群を文書情報項目の集合に追加することも可能です。これらは Acrobat で (Acrobat Reader では無理) 「ファイル」→「プロパティ ...」→「カスタム」で表示・編集することができます。

定義済み文書情報項目群も、カスタム文書情報項目群も、いずれも TET PDF IFilter 内で pCOS パスを用いて指し示すことができます (後述)。

文書レベルと画像レベルの XMP プロパティ XMP (*Extensible Metadata Platform*¹) は、メタデータを指定するためのフレームワークです。XMP はたとえば、PDF/A 準拠のために必須であり、対応アプリケーションは増加しつつあります。XMP メタデータは、多数のプロパティを内容とするスキーマ群の中に編成されています。プロパティは、1 個の名前空間接頭辞と、そのプロパティ名を用いて指し示されます。

XMP メタデータは通常、文書全体に関連づけられています。しかし、PDF では XMP を個別のページ・画像などのオブジェクトに関連づけることも可能になっています。実際には、この機能は主にラスタ画像について用いられています。たとえば、デジタル画像に、その撮影者・著作権表記・シーン詳細などの情報を帯びさせることができます。画像関連の XMP メタデータの重要な発生源の 1 つは Adobe Photoshop です。Photoshop から PDF を生成した場合、あるいは Photoshop から生成された画像をいろいろな Adobe のワークフローで使った場合、その XMP メタデータが通常、その PDF 文書内に入れ込まれており、TET PDF IFilter でインデックスすることが可能です。

XMP 仕様は、すべての定義済みの XMP スキーマ群とプロパティ群の記述を含んでいます。XMP について詳しい情報は下記にあります：

www.adobe.com/devnet/xmp

文書レベルか画像レベルの XMP プロパティ群は (これ以外のいろいろなオブジェクトに関連づけられた XMP も)、2 段階で指し示すことができます：pCOS パスは、関連する PDF オブジェクトを特定し、2 部分の XMP プロパティ名は、求める XMP スキーマとプロパティを指し示します。

XMP メタデータは、Acrobat や、Adobe Creative Suite のそれ以外のいろいろな製品で、Custom XMP パネル群で表示・編集することができます。これらは、特定の XMP メタデータのためのカスタムユーザーインターフェースの記述です。サンプル XMP パネル群が TET

1. www.adobe.com/products/xmp を参照。

PDF IFilter とともにインストールされています。詳しい情報とインストール方法説明は下記にあります：

www.pdflib.com/knowledge-base/xmp-metadata/xmp-panels/

拡張 pCOS パス pCOS (*PDFlib Comprehensive Object Syntax*) インタフェースは、PDF 文書のページ寸法・メタデータ・インタラクティブ要素といった、ページ内容を記述しないすべてのセクションから任意の情報を取得するためのシンプルでエレガントな機構を提供します。pCOS インタフェースの利用例や、pCOS パス文法の説明は、別の文書として入手可能な pCOS パスリファレンスに含まれています。さらなる利用例が、www.pdflib.com/pcos-cookbook/ にある pCOS Cookbook にあります。

pCOS は、TET PDF IFilter で、文書に関する情報を指し示すために用いることができます。pCOS パスは、その PDF 文書のしおり・フォント名・ページ・サイズといったいくつかの側面を表現し、また、文書情報項目や XMP メタデータストリームを指し示すこともできます(ただし XMP ストリーム内のプロパティを指し示すことはできません)。pCOS API 関数は TET PDF IFilter では利用できませんが、文書に関する情報をインデックスするために、いろいろな pCOS パスを XMP 構成ファイル内の表現として与えることが可能です。

表 3.1 に、よく用いられるいくつかの PDF オブジェクトに対する pCOS パスを挙げます。多くの pCOS オブジェクトは配列によって表現されており、これは配列添字を角カッコ内に必要とします。たとえば `pages[o]` は先頭ページを意味します(ページは 0 から始めてカウントされるのです)。TET 製品(TET PDF IFilter の基礎を成す)が対応しているすべての pCOS パスだけではなく、IFilter は、pCOS パスに対する右記の文法拡張に対応しています：「*」(アスタリスク) ワイルドカードキャラクタを、配列または辞書の添字のかわりに用いることが可能です。これは、TET PDF IFilter が、その配列添字の可能な値すべてをなめて、それぞれのオブジェクト値をすべて、インデックス作成処理に含めることを意味します。1 個の pCOS パス内で、任意の数のワイルドカードが使えます。

表 3.1 いろいろな PDF オブジェクトに対する pCOS パス

pCOS パス	型	説明
<code>length:pages</code>	数値	その文書内のページ数
<code>/Info/Title</code>	文字列	標準文書情報フィールド Title
<code>/Info/ArticleNumber</code>	文字列	カスタム文書情報フィールド ArticleNumber (文書情報項目は任意の名前を用いることができます)
<code>/Root/Metadata</code>	ストリーム	その文書のメタデータを持つ XMP ストリーム
<code>images[*]/Metadata</code>	ストリーム	その文書内の全画像に対する XMP メタデータストリーム群
<code>fonts[*]/name</code>	文字列	フォントの名前。項目数は <code>length:fonts</code> で取得できます
<code>length:fonts</code>	文字列	文書内のフォント数
<code>length:images</code>	文字列	文書内の画像数
<code>fonts[*]/embedded</code>	論理値	フォント内の埋め込みステータス
<code>pages[*]/width</code>	数値	ページの表示領域の幅
<code>pages[*]/annots[*]/A/URI</code>	文字列	全ページ上の Web リンク群のターゲット URL
<code>pages[*]/Resources/Properties[*]/Metadata</code>	ストリーム	その文書内の全ページ上の配置アートワークに対する XMP メタデータ
<code>fields[*]/V</code>	文字列	その文書内の全ページ上の全フォームフィールドのテキスト内容(値)

メタデータ源についてのXML構成 メタデータプロパティ群の源は、XMP 構成ファイルの *Source* 要素 (*Property* 要素の子) 内で構成することも可能です。*pdfObject* 属性は、PDF オブジェクトに対する拡張 pCOS パスを内容とし、一方、*xmpName* 属性は、XMP プロパティのスキーマ接頭辞と名前を内容とします：

```
<Source pdfObject="/Info/ArticleNumber"/>
<Source xmpName="acme:number"/>
```

3.2 メタデータの編成

メタデータは下記のように階層構造に編成されています：

- ▶ **プロパティ**は、メタデータのための基本的な構成単位です。Windows オペレーティングシステムと IFilter インタフェースの中のプロパティ群は、一意な数値識別子によって編成されています（後述）。
- ▶ **プロパティセット**は、プロパティ群のグループで構成され、そのプロパティ群は通常、何らかの論理的な関係を持っています。1つの集合の中のすべてのプロパティは同一の GUID（後述）を共有しています。プロパティセット群は XML 構成ファイル内で指定することができます。
- ▶ **プロパティセットコレクション**は、プロパティセット群で構成され、TET PDF IFilter は、いくつかの定義済みのプロパティセットコレクションを実装しています。これらを用いると、複数のプロパティセットを一緒にまとめて有効にしたり無効にしたりすることができます。追加のプロパティセットコレクションを構成する必要はありません。

プロパティ識別と GUID プロパティは、IFilter インタフェース内で、2 部分から成る識別子によって識別されます：

- ▶ 1 番目の部分は**グローバル一意識別子 GUID (ユニバーサル固有識別子 UUID と呼ばれることもあります)**です。これは、一意な 128 ビットの識別子であり、1つのプロパティセット内のすべてのプロパティに対して同一の値を持ちます。GUID の詳細は下記にあります：

www.itu.int/ITU-T/studygroups/com17/oid/X.667-E.pdf

GUID を作成するにはさまざまなツールが入手可能です。オンラインサービスを使うことも可能で、たとえば下記で利用可能です：

www.itu.int/ITU-T/asn1/uuid.html

GUID はたとえば右記のような形です：*7a737220-ocdo-11dd-bd75-0002a5d5c51b*。

- ▶ 2 番目の部分は、プロパティを、そのプロパティセット内部で一意に識別します。これは、**識別子**という、または略して **ID** という正の整数から成る場合があります。1つのセットの中のプロパティ識別子群は値 2 で始まる必要がありますが、それ以外については任意です。プロパティ識別子には、すべての IFilter クライアントが対応しています。あるいは、この 2 番目の部分は平文名から成る場合もあります。ID のかわりに名前を用いることは非推奨であり、Windows Search などいくつかの IFilter クライアントはこれに対応していません。しかし、SharePoint などこれに対応している IFilter クライアントにおいては、これで構成がより便利になる可能性があります。GUID+ 名前方式を有効にすることについては、43 ページの「プロパティの GUID+ 名前処理のための XML 構成」を参照してください。

GUID+ID または GUID+ 名前の組み合わせは、メタデータプロパティクエリを検索製品内で構成するために必須です。メタデータプロパティのこの他の側面は 42 ページの 3.4「カスタムメタデータプロパティ」で解説します。

3.3 定義済みメタデータプロパティ

下記のプロパティセットコレクションが TET PDF IFilter に内蔵されています：

- ▶ **シェルプロパティ群**は、Windows に知られており、ユーザーフレンドリな名前を持っています。TET PDF IFilter は、PDF 文書に同等物を持つすべてのシェルプロパティに内容を与えます。よくある例は **System.Author**・**System.Title**・**System.Document.DateCreated** です。シェルプロパティの一覧と説明は下記にあります：

[msdn.microsoft.com/en-us/library/windows/desktop/dd561977\(v=vs.85\).aspx](http://msdn.microsoft.com/en-us/library/windows/desktop/dd561977(v=vs.85).aspx)

- ▶ **PDF プロパティ群**は、PDF 文書に特有です。これらは pCOS パスから内容を与えられます。たとえば PDF バージョン番号・しおり内容・ページサイズです。
- ▶ **文書 XMP プロパティ群**は、XMP 2005 仕様中のすべての定義済みの文書プロパティを網羅しています。同仕様については下記を参照してください：

www.adobe.com/devnet/xmp

- ▶ **画像 XMP プロパティ群**は、文書内の画像群に添付されている XMP プロパティを網羅しています。画像プロパティ群も XMP 2005 仕様に由来します。
- ▶ **内部プロパティ群**は、業務用途ではなく開発・デバッグ支援として意図されたものです。これらは、ソフトウェアバージョン番号やインデックス作成時刻を含んでいます。

定義済みプロパティ群をその GUID・ユーザーフレンドリ名とともに示した全一覧が付章 A「定義済みメタデータプロパティ」にあります。

これらのコレクションの中のプロパティ群は、個別に構成する必要はありません。ただし、すべてのコレクションがデフォルトで有効になっているわけではないので（後述）、必要に応じて有効にする必要があります。定義済みのプロパティを、対応する GUID+ID 組み合わせをカスタムプロパティに割り当てることによって上書きすることも可能です。

プロパティセットコレクションを有効にするための XML 構成 定義済みプロパティセットコレクションの一部ないし全部を、*Metadata/PropertySetCollection* 要素内の対応する属性によって有効にすることができます：

```
<PropertySetCollection
  shell="true"
  pdf="true"
  documentXmp="true"
  imageXmp="true"
  internal="true"/>
```

デフォルトでは、*Shell*・*Internal* プロパティセットコレクションは有効になっており、*documentXMP*・*imageXMP*・*PDF* プロパティセットコレクションは無効になっています。カスタムプロパティは、1 個ないし複数のカスタムプロパティが指定された場合、暗黙的に有効になります。

あらゆる定義済みとカスタムのメタデータプロパティの処理を、*Filtering* 要素の *metadataHandling* 属性で完全に無効にすることも可能です：

```
<Filtering indexNestedPdf="true" metadataHandling="ignore">
```

3.4 カスタムメタデータプロパティ

カスタムメタデータプロパティは、定義済みプロパティのほかに追加するプロパティであり、企業・組織・業界などの個別の要請に応えるためのものです。TET PDF IFilter は、カスタムプロパティに対する完全制御を与えます。それらは、TET PDF IFilter によって生成されて検索エンジンによってインデックスされるように構成ファイル内で指定することができます。

カスタムデータプロパティを計画 カスタムプロパティを指定するためには、下記の諸側面を考慮する必要があります (GUID・識別子・フレンドリ名について詳しくは 40 ページの「プロパティ識別と GUID」を参照してください) :

- ▶ 1 個ないし複数のプロパティをプロパティセット内にグループ化することができます。各プロパティセットは、GUID という一意な 128 ビット識別子を必要とします。
- ▶ プロパティ識別子は、そのプロパティをそのプロパティセット内で識別する一意な整数です。セット内でプロパティ識別子は値 2 から始まります。IFilter クライアントによっては、この識別子をフレンドリ名へ置き換えることも可能です。定義済みのプロパティ群を、対応する GUID+ID 組み合わせを割り当てることによって上書きすることも可能です。
- ▶ プロパティに対するフレンドリ名は、識別子が利用可能な場合にはオプションであり、そうでない場合には必須です。これは任意の名前をとることができますが、構成ファイル内で一意な名前であればなりません。IFilter クライアントによっては、これは識別子のかわりに用いることができますが、フレンドリ名はすべての IFilter クライアントで使えるわけではありません。
- ▶ プロパティ源 : プロパティは、37 ページの 3.1 「PDF 内のメタデータ源」に従って、文書メタデータから、あるいは一般的 PDF 情報から、内容を与えられることができます。
- ▶ プロパティのデータ型 : *Int32* (32 ビット整数) ・ *Double* (倍精度の浮動小数点数) ・ *Boolean* (true/false) ・ *DateTime* (後述) ・ *String*。
- ▶ 優先権規則 : そのプロパティに対して複数のデータ源があるときは、最初に得られる空でないデータ源が優先権を有する (すなわち、それより後の源群は無視される) か、あるいはすべての空でない源からのデータが収集されるかを、指定することが可能です。
- ▶ そのプロパティがベクトルとして出力されるかどうか、すなわち、単値ではなく複数の値が IFilter インタフェースへ配列構造で受け渡されるかどうかを指定します (44 ページの 3.5 「多値プロパティ」を参照)。
- ▶ プロパティがフルテキストの一部としてインデックスされる場合にそのプロパティ名の頭に付加される接頭辞 (45 ページの 3.6 「メタデータプロパティをテキストとしてインデックス」を参照)。

DateTime プロパティ型 時点を指定するには、プロパティに対して *DateTime* データ型を用いることができます。*DateTime* プロパティ群について生成される出力はつねに、IFilter インタフェースによって必要とされる形式になっていますが、入力は、そのプロパティの源に応じてさまざまな形式に対応しています :

- ▶ データ源が pCOS パスの場合、たとえば標準文書情報フィールドや、*/Info/CustomDate* といったカスタム文書情報フィールドの場合には、その値は、ISO 32000-1 の 7.9.7 項で仕様定義されている標準 PDF 日付形式であると見なされます。PDF 日付文字列の一般形式は *D:YYYYMMDDHHmmSSOHH'mm* です。なお、PDF 日付形式は (後述の XMP 形式とは異なり)、秒の小数には対応していません。いくつか例を挙げます :

D:201109231858	GMT
D:201109231058-08'00'	同じ時点を米国太平洋標準時で
D:201109231958+01'00'	同じ時点を中央欧州標準時で

- ▶ データ源が XMP プロパティの場合 (例: `xmp:ModifyDate`) には、その値は、XMP 2005 リファレンスで基本値型 `Date` について仕様定義されている形式であると見なされます。これは、ISO 8601¹ で仕様定義されている形式と同等です。XMP 日付文字列の一般形式は `YYYY-MM-DDThh:mm:ss.sTZD` であり、ここで `TZD` はタイムゾーン指定子 (`Z` か `+hh:mm` か `-hh:mm`) です。なお、いくつかの部分はオプションです: XMP 日付は精度に応じて 6 段階の粒度をサポートしており、一方、PDF 日付形式はただ 1 段階の粒度しかサポートしていません。いくつか例を挙げます:

2011-09-23T18:58:30Z	GMT
2011-09-23T13:58:30-05:00	同じ時点を米国東部標準時で
2011-09-23T19:58:30+01:00	同じ時点を中央欧州標準時で

カスタムプロパティのための XML 構成 1 個ないし複数のカスタムプロパティを、`PropertySet` 要素内で指定することができます。ここでは各 `Property` 要素が、そのセット内の各プロパティを記述します:

```
<PropertySet guid="33333333-5354-4c72-992C-5DF2AA4E7CBA">
  <Property friendlyName="MailTo" identifier="2" type="String">
    <Source xmpName="acme:mailto"/>
  </Property>
</PropertySet>
```

複数の PDF 源を、同一の Windows プロパティへマップすることが可能です。`Property` 要素の存在は自動的に、その指定されたプロパティに対する処理を有効にします。ただし、あらゆる定義済みとカスタムのメタデータプロパティの処理を、`Filtering` 要素の `metadataHandling` 属性で完全に無効にすることも可能です:

```
<Filtering metadataHandling="ignore">
```

プロパティの GUID+ 名前処理のための XML 構成 TET PDF IFilter は、プロパティの識別子が利用可能な場合には、GUID+ID を用いて IFilter インタフェース内でプロパティを識別します。`identifier` 属性を持たない、`friendlyName` 属性のみを持つカスタムプロパティについては、かわりに GUID+ 名前によって識別されます。定義済みプロパティについても GUID+ 名前処理を有効にする (そしてグローバルに GUID+ 名前処理を強制する) ためには、`Filtering` 要素の `useIdentifier` 属性を用いることができます:

```
<Filtering useIdentifier="false">
```

定義済みプロパティに対して用いられる名前は、表 3.2 で示すプロパティ接頭辞から頭の `TET_` と末尾のアンダースコア「`」`を削ったものになります (例: `System_Author.pdfversion.dc_contributor.photoshop_DateCreated`)。

環境によっては、とりわけ SharePoint では、プロパティに対して GUID+ 名前処理を使うことは、GUID+ID よりも便利です。

1. www.w3.org/TR/NOTE-datetime を参照。

3.5 多値プロパティ

メタデータプロパティは、1個ないし複数の値を内容とすることができます。単値プロパティは、その文書を全体として記述する1個の単値から成ります。単値プロパティの例としては、作成日（プロパティ源：`xmp:CreateDate`・`/Info/CreationDate`）や一意な文書識別子（`dc:identifier`）があります。

多値プロパティは、文書あたり複数回出現する可能性があります。多値プロパティの例としては、文書作成者群や文書キーワード群のリストがあります。プロパティが複数値を持つにはいくつかの理由がありえます：

- ▶ そのプロパティの源が XMP コンテナ型であり、それゆえ、1度に複数の項目を保持している。例：`dc:creator` は XMP 内で型 `Seq` を持ちます。
- ▶ そのプロパティが、ワイルドカードを持つ多値な pCOS パスから内容を与えられる。ここでワイルドカードは、任意の数の個別の項目へ展開されます。例：`bookmarks[*]/Title`。
- ▶ そのプロパティが、複数の源から内容を与えられ、かつ、そのプロパティの `precedence` 属性が値 `try-all` を持っている。例：`pdf:Keywords`・`/Info/Keywords`。ただし、デフォルトの `precedence=first-wins` は、最初の空でないプロパティ源のみを処理します。

プロパティのベクトル処理 デフォルトでは、TET PDF IFilter はプロパティ定義内のすべての関連する源を処理し（`precedence` 属性に従って）、そして利用可能なかぎりの空でないプロパティ値を出力します。いいかえれば、それぞれの値は IFilter クライアントに対して1個の値として返されます。複数のプロパティ値が不特定の順序で返されます。

あるいは、多値のプロパティを IFilter クライアントに対してベクトルとして与えることも可能です。これは、1個ないし複数の値を保持できる1個の配列実体が出力されることを意味します。プロパティに対するベクトル処理の関連する側面がいくつかあります：

- ▶ SharePoint は、多値プロパティに、それがベクトル実体として処理される場合にのみ対応しています。
- ▶ IFilter クライアントによっては、たとえば Windows Search など、多値ベクトルプロパティ内の1個ないし複数の値を1回のクエリで検索できるベクトルクエリに対応しているものもあります。
- ▶ インデックスサービスは、型 `String` を持つプロパティに対してのみベクトルに対応しています。

2つの別個の概念があることに留意してください：**多値**とは、そのプロパティの源のありようを表していますが、**ベクトル処理**とは、プロパティ値群が IFilter クライアントへ転送される方式を表しています。ベクトル処理は、多値プロパティに対して、それが1個の値のみを内容とする場合であっても、適用することができます。

定義済みプロパティのうちのいくつかは多値です（付章 A「定義済みメタデータプロパティ」を参照）。

ベクトルプロパティのための XML 構成 カスタムプロパティに対するベクトル処理は、`Property` 要素の `emitAsVector` 属性で有効にすることができます：

```
<Property friendlyName="MailTo" type="String" precedence="try-all" emitAsVector="true">
  <Source xmpName="acme:mailto"/>
  <Source xmpName="gov:mailto"/>
</Property>
```

3.6 メタデータプロパティをテキストとしてインデックス

多くのテキスト検索エンジンは、何らかの方式でのプロパティクエリに対応しています。しかし、プロパティに対するクエリは、SQL Serverのような、フルテキスト検索にのみ対応している検索製品では、可能でない場合もあります。あるいは別のシナリオとして、明示的にプロパティを検索することは、もしもその検索語のあらゆる出現を、それが文書内容の中に出現するか、それとも何らかのプロパティの中に出現するかにかかわらず検索したい場合には、望ましくないかもしれません。いずれの場合においても、TET PDF IFilter に対して、すべてのプロパティをフルテキストインデックスの中へ入れ込むよう指示することが可能です。プロパティを文書内容本体から区別するために、TET PDF IFilter は、文字列をプロパティ値の頭に付加することも可能です。これによって、検索エンジンがプロパティ検索に直接は対応していない場合においても、プロパティをプロパティであると見分けることが容易になります。付章 A「定義済みメタデータプロパティ」に挙げる定義済みプロパティに対しては、固定された接頭辞が用いられます。

プロパティをテキストとしてインデックスすれば、もしそれをしなかったらプロパティ検索に一切対応していない環境において、限定的なプロパティ検索が可能になる一方で、そのプロパティ検索は限定的なものであるという事実は知っておくべきです。たとえば、論理型などの表現がプロパティ値として得られません。

メタデータをテキストとしてインデックスするための XML 構成 メタデータプロパティをテキストとしてインデックスするには、*Filtering* 要素の *metadataHandling* 属性を *propertyAndText* (プロパティをメインテキスト内へ透明な形で混ぜ込みたい場合) か *propertyAndPrefixedText* (プロパティを接頭辞で識別したい場合) に設定します：

```
<Filtering metadataHandling="PropertyAndPrefixedText">
```

文書をフィルタする際にカスタムプロパティの頭に付加されるオプションな接頭辞は、カスタムプロパティに対して *Property* 要素の *textIndexPrefix* 属性で指定することができます：

```
<Property friendlyName="Title" identifier="7" textIndexPrefix="TITLE_">
...
</Property>
```

定義済みプロパティの頭に付加される接頭辞は、下記の方式に従って構築されます：

```
TET_<プロパティ名>_
```

ここでプロパティ名の中のピリオドキャラクタ「.」はアンダースコアキャラクタ「_」へ置き換えられます (表 3.2 の例を参照してください)。

表 3.2 メタデータプロパティをテキストとしてインデックスするための接頭辞

プロパティセット コレクション	サンプルプロパティ名	メタデータをテキストとして インデックスするための接頭辞
Shell	System.Author	TET_System_Author_
TET	PDFlib.TETPDFIFilter.pdfversion	TET_pdfversion_
XMP	dc:contributor	TET_dc_contributor_
Image	photoshop:DateCreated	TET_photoshop_DateCreated_

シナリオ 1：メタデータプロパティをメインテキスト内へ透明な形で混ぜ込む もしメタデータプロパティが、ターゲット文書（群）を特定する十分に弁別的なテキストを内容として持っているならば、プロパティをフルテキストインデックスへ入れ込んで、それを標準的なフルテキストクエリへ入れ込めば充分でしょう。たとえば、特定の記事番号をクエリする場合には、その番号がその文書のメインテキスト内に出現しようが、それともメタデータプロパティ内に出現しようが、探しているその記事番号について語っている文書を 1 個に絞り込めるかぎりは、どちらでもかまわないのです。いいかえれば、テキストがメインテキスト内に出現しようが、それとも何らかのメタデータプロパティ内に出現しようが、どちらでもかまわない場合には、単にプロパティをフルテキストとしてインデックスすることを有効にする必要があります、それを越える手順は必要ありません。

メタデータをメインテキスト内へ透明な形で混ぜ込むには、下記の XML 構成を用います：

```
<Filtering metadataHandling="propertyAndText">
```

シナリオ 2：メタデータをメインテキストから区別 別の場合においては、テキストがメイン文書内に出現しているのか、それとも何らかのメタデータプロパティ内に出現しているのかに、意味がある場合もあるでしょう。たとえば、作成者が *Doyle* である文書を検索するか、それともメインテキスト内に語 *Doyle* を含む文書を検索するかは、互いに大きな違いです。このシナリオにおいては、プロパティをフルテキストとしてインデックスすることを有効にするだけではなく、メイン文書内容中のテキストとメタデータプロパティ内のテキストとを区別できるよう、各プロパティに対して適切な接頭辞をも入れ込む必要があります。

定義済みプロパティ *System.Author* の値の頭には、接頭辞 *TET_System_Author_* が付加されます。たとえば、*System.Author=Doyle* に対するプロパティベースの検索を、*TET_System_Author_Doyle* に対するフルテキスト検索でエミュレートすることができます。*System.Author* は定義済みプロパティですので、これに対応する XML 構成はプロパティ独自の項目を一切必要とせず、単にプロパティを接頭辞付きテキストとしてインデックスすることを有効にする必要があります：

```
<Filtering metadataHandling="propertyAndPrefixedText">
```

記事番号 *XY123456* を持つ文書に対するプロパティベースの検索を、*ArticleNumber_XY123456* に対するフルテキスト検索でエミュレートするには、下記の XML 構成を用います：

```
<Filtering metadataHandling="propertyAndPrefixedText">
```

```
<PropertySet guid="404e8a40-2e85-11dd-97f6-0002a5d5c51b">
  <Property identifier="2" textIndexPrefix="ArticleNumber_">
    <Source pdfObject="/Info/ArticleNumber"/>
  </Property>
</PropertySet>
```

3.7 メタデータを採りページ内容を無視

場合によっては、ページ内容ではなくメタデータプロパティ群のみに基づいて検索を行うほうが望ましいこともあります。すなわち、インデックス作成を行う過程でページ内容を完全に無視して、メタデータプロパティ群にのみ依拠するということです。こうしたシナリオとして考えられるのは、下記のうちの1つないし複数です：

- ▶ ユーザーが発するクエリへの制御を強化：メタデータプロパティ群に基づいて正確な結果が得られるときに、長い結果リストを調べあげるのは無駄です。
- ▶ 検索したい文書群がだいたい同一の語群を内容としているが、ただしそれらが異なった組み合わせになっている。たとえば請求書などの取引文書。
- ▶ ページ内容本体があらかじめよくわかっているので、検索上実際の意味を持たない。しかし、特定の文書に対して集められたページ群が関心対象である。たとえば、可変の数の契約と紐付けられた保険取引では、契約の正確な文言ではなく、契約文書の数と種類が検索上意味を持つ。
- ▶ 検索したい文書群が検索可能なテキストを全く内容としていない。たとえばスキャンされた、OCRを全く実行されていない文書。
- ▶ インデックスに何ら意味のあるありようで貢献しない情報を内容とする文書群。たとえば数のみを内容とする長い財務文書や、テキストのない（ないし図中のキャプション内のテキストのみの）技術図面。
- ▶ 上述のいずれかの場合におけるパフォーマンス最適化：内容が検索の助けにならないとわかっているなら、文書をインデックスするのは無駄です。

ページ内容のインデックス作成を無効にするためのXML構成 ページ内容のインデックス作成を完全に無効にするには、*Filtering* 要素の *indexPageContents* 属性を *false* に設定します：

```
<Filtering indexPageContents="false" metadataHandling="property">
```


4 各種 IFilter クライアントでメタデータ処理

4.1 Windows Search でメタデータ

プロパティ記述ファイルを作成 Windows Search は、定義済みとカスタムのメタデータプロパティ群の記述を保持している Windows プロパティシステムにアクセスします。Windows Search でカスタムメタデータプロパティを使うには、プロパティ群の名前・データ型・GUID などのプロパティ属性を指定するプロパティ記述ファイル (*.propdesc* という名前のことが多い) を用意する必要があります。このプロパティ記述は、XML 構成ファイル内の対応するプロパティ記述と整合している必要があります。プロパティ記述は、下記で記述されている文法に従った XML ファイルとして指定される必要があります：

[msdn.microsoft.com/en-us/library/bb773879\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/bb773879(VS.85).aspx)

サンプルプロパティ記述ファイル群が、TET PDF IFilter とともにインストールされています。下記の (不完全な) 例で、いくつかのプロパティ記述を演示します：

```
<propertyDescription name="PDFlib.TETPDFIFilter.fontcount"
  formatID="{5eac0060-1ba4-11dd-92c4-0002a5d5c51b}" propID="2">
  <typeInfo type="Int32" isInnate="true" isVisible="true" isQueryable="true"/>
  <labelInfo label="fontcount" sortDescription="LowestHighest"/>
  <searchInfo inInvertedIndex="true" isColumn="true" columnIndexType="OnDisk"/>
</propertyDescription>
<propertyDescription name="PDFlib.TETPDFIFilter.weblink"
  formatID="{5eac0060-1ba4-11dd-92c4-0002a5d5c51b}" propID="6">
  <typeInfo type="String" isInnate="true" multipleValues="true" isVisible="true"
    isQueryable="true"/>
  <labelInfo label="weblink" sortDescription="AToZ"/>
  <searchInfo inInvertedIndex="true" isColumn="true" columnIndexType="OnDisk"/>
</propertyDescription>
```

多値プロパティに対しては *multipleValues="true"* 属性が必須であることに留意してください (44 ページの 3.5 「多値プロパティ」を参照)。

定義済みプロパティ すべての定義済みプロパティ (付章 A 「定義済みメタデータプロパティ」を参照) のためのプロパティ記述ファイル *predefined_properties.propdesc* が、TET PDF IFilter とともにインストールされています。そのプロパティ定義は TET PDF IFilter に内蔵されていますが、それをクエリで使えるようにするにはまず、使いたいプロパティセットコレクション群を XML 構成ファイル内で有効にして (41 ページの 3.3 「定義済みメタデータプロパティ」を参照)、かつ、プロパティ記述を登録する (49 ページの「メタデータプロパティ群を Windows に登録」を参照) 必要があります。

Windows Search のための XML 構成 表 4.1 に、Windows Search を使用する際の TET PDF IFilter のための XML 構成に関連する要請・推奨事項を挙げます。

メタデータプロパティ群を Windows に登録 Windows Search でプロパティクエリを使うには、対応するプロパティ記述を Windows プロパティシステムに登録する必要があります。プロパティ記述を登録するには、TET PDF IFilter とともにインストールされているコ

表 4.1 Windows Search のための XML 構成

要素	属性	要請・推奨事項
<i>Filtering</i>	<i>uselidentifier</i>	true にする必要があります。なぜなら、Windows Search はプロパティ識別のためには GUID+ID 方式にのみ対応しており、GUID+ 名前には対応していないからです（43 ページの「プロパティの GUID+ 名前処理のための XML 構成」を参照）。
<i>Property</i>	<i>identifier</i>	必須。なぜなら、Windows Search はプロパティ識別のためには GUID+ID 方式にのみ対応しており、GUID+ 名前には対応していないからです（43 ページの「プロパティの GUID+ 名前処理のための XML 構成」を参照）。
<i>PropertySetCollection</i>	<i>shell</i>	Windows Search に知られているシェルプロパティ群に対応するには true にする必要があります（なお、true がもともとデフォルトです）。

マンドラインプログラム *registerpropdesc.exe* を使います。ただしこのツールは、Windows Search がシステムにインストールされている場合にのみ動作します。プロパティ記述ファイルの名前をこのツールに与えると、それが Windows プロパティシステムで登録されます（パスが空白キャラクタを含む場合には、必ずパスをダブルクォーテーションでくくってください）：

```
registerpropdesc "predefined_properties.propdesc"
```

注記 6 ページの「Windows Vista/7/8 で特権コマンドを実行」の項を、もしこれらのシステムを使用している場合には読んでください。

Windows プロパティシステムは、プロパティ記述ファイルの名前を、下記のレジストリキー群（および、このキーの末尾要素の数を増やしたもの）に格納しています：

```
HKLM\SOFTWARE\Microsoft\Windows\CurrentVersion\PropertySystem\PropertySchema\0000
```

この *registerpropdesc* ツールは、プロパティファイルがうまく登録できなかつたとき（たとえばプロパティの重複を検出した場合など）には、エラーメッセージと *HRESULT* エラー値を出力します。エラーが起きたとき、詳細を調べるには、アプリケーションイベントログをチェックします：

- ▶ 「スタート」→「設定」→「コントロールパネル」→「管理ツール」→「イベント ビューア」
- ▶ 左のペーンで「アプリケーション」をクリック。
- ▶ プロパティ登録で問題があった場合には、ソース *Microsoft-Windows-propsys* の項目があります。その項目を内容としている行をダブルクリックして、エラーメッセージを調べます（例：*Omitted duplicate property*）。

あるいは、*registerpropdesc.exe* が出力する *HRESULT* 値を用いて問題を分析することも可能です。*HRESULT* 値と説明の一覧は下記にあります：

msdn.microsoft.com/en-us/library/cc231198.aspx

prop.exe という共有ソースのコマンドラインツールは、*registerpropdesc* と同様の機能を提供し、かつそれ以外にも Windows プロパティシステムを扱う機能群を提供します：

prop.codeplex.com

プロパティ記述を登録するための要請事項 Windows Search のためのプロパティ記述を登録する際には、下記の重要な要請事項に留意してください：

- ▶ カスタムプロパティを検索で利用可能にするためには、プロパティ記述を登録して、検索サービスを停止させて再起動し、Windows Search オプションで再構築を強制する必要があります：

```
registerpropdesc "predefined_properties.propdesc"  
net stop wsearch  
net start wsearch
```

...カタログを再構築（7ページの「Windowsデスクトップサーチを開始/停止」を参照）...

- ▶ 登録するプロパティは、既存の Windows プロパティ（付章 A「定義済みメタデータプロパティ」に挙げているシェルプロパティを含め）と GUID+ID または名前で衝突してはいけません。Windows プロパティの一覧は下記にあります：

[msdn.microsoft.com/en-us/library/windows/desktop/dd561977\(v=vs.85\).aspx](https://msdn.microsoft.com/en-us/library/windows/desktop/dd561977(v=vs.85).aspx)

- ▶ *registerpropdesc* ツールを走らせるには、*HKEY_LOCAL_MACHINE* レジストリハイブへの適切な書き込み権限が必要です。
- ▶ レジストリはプロパティ記述ファイルの名前を保持しているだけです（その内容は保持していないので）、*propdesc* ファイルはその登録された場所にありつづけている必要があります。そうでないとプロパティ記述が検索で利用可能でなくなってしまう。
- ▶ *propdesc* ファイルは、すべてのユーザから読み込み可能である必要があります。
- ▶ 複数のプロパティファイルを登録することも可能です。ただし、2つのプロパティ記述が同一のプロパティ（GUID+ID によって特定される）の記述を含んではいけません。

-u オプションを用いると、以前に登録されているプロパティ記述を除去することができます：

```
registerpropdesc -u "c:\property\acme.propdesc"
```

メタデータプロパティを検索 カスタムメタデータプロパティを構成した後は、プロパティを検索できます。表 4.2 にプロパティクエリの例を挙げます。Windows Search のためのすべての定義済みプロパティの名前を付章 A「定義済みメタデータプロパティ」に挙げてあります。

シェルプロパティ群は、表 1.1 に挙げるサンプルに従ってインタラクティブにクエリすることも可能です。インタラクティブ検索のための Advanced Query Syntax (AQS) は、カスタムプロパティには Windows XP/Vista では対応しておらず、Windows 7/8 でのみ対応しています。Windows XP/Vista では、カスタムメタデータをインタラクティブ検索に含めるためには、メタデータプロパティをテキストとしてインデックスする（45 ページの 3.6「メタデータプロパティをテキストとしてインデックス」を参照）か、あるいは SQL ベースの検索を使う（下記参照）必要があります。

表 4.2 Windows Search 3.0 以上のためのメタデータクエリの例

検索語例	説明
author:Doyle	作成者が Doyle を含んでいる
author:"Conan Doyle"	作成者が語列 Conan Doyle を含んでいる
author:Doy	作成者が Doy で始まる

表 4.2 Windows Search 3.0 以上のためのメタデータクエリの例

検索語例	説明
size: < 500000	文書サイズが 500,000 バイト未満
date: <=4/7/12	更新日が 2012 年 4 月 7 日以前（日付形式はシステム設定に依 存します）
System.Document.DateCreated: = 09/23/2002	作成日が 2002 年 9 月 23 日（日付形式はシステム設定に依 存します）
System.MIMETYPE: ="application/pdf"	すべての PDF 文書をインデックス内にリスト
System.Document.PageCount: = 144	144 ページある文書
PDFlib.TETPDFIFilter.bookmark:Index	（Windows XP/Vista では動作しません）文書がテキスト Index を持つしおりを含んでいる

メタデータプロパティに対する SQL クエリ SQL クエリは、定義済みだけでなくカスタムのクエリも検索できます。いくつかの例を以下に示します。これらの例は、XML 構成ファイル内ですべての定義済みプロパティのインデックス作成が有効にされていること、および、*predefined_properties.propdesc* プロパティ記述が登録されていることを前提にしています。ここでは、ADO (*ActiveX Data Objects*) と PowerShell スクリプトを用いて SQL ベースのクエリを送信します。しかし、これ以外の ADO または ADO.NET 環境でも SQL ステートメントを使うことができます。Windows Search のための SQL 文法拡張の説明は下記にあります：

[msdn.microsoft.com/en-us/library/bb231256\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/bb231256(VS.85).aspx)

プロパティ検索は 2 方向に適用できます：

- ▶ 特定のプロパティ値をクエリ：どの文書が *Doyle* を作成者として持っているか？
- ▶ 1 個ないし複数のファイルの中で特定のプロパティの値をクエリ：この文書の作成者は誰か？

PowerShell で SQL クエリを送信 Windows PowerShell は、Windows 7/8・Windows Server 2008 に内蔵されており、Windows XP SP2・Windows Server 2003・Windows Vista にインストール可能です。サンプル PowerShell クエリスクリプト群が、TET PDF IFilter とともにインストールされています。下記の PowerShell スクリプトは、すべての文書についてすべてのしおりをリストします：

```
$objConnection = New-Object -comobject ADODB.Connection
$objRecordset = New-Object -comobject ADODB.Recordset
$objConnection.Open("Provider=Search.CollatorDSO;Extended
Properties='Application=Windows';")

$objRecordSet.Open(
"SELECT System.ItemPathDisplay, `\"PDFlib.TETPDFIFilter.bookmark`\" FROM SYSTEMINDEX ",
$objConnection)

While ($objRecordset.EOF -ne $True) {
    $private:item = $objRecordset.Fields.Item("System.ItemPathDisplay")
    Write-Output $item.Value
    $item = $objRecordset.Fields.Item("PDFlib.TETPDFIFilter.bookmark")
    Write-Output $item.Value
    $objRecordset.MoveNext()
}
```

下記の PowerShell スクリプトは、少なくとも 1 個のしおりがテキスト *alpha* を含んでいる文書をすべてリストします：

```
$objConnection = New-Object -comobject ADODB.Connection
$objRecordset = New-Object -comobject ADODB.Recordset
$objConnection.Open("Provider=Search.CollatorDSO;Extended
Properties='Application=Windows';")

$objRecordset.Open("SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE " +
    "`"PDFlib.TETPDFIFilter.bookmark`" = SOME ARRAY ['alpha']", $objConnection)

While ($objRecordset.EOF -ne $True) {
    $private:item = $objRecordset.Fields.Item("System.ItemPathDisplay")
    Write-Output $item.Value
    $objRecordset.MoveNext()
}
```

VBScript で SQL クエリを送信 下記の VBScript コードは、すべての文書を、その文書を作成したアプリケーションの名前とともにリストします。残念ながら、VBScript クエリで使うことができるのはシェルプロパティの一部だけです。他のプロパティはクエリできません：

```
On Error Resume Next

Set objConnection = CreateObject("ADODB.Connection")
Set objRecordSet = CreateObject("ADODB.Recordset")

objConnection.Open "Provider=Search.CollatorDSO;Extended
Properties='Application=Windows';"

objRecordSet.Open "SELECT System.ItemPathDisplay, System.ApplicationName FROM
SYSTEMINDEX", objConnection

objRecordSet.MoveFirst

Do Until objRecordset.EOF
    Wscript.Echo objRecordset.Fields.Item("System.ItemPathDisplay")
    Wscript.Echo objRecordset.Fields.Item("System.ApplicationName")
    Wscript.Echo ""
    objRecordset.MoveNext
Loop
```

SQL で複雑なプロパティクエリ 下記のサンプル群は、関連する SQL ステートメントのみを内容としており、あらゆる SQL 環境で使用可能です。PowerShell スクリプトでこれらのステートメントを使うには、適切なクォーティングを適用する必要があります。たとえば `"PDFlib.TETPDFIFilter.width"` ではなく `"PDFlib.TETPDFIFilter.width"` としなければなりません。下記の多くの例は、ベクトルプロパティに対して配列クエリを用いています (44 ページの 3.5 「多値プロパティ」を参照)。配列クエリのための文法の詳細は下記にあります：

[msdn.microsoft.com/en-us/library/bb231264\(VS.85\).aspx](https://msdn.microsoft.com/en-us/library/bb231264(VS.85).aspx)

- ▶ 作成者が *Doyle* を含んでいる文書をすべてリスト：

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE CONTAINS("System.Author",
'Doyle')
```

- ▶ 作成者が *Rudy* で始まる文書をすべてリスト：

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE CONTAINS("System.Author",
'"Rudy*"')
```

- ▶ PDF/A-1a に準拠している文書をすべてリスト :

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE "PDFlib.TETPDFIFilter.pdfa" =
'PDF/A-1:2005'
```

- ▶ 少なくとも 1 個のしおりが *alph* で始まる文書をすべてリスト :

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE "PDFlib.TETPDFIFilter.bookmark"
LIKE SOME ARRAY ['alph%']
```

- ▶ フォント *Bembo* と *TimesNewRoman* の少なくとも 1 つを含んでいる文書をすべてリスト :

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE "PDFlib.TETPDFIFilter.font" =
SOME ARRAY ['Bembo', 'TimesNewRoman']
```

- ▶ *Bembo* フォントと *Bembo-Bold* フォントを両方含んでいる文書をすべてリスト :

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE "PDFlib.TETPDFIFilter.font" =
SOME ARRAY ['Bembo'] AND "PDFlib.TETPDFIFilter.font" = SOME ARRAY ['Bembo-Bold']
```

- ▶ 少なくとも 1 つのページが *width*=595 である文書をすべてリスト。595 をシングルクォーテーションキャラクタでくくっていることに留意してください。これは、*width* が型 *Double* だからです。型 *Int32* に対してはこのクォーテーションは必要ありません。

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE "PDFlib.TETPDFIFilter.width" =
SOME ARRAY ['595']
```

- ▶ 少なくとも 1 つのページが *width*=200 であり、かつ少なくとも 1 つのページが *height*=150 である文書をすべてリスト :

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE "PDFlib.TETPDFIFilter.width" =
SOME ARRAY ['200'] AND "PDFlib.TETPDFIFilter.height" = SOME ARRAY ['150']
```

- ▶ 少なくとも 1 個のテニス画像 (Photoshop カテゴリ *TEN* = テニス) を持つ文書をすべてリスト :

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE
"PDFlib.TETPDFIFilter.images.photoshop.SupplementalCategories" = SOME ARRAY ['TEN']
```

- ▶ PDF バージョンが 1.6 より高い文書をすべてリスト (PDF バージョンは TET PDF IFilter によって番号 × 10 を内容とする文字列として返されます。例: 16 なら PDF 1.6) :

```
SELECT System.ItemPathDisplay FROM SYSTEMINDEX WHERE "PDFlib.TETPDFIFilter.pdfversion"
> '16'
```

4.2 SharePoint・Search Server でメタデータ

Search Server におけるメタデータプロパティ処理は、SharePoint におけるメタデータ処理とほぼ同等です。この節では、この両方の製品のための構成手順を解説します。いくつかの場合においては、製品ごとの詳細に明示的に言及します。

SharePoint におけるメタデータ処理は、Enterprise Search Administration 名前空間 **Microsoft.Office.Server.Search.Administration** を用いてプログラマ的に制御できます。SharePoint におけるメタデータプロパティの処理については下記の諸概念が用いられます：

- ▶ **クロール対象プロパティ**は、TET PDF IFilter によって、PDF 文書をインデックス（クロール）する際に作成されます。クロール対象プロパティは複数の値を持つことができます。それは TET PDF IFilter 構成ファイルを通じて制御されます。Sharepoint は、クロール対象プロパティを他のフィールドと同様に検索しますが、プロパティクエリのための高度な機能群は利用できません。
- ▶ **管理プロパティ**は、検索クエリと高度な検索で用いることができ、検索結果に表示されます。

「高度な検索」ページの管理プロパティに関する詳細は下記にあります：

msdn.microsoft.com/en-us/library/bb428648.aspx

msdn.microsoft.com/en-us/library/bb608302.aspx

TET PDF IFilter を SharePoint とともに使うための XML 構成 表 4.3 に、SharePoint を使用する際の TET PDF IFilter のための XML 構成に関連する要請・推奨事項を挙げます。

表 4.3 SharePoint のための XML 構成

要素	属性	要請・推奨事項
<i>Filtering</i>	<i>useIdentifier</i>	定義済みプロパティを含めすべてのプロパティの GUID+ 名前処理を強制するには false に設定します。
<i>Property</i>	<i>emitAsVector</i>	複数の値を持ちうるプロパティ群に対しては true にする必要があります（44 ページの 3.5 「多値プロパティ」を参照）。
<i>Property</i>	<i>friendlyName</i>	GUID+ 名前処理を使うと、リスト内でプロパティを見つけやすくなりますので、この属性の使用を推奨します（40 ページの「プロパティ識別と GUID」を参照）。

カスタムメタデータプロパティを構成 カスタムメタデータプロパティをインデックス作成のために用意するには下記のように操作します：

- ▶ フルクロールを走らせます。クロールの間、TET PDF IFilter は、XML 構成ファイル内で構成されている、かつ 1 個ないし複数の文書内で見つかったプロパティ群をレポートします。これは、SharePoint が新たに見つかったクロール対象プロパティ群をピックアップするために必要です。
- ▶ SharePoint は、すべての新しいプロパティカテゴリに対して合成名を生成します（例：「*Category 1*」・「*Category 2*」など）。TET PDF IFilter とともにインストールされている PowerShell スクリプトは、カテゴリ名とその GUID のリストを生成します。このスクリプトを下記のように走らせます：

```
list_categories.ps1 <サイトURL>
```

結果リスト出力の中で、カテゴリをその GUID で特定することができます。この GUID を、TET PDF IFilter のための XML 構成ファイル内の *PropertySet/@guid* のそれと（カスタムプロパティを作成した場合）、あるいは、付章 A「定義済みメタデータプロパティ」のプロパティセット GUID と（定義済みプロパティセットコレクションの中のプロパティを用いる場合）比較してください。

- ▶ SharePoint によって生成された合成カテゴリ名を、そのカテゴリを *list_categories.ps1* スクリプトによって返された GUID によって特定してから、ユーザーフレンドリな名前に変更：

SharePoint : 「スタート」 → 「SharePoint 3.0 サーバーの全体管理」をクリック。左側のコラム内で「共有サービス管理 : *SharedServices1* (ないしこのようなもの)」 → 「検索の設定」 → 「メタデータ プロパティのマッピング」をクリック。また左側のコラム内で「クロールされたプロパティ」をクリック。「クロールされたプロパティのビュー」内で、カテゴリのために表示されるドロップダウンメニューを開き、「カテゴリの編集」をクリック。「カテゴリの編集」ダイアログで、そのカテゴリの名前を変更。

- ▶ SharePoint : 「スタート」 → 「SharePoint 3.0 サーバーの全体管理」をクリック。左側のコラム内で「共有サービス管理 : *SharedServices1* (ないしこのようなもの)」 → 「検索の設定」 → 「メタデータ プロパティのマッピング」をクリック

Search Server : 「スタート」 → 「すべてのプログラム」 → 「Microsoft Search Server」 → 「Search Server 2008 Administration」をクリック。すると、「Metdata Property Mappings」と題されたページが開きます。

- ▶ 「New Managed Property」をクリック (図 4.1 参照)。
- ▶ プロパティの名前・説明・データ型を入力し、右下隅の、フィールド「Crawled properties mapped to this managed property」の右にあるボタン「Add Mapping」をクリック。

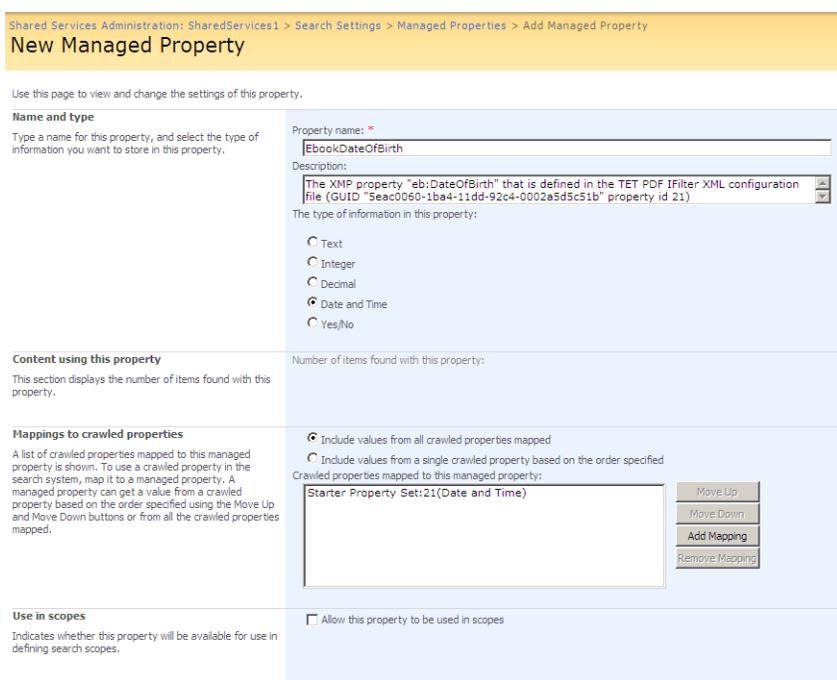


図 4.1
SharePoint で管理
プロパティを追加

- ▶ ドロップダウンメニュー「*Select a category*」は、「*Select a crawled property*」と題されたリストボックスに表示されているクローリング対象プロパティ群をフィルタします。このリストボックスは、プロパティと、そのカテゴリ名とプロパティ識別子またはフレンドリ名のリストを表示します。識別子かフレンドリ名かの選択は、TET PDF IFilter の XML 構成ファイル内のプロパティ定義に依存しますので、付章「カスタムメタデータプロパティを構成」を参照してください。新規にクローリングされたプロパティを選択し、その管理プロパティに名前を名前を割り当てて、それを保存します。

サンプル PowerShell スクリプト出力 *litwaredemo* サイトのための PowerShell スクリプト群を走らせるシーケンスを下記に挙げます：

```
PS C:\> & "C:\Program Files\PDFFlib\TET PDF IFilter 4.1\IFilter clients\Sharepoint\list_property_categories.ps1" http://litwaredemo
```

```
"Category 1" 007867f0-c59b-43fc-ab1e-8eee77057254
"Category 2" 1e3ee840-bc2b-476c-8237-2acd1a839b22
"Category 5" c60e822a-074f-4bd5-9889-6ebd372f2000
"Category 6" 17eb8447-fc9b-4d4d-81df-31e9aa770cbf
"Category 7" 5eac0060-1ba4-11dd-92c4-0002a5d5c51b
"Dublin Core" d92bb3ca-ce2b-4b9b-972a-5bf54b468171
...
```

生成されるリストは、「*Category 7*」が、*starter_advanced_search.xml* で定義されているプロパティセットに属していることを示しています。これを上述のようにして「*Starter Property Set*」に名前変更します。

管理プロパティへの検索のための SharePoint XML を作成 このステップでは、管理プロパティを「高度な検索」に追加するために必要な XML を作成します。次の項で、この XML を実際に適用する方法を説明します。以下、例とともにこの XML を説明します：

それぞれのプロパティについて、*PropertyDefs* 要素の子として *PropertyDef* 要素を作成する必要があります：

```
<PropertyDef Name="EbookDateOfBirth" DataType="datetime" DisplayName="Ebook Date of Birth"/>
```

Name 属性はプロパティ名と対応している必要があり、*DataType* 属性はそのプロパティの型を表 4.4 に従って記述し、*DisplayName* はユーザインタフェースで表示される任意の名前を内容とします。

新規の管理プロパティを、PDF 文書群に対する検索で利用可能にするには、そのプロパティを *ResultType* 要素に追加します：

```
<ResultType DisplayName="PDF Documents" Name="pdfdocuments">
  <Query>FileExtension='pdf'</Query>
  <PropertyRef Name="Author"/>
  <PropertyRef Name="Description"/>
  <PropertyRef Name="FileName"/>
  <PropertyRef Name="Size"/>
  <PropertyRef Name="Path"/>
  <PropertyRef Name="Created"/>
  <PropertyRef Name="Write"/>
  <PropertyRef Name="CreatedBy"/>
  <PropertyRef Name="ModifiedBy"/>
</ResultType>
```

```
<PropertyRef Name="EbookDateOfBirth"/>
</ResultType>
```

表 4.4 SharePoint のためのプロパティデータ型一覧

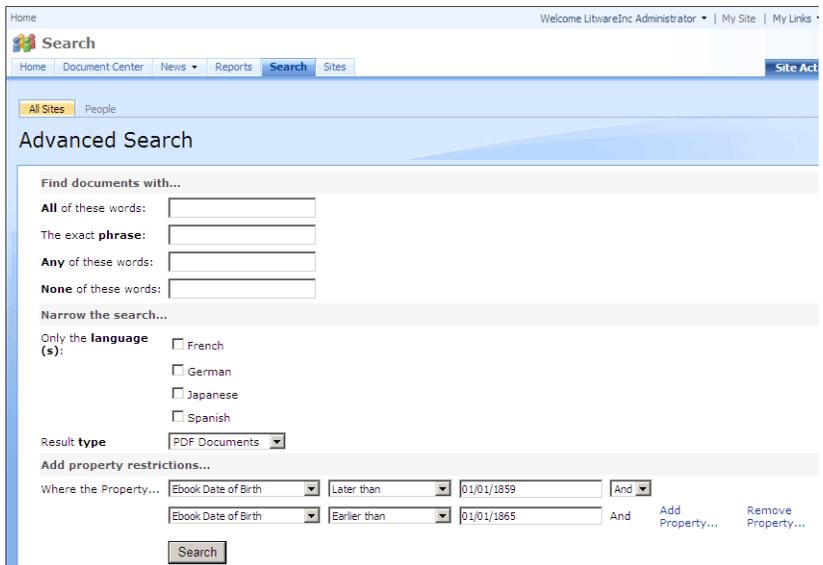
TET PDF IFilter におけるデータ型 SharePoint のためのデータ型	
Int32	integer
Double	decimal
Boolean	boolean
DateTime	datetime
String	text

スタータセット内のすべてのプロパティのための完全な XML ファイル (*starter_advanced_search.xml*) が、TET PDF IFilter とともにインストールされています。これは、SharePoint 評価のための Microsoft の Virtual PC イメージで作成されています。この XML を評価・検証のために使うことができます。業務サイト用の XML は、そのサイトに存在する XML と、新たに構成されたプロパティ群に基づいて、注意深く構築しなければならないことに留意してください。

カスタムメタデータプロパティを検索 高度なメタデータ検索を実装するためには、TET PDF IFilter を、プロパティをインデックスするよう構成する必要があります (42 ページの 3.4 「カスタムメタデータプロパティ」を参照)。これは TET PDF IFilter に対して、クロール対象プロパティを作成するよう指示します。その後、新規の管理プロパティを、下記のようにして「高度な検索」ページで利用可能にする必要があります：

- ▶ SharePoint サイトにログオンし、「高度な検索」へ移動。今したいことは、この検索フォーム上で、ページ下端の「Add property restrictions...」の下に新しい管理プロパティを追加することです。

図 4.2
SharePoint の
「高度な検索」ページで
カスタムプロパティ



4.3 SQL Server でメタデータ

SQL Server は、メタデータプロパティのインデックス作成や検索に対応していません。ゆえにデフォルトでは、あらゆるプロパティは無視されます。メタデータクエリを実装するには、メタデータプロパティ群をテキストとしてインデックスすることを推奨します (45 ページの 3.6 「メタデータプロパティをテキストとしてインデックス」を参照)。

SQL Server のための XML 構成 表 4.6 に、SQL Server を使用する際の TET PDF IFilter のための XML 構成に関連する要請・推奨事項を挙げます。

表 4.6 SQL Server のための XML 構成

要素	属性	要請・推奨事項
<i>Filtering</i>	<i>metadataHandling</i>	プロパティをテキストとしてインデックスすることを有効にするには、この属性を <i>propertyAndText</i> か <i>propertyAndPrefixedText</i> に設定します (45 ページの 3.6 「メタデータプロパティをテキストとしてインデックス」を参照)。
<i>Property</i>	<i>textIndexPrefix</i>	明示的にプロパティを検索したい場合は接頭辞を設定します。

メタデータプロパティを検索 SQL Server にはプロパティを検索するための専用機能は何もありませんので、フルテキストクエリの中でメタデータプロパティをクエリする必要があります。プロパティをテキストとしてインデックスすることを有効にした後は、作成者が *Arthur Conan Doyle* である文書群を下記のように検索できます：

```
SELECT name FROM DocumentTable WHERE CONTAINS(*,"TET_System_Author_Arthur Conan Doyle")
GO
```

作成者が *Arthur* で始まる文書群をクエリするには下記のステートメントを用います：

```
SELECT name FROM DocumentTable WHERE CONTAINS(*,"System_Author_Arthur*")
GO
```

4.4 インデックスサービスでメタデータ

列定義ファイルを作成 インデックスサービスでカスタムメタデータプロパティを使うには、プロパティ群の名前・データ型・GUID を指定する列定義ファイル (*msscolmn.txt* という名前のことが多い) を準備する必要があります。そのプロパティ項目は、XML 構成ファイル内の対応するプロパティ記述と整合している必要があります。列定義は、下記の文法に従って指定する必要があります：

```
"FriendlyName" (Type) = PropertySetGuid PropertySpec
```

それぞれの行は、下記の構成要素の全部ないし一部を内容とします：

- ▶ **FriendlyName** (ダブルクォーテーションマークでくくる必要があります)：そのプロパティのエイリアス名 (空白類文字を含んではいけません)：XML 構成ファイル内の **Property** 要素の **name** 属性と (または付章 A「定義済みメタデータプロパティ」に挙げている定義済みプロパティのいずれか 1 つの名前と) 整合する必要があります。
- ▶ **Type** (オプション。括弧でくくる必要があります)：表 4.7 に示す型指定を用いて、XML 構成ファイル内の **Property** 要素の **type** 属性と整合する必要があります。この項目がない場合は、デフォルト型 **DBTYPE_WSTR** | **DBTYPE_BYREF** (すなわち String) が用いられます。
- ▶ **PropertySetGuid**：XML構成ファイル内の**PropertySet**要素の**guid**属性と整合する必要があります。
- ▶ **PropertySpec**：XML構成ファイル内の**Property**要素の**friendlyname**属性か**identifier**属性と整合する必要があります。**friendlyname** は、**identifier** 属性がないか、あるいは **Filtering** 要素の **uselidentifier** 属性が **false** である場合にのみ使えます。識別子は 10 進または 16 進 (**ox** で始まる) 記法で指定できます。

空行と、番号記号「#」で始まる行は無視されますので、コメントに使うことができます。以下、いくつかのプロパティ記述の例を演示します (**[Names]** 項目を持つ先頭行は重要です)：

```
[Names]
"fontcount" (DBTYPE_I4) = 5eac0060-1ba4-11dd-92c4-0002a5d5c51b 2
"weblink" (DBTYPE_WSTR | DBTYPE_BYREF) = 5eac0060-1ba4-11dd-92c4-0002a5d5c51b 6
```

表 4.7 インデックスサービスのためのプロパティデータ型一覧

TET PDF IFilter におけるデータ型	msscolmn.txt 内のプロパティ記述におけるデータ型	インデックスサービス (プロパティのインタラクティブ編集) のためのデータ型
Int32	DBTYPE_I4	VT_I4
Double	DBTYPE_R8	VT_R8
Boolean	DBTYPE_BOOL	VT_BOOL
DateTime	DBTYPE_DATE	VT_DATE
String	DBTYPE_WSTR DBTYPE_BYREF ベクトル処理のプロパティの場合： DBTYPE_WSTR DBTYPE_VECTOR	VT_LPWSTR

注記 列定義ファイルは、64 ビットバージョンのインデックスサービスでは動作しないようです。

定義済みプロパティ すべての定義済みプロパティ（付章A「定義済みメタデータプロパティ」を参照）のための列定義ファイルが、TET PDF IFilter とともにインストールされています。これは、次項で説明する手順に従って構成することができます。プロパティ定義はTET PDF IFilterに内蔵されていますが、使いたいプロパティセットコレクションをXML構成ファイル内で有効にする必要があります（41 ページの 3.3「定義済みメタデータプロパティ」を参照）。このすべての定義済みプロパティのための列定義ファイルを使うことは、自分の文書群の中にどのプロパティ群が実際に見つかるかを見る簡便な手段です。なぜなら検出されたプロパティ群だけが管理コンソールに表示されるからです。

インデックスサービスのための XML 構成 表 4.8 に、インデックスサービスを使用する際の TET PDF IFilter のための XML 構成に関連する要請・推奨事項を挙げます。

表 4.8 インデックスサービスのための XML 構成

要素	属性	要請・推奨事項
<i>Property</i>	<i>emitAsVector</i>	多値を持ちうるプロパティ群に対しては、そのプロパティの <i>type</i> 属性が String 以外である場合には、false にする必要があります（44 ページの 3.5「多値プロパティ」を参照）。

カスタムメタデータプロパティを構成 カスタムメタデータプロパティをインデックス作成のために用意するには、以下のように操作します。いくつかのステップは、定義済みプロパティについては不要であり、カスタムプロパティについてのみ必要です：

- ▶ (定義済みプロパティについては不要) XML 構成ファイル内に、TET PDF IFilter のための構成項目群を作成（42 ページの 3.4「カスタムメタデータプロパティ」を参照）。
- ▶ (定義済みプロパティについては、TET PDF IFilter とともにインストールされている列定義ファイルを使用) 61 ページの「列定義ファイルを作成」の説明に従ったプロパティ記述を持つ列定義ファイルを作成。
- ▶ 「スタート」→「ファイルを指定して実行...」をクリックし、*regedt32.exe* とタイプしてレジストリエディタを起動。下記の項目へ移動：

HKEY_LOCAL_MACHINE\SYSTEM\CurrentControlSet\Control\ContentIndexCommon

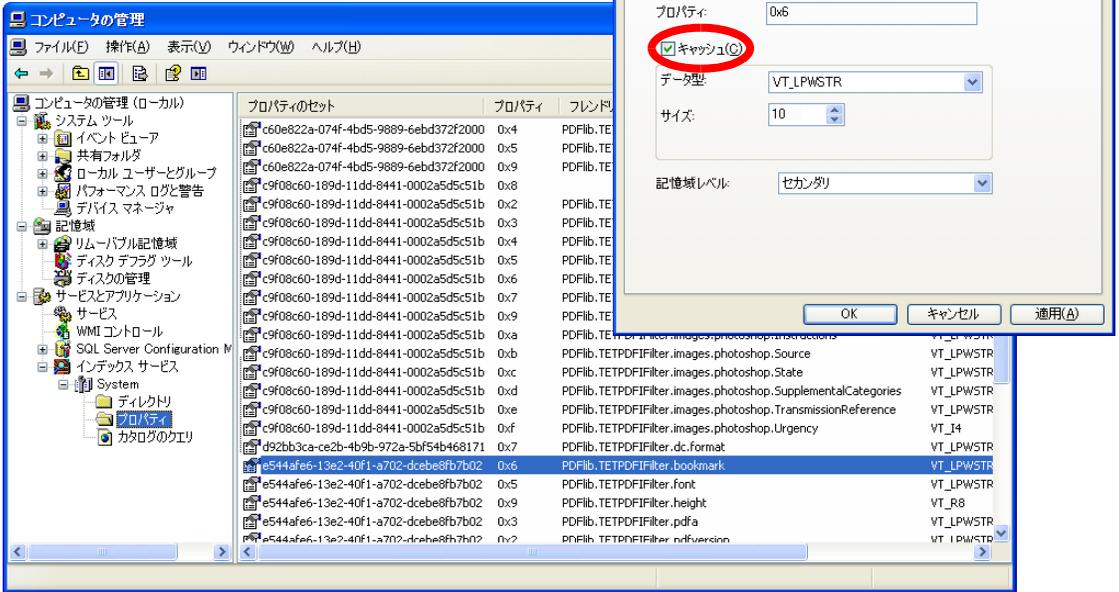
キー *DefaultColumnFile* を見つけます。このキーがすでにある場合は（Windows Search をインストールしてあるときにはこの可能性があります）、その内容としているファイル名に着目します。典型的な値は下記のとおりです：

C:\Documents and Settings\All Users\Application Data\Microsoft\Search\Config\msscolmn.txt

このキーがすでにある場合は、次のステップではその場所にある既存のテキストファイルを編集できます。このキーがない場合は、「編集」→「新規」→「文字列値」をクリックしてそれを作成し、次のステップで作成する構成ファイルを指し示すフルパス名を追加します。

- ▶ 既存のテキストファイル *msscolmn.txt* を編集し（あるいは適宜、任意の場所に新規作成して）、さきに用意しておいたプロパティ定義群を挿入。
- ▶ カタログを再構築。必ず、構成したプロパティそれぞれについて少なくとも 1 個のファイルがインデックスされているようにします。プロパティは、それがインデックス作成処理の間にピックアップされている場合のみ、以下の構成ステップ群で利用可能になります。

Fig. 4.3.
管理コンソールでインデックス
サービスのプロパティリスト



- ▶ 「サービスとアプリケーション」 → 「インデックス サービス」 → 「System」へ移動し、「プロパティ」を右クリックし、「すべてのタスク」 → 「一覧を最新の情報に更新」を選択。インデックスされた文書群の中で見つかったプロパティそれぞれについて、個別の行が表示されます。プロパティを、1 列目にあるその GUID（「プロパティのセット」と題されています）で特定することもできますし、2 列目にある数値識別子またはプロパティ名（「プロパティ」と題されています）で特定することもできます。識別子か名前かの選択は、プロパティの識別にどの方式を選択したか（GUID+名前かGUID+ID、40 ページの「プロパティ識別と GUID」を参照）に依存します。3 列目（「フレンドリ名」と題されています）は、列定義ファイル内で指定されているユーザーフレンドリプロパティ名を内容としています。
- ▶ 使いたいカスタムプロパティをそれぞれダブルクリック。現れるダイアログで、「キャッシュ」チェックボックスを有効にします（図 4.3 参照）。キャッシュされたプロパティ群のみが検索に役立ちます。
- ▶ インデックスサービスを停止させて再起動し、カタログを再構築します。

メタデータプロパティを検索 カスタムメタデータプロパティを構成した後は、プロパティを検索することができます。表 4.9 にプロパティクエリの例を挙げます。インデックスサービスのためのすべての定義済みプロパティの名前を付章 A「定義済みメタデータプロパティ」に挙げています。プロパティベースのクエリのための文法の説明は下記にあります：

[msdn.microsoft.com/en-us/library/ms690541\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms690541(VS.85).aspx)

表 4.9 インデックスサービスのためのメタデータクエリの例

検索語例	説明
@DocAuthor Doyle	作成者が Doyle を含む
@DocAuthor Conan Doyle	作成者が Conan Doyle を含む
@DocAuthor Doy*	作成者名の中に、Doy で始まる語がある
@size < 500000	文書サイズが 500,000 バイト未満
@contents Mow*	文書内容が、Mow で始まる語を含む。なお、「*」ワイルドカードは、テキストの左へは働かず、右へのみ働きます。
@Write <= 2012/04/07	更新日が 2012 年 4 月 7 日以前（日付形式は固定）
@System.MIMETYPE "application/pdf"	インデックス内のすべての PDF 文書をリスト

5 トラブルシューティング

5.1 TET PDF IFilter が全く動かない

TET PDF IFilter が全く動作していないように見えるときは、下記の項目をチェックしてください。

TET PDF IFilter が正しく登録されているか？ コマンドラインツール *FiltReg.exe* を使って、TET PDF IFilter の正しい登録をチェックすることができます。このプログラムは、さまざまな IFilter に関連しているすべてのファイル拡張子をリストします。具体的には、関連づけられた IFilter DLL のファイル拡張子と名前を印字します。このプログラムは、Microsoft Visual Studio 8 とともにインストールされるほか、Windows 7 用 Windows SDK にも含まれています。詳しくは下記を参照してください：

[msdn.microsoft.com/en-us/library/ms692537\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms692537(VS.85).aspx)

注記 64 ビットエディションの TET PDF IFilter をテストするには、64 ビット版の *FiltReg.exe* が必要です。64 ビット版を得るには、Windows SDK を 64 ビット機にインストールする必要があります。そうでない場合には 64 ビットツール群はインストールされません。

TET PDF IFilter が正しく登録されている場合、*FiltReg.exe* の出力は下記のような行を含んでいます：

Filters loaded by extension:

```
...  
.pdf --> PDFlib TET PDF IFilter 32-bit (C:\Program Files\PDFlib\TET PDF IFilter  
4.1\bin\TETPDFIFilter.dll)
```

TET PDF IFilter が正しく登録されていない場合には、これを手動で登録する必要があります（6 ページの「手動インストール」を参照）。

Windows デスクトップサーチをインストールしてある場合には、正しい IFilter 登録を右記のようにテストできます：Windows タスクバーの通知領域で、インデックス作成アイ

コン  を右クリックし、「デスクトップサーチオプション...」→「詳細設定」→「ファイルの種類」をクリック。すると、ファイルタイプと、関連付いたフィルタの長いリストが生成されます。このリストの中で、「拡張子」列の中の「pdf」へスクロールします。「フィルタの説明」列の中の対応する項目に「PDFlib TET PDF IFilter 32-bit」（64 ビット版なら 64-bit）と書いてあるはずですが。

IFilter Explorer IFilter 登録を分析するためのもう 1 つの有用なツールは Citeknet の IFilter Explorer です。これは、システム内のすべての登録されている IFilter を調査するほか、SQL Server・SharePoint といった製品独自のレジストリ項目群をも考慮に入れます。IFilter Explorer は下記でダウンロード可能です：

www.citeknet.com/Products/IFilters/IFilterExplorer/tabid/62/Default.aspx

TET PDF IFilter と Adobe Acrobat TET PDF IFilter の後に Adobe Reader か Adobe Acrobat をインストールすると（あるいは Acrobat の自動修復モードを走らせると）、それらは TET PDF IFilter のレジストリ項目群を上書きします。この状況を直すには、TET PDF IFilter を

修復モードで走らせるか、6 ページの「手動インストール」の項に従って TET PDF IFilter DLL を手動で登録します。

ライセンスキーは利用可能か？ TET PDF IFilter は、Windows 2000/XP/Vista/7/8 では商用ライセンスキーなしで使えますが、Windows Server 2003/2008 ではライセンスキーが必要です。サーバシステムを使用していて PDF インデックス作成が動作していないように見える場合は、TET PDF IFilter のライセンスキーを見つけれないでいる可能性があります。この場合には、TET PDF IFilter は評価モードで動作しますので、小規模文書に制限されます。

この状況は、Windows イベントログをチェックすることで検知できます（70 ページの「アプリケーションイベントログ」を参照）。ライセンスキーの問題が起きた場合には、ソースが「*TET PDF IFilter*」で分類が「*TET Error*」の項目があります。そのエラーを内容としている行をダブルクリックしてエラーメッセージを調べてください。下記のテキストは、有効なライセンスキーが見つからなかったことを示しています：

```
TET API Error in TetIFilter::Init: open_document_mem:  
Invalid license key (error number 1986)
```

このメッセージを見つけたときは、レジストリにライセンスキーを入力する必要があります（6 ページの「手動インストール」を参照）。

5.2 TET PDF IFilter の動作上の問題

TET PDF IFilter が期待どおりに動作していないように見える場合には、以下に説明する分析手法が役立つ可能性があります。

問題を起こす文書を特定 使っている IFilter クライアントによって、ログ記録の項目は問題のファイルの名前を含んでいることもあればないこともあり、また、ファイル名は有用であることもあればそうでないこともあります。たとえば、Indexing Server はファイル名を含みますが、Windows デスクトップサーチはそれを含まません。一方で、SharePoint は、文書を HTTP を通じてダウンロードして一時ローカルコピーを生成します。そのイベントログはこの一時ローカルファイル名を内容としており、これは元の URL とは無関係となっています。問題の PDF 文書を特定する助けとして、イベントログ項目はファイルサイズをバイト単位で含んでいます。検索エンジン自体を使って問題の文書をすばやく特定することも可能です。

- ▶ Windows デスクトップサーチでは、下記のクエリ表現が使えます (12345 はファイルサイズをバイト単位で表したものと)：

```
size: = 12345
```

- ▶ インデックスサーバでは、下記のクエリを使います (12345 はファイルサイズをバイト単位で表したものと)：

```
@size = 12345
```

- ▶ SharePoint では、ファイルをフィルタしようとして失敗した試みを、SharePoint クロールログ内で特定することが可能です (Shared Services Administration : 「**SharedServices**」 → 「**Search Settings**」 → 「**Crawl Log**」)。ここにリストされている PDF 文書に関するエラーは、Windows アプリケーションイベントログ内で TET PDF IFilter によって発せられているエラーに対応しています。このクロールログ内のファイルのファイルサイズを、イベントログ内の項目と比較すれば、問題を起こしている文書を特定することが可能です。

ロックされている PDF 文書はインデックスされない 何らかのアプリケーションが PDF ファイルをロックしているとき、TET PDF IFilter はその文書をインデックスできません。特に、ファイルは Acrobat で開かれている間はロックされています。IFilter クライアントは、ロックされていた文書を後で再試行するかもしれませんが、ロックされている文書が解放されるまではインデックスは不完全になります。ですので、インデックス作成中は PDF 文書を Acrobat で表示することを避けることを推奨します。

5.3 PDF 文書がインデックスされないか不完全

SharePoint Server と Search Server には、大容量文書のインデックス作成に影響を与えるさまざまな制約があります。これらの制約は Microsoft の文書ではよく説明されていないので、以下の諸注意は、Microsoft のサポート記事とブログに基づいて情報を収集しているものです。これらの諸注意は公認のものではありませんので、疑点があれば Microsoft に問い合わせてください。

最大ファイルサイズ レジストリ値 *MaxDownloadSize* は、クロールされインデックスされる文書群の最大ファイルサイズを決定します（単位：MB、デフォルト：16）。

最大膨張係数 レジストリ値 *MaxGrowFactor* は、1 個の文書に対する、インデックスされるテキストの最大量を決定するために *MaxDownloadSize* 値に乘じられる係数を内容としています。この係数が必要な理由は、テキストがファイル内部で圧縮されている場合があるからです。PDF 文書でもたいていはそうです（単位：なし、デフォルト：4）。

デフォルト設定では、文書あたり最大 64MB の抽出テキストがインデックスされます。正確な製品とバージョンによって、*MaxDownloadSize*・*MaxGrowFactor* レジストリ値は、Windows レジストリ内の下記のキーの下に見つかる可能性があります：

```
HKEY_LOCAL_MACHINE\SOFTWARE\Microsoft\Office Server\<バージョン>
  \Search\Applications\<GUID>\Gathering Manager
```

```
HKEY_LOCAL_MACHINE\SOFTWARE\Microsoft\Shared Tools\Web Server Extensions\<バージョン>
  \Search\Applications\<GUID>\Gathering Manager
```

```
HKEY_LOCAL_MACHINE\Software\Microsoft\SPSSearch\Gathering Manager
```

ここで *GUID* はインストールごとに異なります。

チャンクバッファサイズ もう 1 つの制約は、文書あたりの、インデックスされることができる一意な語の総数に影響を与えます。値 *CB_ChunkBufferSizelnMegaBytes* は、文書あたりの、一意な語のコレクションのために予約されるスペースを決定します（単位：MB、デフォルト：8）。

文書用予約バイト数 *CB_MinBytesReservedForDoc* 値は *CB_ChunkBufferSizelnMegaBytes* 値に依存します。これは、*CB_ChunkBufferSizelnMegaBytes* の値よりも 2 MB 少なくするべきです。ただしデフォルト値どうしの間ではこの関係が成り立っていません（単位：バイト、デフォルト：3,145,728）。

正確な製品とバージョンによって、*CB_ChunkBufferSizelnMegaBytes*・*CB_MinBytesReservedForDoc* レジストリ値は、Windows レジストリ内の下記のキーの下に見つかる可能性があります：

```
HKEY_LOCAL_MACHINE\SOFTWARE\Microsoft\Office Server\<バージョン>
  \Search\Global\Gathering Manager
```

```
HKEY_LOCAL_MACHINE\SOFTWARE\Microsoft\Shared Tools\Web Server Extensions\<バージョン>
  \Search\Global\Gathering Manager
```

Microsoft Office SharePoint Server 2007のためにこれらの値を説明しているMicrosoftサポート記事が下記にあります：

support.microsoft.com/kb/970776/EN-US

5.4 デバッグ機能

検索結果が自分の期待にそわず、インデックスされた文書群から抽出されたテキスト内容に問題があると推測している場合には、以下に説明するデバッグツールが有用である可能性があります。

アプリケーションイベントログ TET PDF IFilter は、Windows イベントログの中に、さまざまなイベントについて項目を作成します。このアプリケーションイベントログは下記のようにチェックできます：

- ▶ Windows XP : 「スタート」 → 「設定」 → 「コントロール パネル」 → 「管理ツール」 → 「イベント ビューア」を選択。左のペーンで「アプリケーション」をクリック。
- ▶ Windows Vista/7/8 : 「スタート」をクリックし、「検索の開始」ボックスの中に「イベント ビューア」とタイプし、「イベント ビューア」プログラムをクリック。「イベント ビューア」ウィンドウの中で「Windows ログ」 → 「アプリケーション」をクリック。
- ▶ TET PDF IFilter のフィルタリングイベントは、ソースが「TET PDF IFilter」の項目を生成します。エラーを内容とする行をダブルクリックして、エラーメッセージを調べます。

アプリケーションイベントログ内の項目群は、レジストリ値

`HKEY_LOCAL_MACHINE\SOFTWARE\PDFlib\TET PDF IFilter4\4.1\logging`

を表 5.1 に従った DWORD 値に設定することによって、さまざまなクラスのイベントについて有効にすることができます。このログ記録レベルはデフォルトでは 1 に設定されています。なお、1 個の PDF 文書からイベントログ内に複数のエラーメッセージが出る場合があります。たとえば、損傷したページがそれぞれ項目を 1 つずつ生成するかもしれません。

表 5.1 Windows イベントログのためのログ記録レベル一覧

レベル	要約	ログ記録されるイベント
0	サイレント	なし：あらゆるエラーメッセージが抑制されます
1 (default)	エラー	すべての失敗した TET 関数呼び出しと TET が発生させた例外。例：ライセンスキーが無効が見つからない、ユーザパスワードを必要とする暗号化 PDF、修復できない深刻に損傷した PDF 文書。レジストリ読み取りの問題。構成ファイル解析の問題。
2	アクティブ ティ	1に、すべての Load()・Init() IFilter インタフェース呼び出しを加えたもの。XML 構成ファイル読み取り
3	詳細	2に、テキストと LCID を含むプロパティを抽出するための IFilter インタフェース呼び出しに関する詳細を加えたもの。このレベルは膨大なログ項目を生成します。

TET PDF IFilter は文書からどのプロパティとテキストを出力したか？ ある特定の文書から TET PDF IFilter が抽出している正確なテキストを見るには、Windows SDK に入っているツール *FiltDump.exe* が使えます。この場合も、64ビット TET PDF IFilter DLL をテストするにはこのツールの 64ビット版が必要です。詳しくは下記を参照してください：

[msdn.microsoft.com/en-us/library/ms692535\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms692535(VS.85).aspx)

-o オプションを用いて、*FiltDump.exe* の出力を UTF-16 符号化されたファイルヘリダイレクトさせることも可能です。こうすると、TET PDF IFilter によって出力されるテキストに対する正確な Unicode テキストと、検知されたロケール (LCID) を見るのが可能になります。サンプル呼び出し：

```
FiltDump.exe -o udhr_japanese.txt udhr_japanese.pdf
```

サンプル出力：

```
FILE: udhr_japanese.pdf
IFILTER: CLSID == {47A1AF35-C345-475D-AE68-EB07E948BD07}
IFILTER: Using IPersistStream
IFILTER: IFilter->Init returned IFILTER_FLAGS_OLE_PROPERTIES flag
```

```
CHUNK: -----
Attribute = {007867F0-C59B-43FC-AB1E-8EEE77057254}\3 (Unknown)
idChunk = 1
BreakType = 2 (Sentence)
Flags (chunkstate) = (Value)
Locale = 1031 (0x407)
IdChunkSource = 1
cwcStartSource = 0
cwcLenSource = 0
```

```
VALUE: -----
Type = 31 (0x1f), VT_LPWSTR
Value = "4.0"
```

```
CHUNK: -----
Attribute = {007867F0-C59B-43FC-AB1E-8EEE77057254}\4 (Unknown)
idChunk = 3
BreakType = 2 (Sentence)
Flags (chunkstate) = (Value)
Locale = 1031 (0x407)
IdChunkSource = 3
cwcStartSource = 0
cwcLenSource = 0
```

```
VALUE: -----
Type = 64 (0x40), VT_FILETIME
Value = "2010/06/10:08:28:04.587"
```

```
CHUNK: -----
Attribute = {B725F130-47EF-101A-A5F1-02608C9EEBAC}\19 (System.Search.Contents)
idChunk = 11
BreakType = 2 (Sentence)
Flags (chunkstate) = (Text)
Locale = 9 (0x9)
IdChunkSource = 11
cwcStartSource = 0
cwcLenSource = 0
```

```
TEXT: -----
UDHR - Japanese
```

```
CHUNK: -----
Attribute = {B725F130-47EF-101A-A5F1-02608C9EEBAC}\19 (System.Search.Contents)
idChunk = 12
BreakType = 2 (Sentence)
Flags (chunkstate) = (Text)
Locale = 17 (0x11)
IdChunkSource = 12
```

```
cwcStartSource = 0
cwCLenSource = 0
```

TEXT: -----

...この文書のテキスト内容...

TET カーネルログ記録 XML 構成ファイル内で適切な TET オプション群を設定することによって、詳細な TET カーネルログ記録を有効にすることも可能です (必ず XML 構成ファイルのファイル名を指定してください。73 ページの 6 章「XML 構成ファイル」を参照):

```
<Tet>
  <TetOptions>logging={filename=C:\debug.log classes={pcos=2}}</TetOptions>
</Tet>
```

こうすると、TET 関数への内部呼び出しに関する詳細やエラーメッセージなどを含むログ・ファイルが生成されます。必ず、TET PDF IFilter を呼び出すサービスから書き込み可能なファイル名を用いてください。また、TET ログ記録は大量の出力を生成し、フィルタリング処理の速度を落とすことに留意しておいてください。

6 XML 構成ファイル

6.1 構成ファイルを使用

TET PDF IFilter の動作は XML 構成ファイルで制御できます。サンプル構成ファイルが TET PDF IFilter とともにインストールされています。

構成ファイルの場所を指定 構成ファイルは、その構成ファイルのフルパス名を持つ文字列値を内容とする下記のレジストリキー内で指定できます：

```
HKEY_LOCAL_MACHINE\SOFTWARE\PDFlib\TET PDF IFilter4\configfile
```

注記 XML 構成ファイルは、TET PDF IFilter のインストールディレクトリの外に置くことを推奨します。こうすれば、TET PDF IFilter のアップデート版をインストールした後などにインストールディレクトリが変わっても構成は生き残ります。

このレジストリ項目が存在しない場合や、空文字列を内容としている場合には、デフォルト構成が用いられます。レジストリ項目で指定されている構成ファイルを開くことができない場合には、アプリケーションイベントログに警告が書き込まれ、インデックス作成にはデフォルト構成が用いられます。構成ファイルの XML 解析が失敗した場合には、イベントログに警告が書き込まれ、インデックス作成は一切行われません。

注記 インストーラは、構成ファイルのためのレジストリ項目を生成しません。これは、必要に応じてユーザーが行う必要があります。

1 台のマシン上の TET PDF IFilter に対しては、ただ 1 つの構成ファイルだけが使えます。ただし、32 ビット版と 64 ビット版は別々の構成ファイルを同一マシン上に持つことも可能です。なぜなら、上記のレジストリ項目はそれぞれ 32 ビットレジストリと 64 ビットレジストリの中で検索されるからです。

構成ファイルを変更したときは、その変更を有効にするには、インデックスを再構築する必要があります。

定義済み XML 構成ファイル いくつかの定義済み構成ファイルが TET PDF IFilter とともにインストールされています：

- ▶ ファイル *default.xml* は、TET PDF IFilter の内部デフォルト設定群を記述しています。これは、カスタマイズした構成ファイルを作成するための開始点として役立つ可能性があります。
- ▶ ファイル *starter.xml* は、TET PDF IFilter とともにインストールされているスタータサンプル群とともに用いることができるプロパティ定義群を内容としています。

XML 名前空間とスキーマ記述 表 6.1 に、XML 構成ファイルで利用可能な要素と属性を挙げます。XML 構成のための名前空間 URI は下記のとおりです：

```
http://www.pdflib.com/XML/TET_PDF_IFilter3/TET_PDF_IFilter_Config-3.0.xsd
```

注記 このスキーマの名前空間 URI とダウンロード場所は、バージョン番号 3 を含んでいますが、これは、カレントのスキーマが、TET PDF IFilter 3.0 で用いられるスキーマと互換であるからです。

XML 構成言語のための XSD スキーマ記述が、TET PDF IFilter とともにインストールされています。また、上記の名前空間識別子として作用する URI にもあります。このスキーマファイルを適切な XML エディタとともに使えば、生成した XML 構成ファイルが TET PDF IFilter が期待する文法を必ず遵守しているようにすることができます。

XML 要素・属性のためのカスタムデータ型 値の説明が提供されている所を除いて、すべての要素は空です。下記のカスタムデータ型が XML 構成ファイル内で用いられます：

- ▶ LCID：16 進か 10 進のローケル識別子。下記参照：

[msdn.microsoft.com/en-us/library/ms776294\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms776294(VS.85).aspx)

値 0x0800 は、カレントシステムデフォルトローケルへ翻訳されます。

- ▶ GUID：ITU-T Rec. X.667 (www.itu.int/ITU-T/studygroups/com17/oid/X.667-E.pdf を参照) に従った一意な 128 ビット識別子の 16 進表記。各部分はダッシュキャラクタ「-」で区切られる必要があります。GUID を作成するツールはさまざまなものが入手可能です。オンラインサービスを使うこともできます。たとえば下記にあるものです：

www.itu.int/ITU-T/asn1/uuid.html

- ▶ pCOSパス：PDFオブジェクトを記述する拡張pCOSパス。pCOSリファレンスと、38ページの「拡張 pCOS パス」の項で説明している pCOS 拡張を参照してください。
- ▶ オプションリスト：PDFlib TET リファレンスマニュアルで仕様定義されている文法に従ったオプションリストを内容とする文字列。
- ▶ 言語識別子：RFC 1766 に従った XMP 言語修飾子、またはその文書内のデフォルト言語を識別する *x-default*。

6.2 XML 要素・属性一覧

表 6.1 に、XML 構成ファイルの要素と属性の詳細を示します。XML 構成ファイルによって制御される効果に関するもっと詳しい情報は、このマニュアルのそれぞれの箇所にあります。個別の IFilter クライアントの場合の要請・推奨設定は、37 ページの 3 章「メタデータをインデックス」内のクライアントごとの節に挙げてあります。

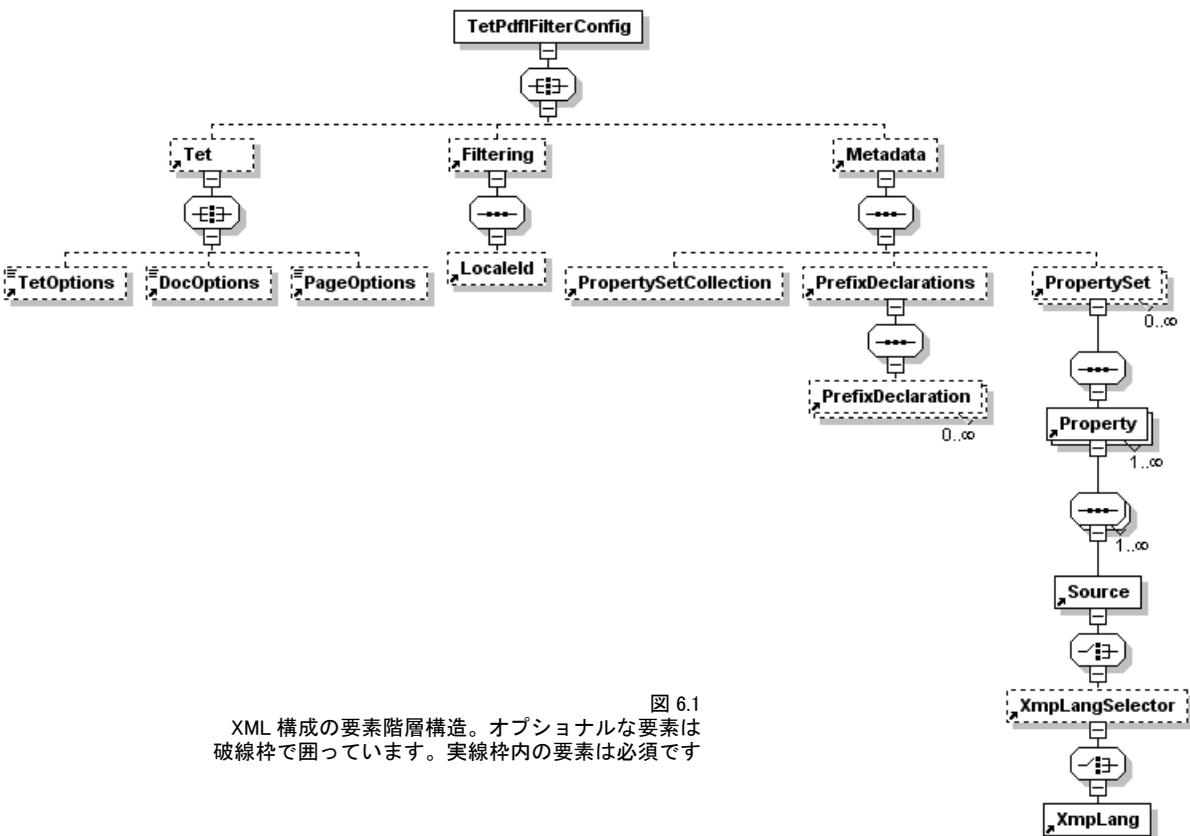


図 6.1 XML 構成の要素階層構造。オプションな要素は破線枠で囲っています。実線枠内の要素は必須です

表 6.1 構成ファイル内の XML 要素・属性一覧

要素	その要素とその属性群の説明
DocOptions 親 : Tet	(0 回または 1 回出現可能) その値は、TET カーネル内の TET_open_document() に対するオプションリストを内容とします。
Filtering 親 : TetPdfFilterConfig	<p>(0 回または 1 回出現可能) PDF フィルタリング処理の詳細を指定します。可能な属性 :</p> <p>indexNestedPdf (論理型。オプション) PDF 添付群を再帰的に処理 (19 ページの 2.1 「PDF 文書のさまざまな領域」を参照)。デフォルト : true</p> <p>indexPageContents (論理型。オプション) PDF ページ群の内容がインデックスされるかどうかを示。ページ内容インデックス作成を無効にすることは、検索がメタデータプロパティ群によってのみ推進されるシナリオにおいて有用でしょう。デフォルト : true</p> <p>metadataHandling (選択肢。オプション) メタデータ処理の種類を選択 (45 ページの 3.6 「メタデータプロパティをテキストとしてインデックス」を参照)。デフォルト : property</p> <p>ignore すべてのメタデータプロパティを落とします。これは、デバッグや、メタデータが必要でない状況においてパフォーマンス最適化のために有用でしょう。</p> <p>property メタデータをプロパティとして処理。</p> <p>propertyAndPrefixedText メタデータをプロパティとして処理することに加え、カスタムプロパティに対しては textIndexPrefix (あれば) で指定されている接頭辞を、定義済みプロパティに対しては 45 ページの表 3.2 に従った接頭辞を頭に付加し、結果をプレーンテキストとして処理。</p> <p>propertyAndText メタデータをプロパティとして処理することに加え、メタデータをプレーンテキストとして処理。</p> <p>useIdentifier (論理型。オプション) Property 要素に identifier 属性と friendlyName 属性が両方ある場合に、プロパティを識別するために identifier 属性が用いられるかどうかを指定。デフォルト : true</p>

表 6.1 構成ファイル内の XML 要素・属性一覧

要素	その要素とその属性群の説明
LocaleId 親 : Filtering	(0 回または 1 回出現可能) ロケール ID 検知を構成 (24 ページの 2.2 「自動言語検知」を参照)。可能な属性 : <ul style="list-style-type: none"> arabic (LCID. オプション) アラビア文字テキストに対する LCID。デフォルト : 0x0401 アラビア語 (SA) chinese (LCID. オプション) 中国文字テキストに対する LCID。デフォルト : 0x0804 中国語 (中華人民共和国) cyrillic (LCID. オプション) キリル文字テキストに対する LCID。デフォルト : 0x0419 ロシア語 (RU) default (LCID. オプション) detection が無効の場合にすべてのテキスト塊に対して用いられるグローバル LCID。デフォルト : 0x0800 (システムロケール) detection (選択肢. オプション) 自動 LCID 検知を制御。デフォルト : auto <ul style="list-style-type: none"> auto 用字系と統計的言語分析に基づいて LCID を決定。 disabled LCID 検知を無効にする。他のすべての属性は、default と useCatalogLang 以外、無視されます。 script (TET PDF IFilter 4.0) 用字系に基づいて LCID を決定。 latin (LCID. オプション) ラテン文字テキストに対する LCID。デフォルト : 0x0409 英語 (US) useCatalogLang (論理型. オプション. TET PDF IFilter 4.0) 文書のカatalogの中の Lang 項目が評価されるかどうかを指定。true の場合、TET PDF IFilter は PDF 文書カATALOG内の Lang 項目をチェックします。存在する場合には、その Lang 項目は LCID へ変換されます。その変換がうまくいった場合には、この LCID は LocaleId/@default 属性の値を上書きします。この LCID がアラビア・中国・キリル・ラテン用字系のうちのいずれか 1 つに属する場合には、それは LocaleId 要素の対応する属性の値を上書きします。デフォルト : true
Metadata 親 : TetPdfFilterConfig	(0 回または 1 回出現可能) メタデータプロパティ群を指定 (42 ページの 3.4 「カスタムメタデータプロパティ」を参照)。存在する場合には、この要素は Filtering と Tet の後に出現する必要があります。
PageOptions 親 : Tet	(0 回または 1 回出現可能) その値は、TET カーネル内の TET_open_page() に対するオプションリストを内容とします。
PrefixDeclaration 親 : PrefixDeclarations	(0 回または 1 回出現可能) Source/@xmpName 内で使える名前空間接頭辞を宣言。可能な属性 : <ul style="list-style-type: none"> prefix (コロン「:」キャラクターを含まない文字列。必須) 名前空間 URI の短縮形として用いられる接頭辞。 uri (URI。必須) 名前空間 URI
PrefixDeclarations 親 : Metadata	(0 回または 1 回出現可能) Xmp 要素群内の XMP プロパティ群に対する名前空間接頭辞群を宣言

表 6.1 構成ファイル内の XML 要素・属性一覧

要素	その要素とその属性群の説明
Property 親 : PropertySet	(0 回または 1 回出現可能) インデックス作成のためのメタデータプロパティを指定 (42 ページの 3.4 「カスタムメタデータプロパティ」を参照)。 identifier と friendlyName のうち少なくとも 1 つが存在する必要があります。両方が与えられている場合には、Filtering/@useIdentifier=false でないかぎり、identifier が IFilter インタフェース内で用いられます。 可能な属性 : identifier (整数 ≥ 2。オプション) PropertySet 内でプロパティを一意に識別する数値。 emitAsVector (論理型。オプション) true の場合、プロパティ値は、値の数にかかわらず、1 個のベクトル実体として出力されます。 false の場合、プロパティは単値として出力されます。複数の源項目が見つかったときは、複数の単値が出力されます。デフォルト : false friendlyName (文字列。オプション) PropertySet 内のプロパティを一意に識別する名前。これはプロパティを文書化するために、あるいは identifier のかわりとして使うことができます。 precedence (選択肢。オプション) 複数の Source 要素に対する優先権を指定 (デフォルト : first-wins) : first-wins 最初の空でない源が用いられます。 try-all すべての空でない源がそのプロパティに寄与します。 textIndexPrefix (文字列。オプション) Filtering/@metadataHandling が propertyAndPrefixedText である場合にプロパティ値の頭に付加される文字列。デフォルト : 空 type (選択肢。オプション) そのメタデータプロパティの Windows データ型。可能な選択肢は Boolean・DateTime・Double・Int32・String です。デフォルト : String
PropertySet 親 : Metadata	(0 回または 1 回出現可能) 同一 GUID を持つプロパティ群のカスタムのセットのフィルタリングを指定 (42 ページの 3.4 「カスタムメタデータプロパティ」を参照)。 存在する場合には、この要素は PropertySetCollection と PrefixDeclarations の後に出現する必要があります。 可能な属性 : guid (GUID。必須) そのプロパティセットに対する一意な 128 ビット識別子を 16 進表記で。
PropertySet-Collection 親 : Metadata	(0 回または 1 回出現可能) 定義済みプロパティセットコレクションのフィルタリングを指定 (41 ページの 3.3 「定義済みメタデータプロパティ」を参照)。プロパティのリストは付章 A 「定義済みメタデータプロパティ」にあります。可能な属性 : documentXmp (論理型。必須) 文書 XMP プロパティ群を出力。デフォルト : false imageXmp (論理型。必須) 画像 XMP プロパティ群を出力。デフォルト : false internal (論理型。必須) TET PDF IFilter の内部プロパティ群を出力。デフォルト : true pdf (論理型。必須) PDF 独自プロパティ群を出力。デフォルト : false shell (論理型。必須) シェルプロパティ群を出力。デフォルト : true

表 6.1 構成ファイル内の XML 要素・属性一覧

要素	その要素とその属性群の説明
Source 親 : Property	(0 回または 1 回出現可能) メタデータプロパティに対する 1 個ないし複数の源を指定。 Property/@precedence が値 first-wins を持っている場合には、要素群の順番は意味を持ちます。下記の属性のうち少なくとも 1 つが与えられる必要があります。 可能な属性 : pdfObject (pCOS パス。オプション) そのプロパティを含んでいる 1 個ないし複数の PDF オブジェクトへの拡張 pCOS パス。デフォルト : /Root/Metadata (すなわち文書レベル XMP) xmpName (スキーマの接頭辞・コロン「:」・プロパティ名から成る文字列。オプション) 完全修飾された XMP プロパティ名。接頭辞は、それが PrefixDeclaration 要素内で宣言されていれば、名前空間 URI のかわりに使えます。この属性は、pdfobject が 1 個ないし複数の XMP ストリームを指し示している場合にのみ用いられます。デフォルト : 空
Tet 親 : TetPdfFilterConfig	(0 回または 1 回出現可能) TET カーネルに対する処理オプション群を指定。オプションリスト文法と可能なオプションの説明は TET マニュアルを参照してください。いくつかのオプションは TET PDF IFilter によって上書きされます。
TetOptions 親 : Tet	(0 回または 1 回出現可能) TET カーネル内の Tet_set_option() に対するオプションリストを指定。
TetPdfFilterConfig 親 : (なし)	(ルート要素としてちょうど 1 回出現する必要があります) XML 構成ファイルのルート要素。可能な属性 : version (文字列。オプション。TET PDF IFilter 4.0) この構成が想定して書かれている TET PDF IFilter のバージョンを指定。TET PDF IFilter 3 の構成ファイルはこの属性に対応していないのでデフォルトは 3.0 です。新しい構成は、適切なバージョン番号 (TET PDF IFilter 4 なら 4.0) を持つこの属性を含む必要があります。デフォルト : 3.0
XmpLang 親 : XmpLangSelector	(XmpLangSelector/@languages=subset の場合にちょうど 1 回出現する必要があります) XMP プロパティの言語を指定。可能な属性 : lang (言語識別子。必須) 言語の名前。現在、x-default が唯一可能な値です。
XmpLangSelector 親 : Xmp	(0 回または 1 回出現可能) インデックス作成のための XMP プロパティの言語変種を選択。これは型 Lang Alt の XMP 源を持つプロパティに対してのみ意味を持ちます。可能な属性 : languages (選択肢) そのプロパティの言語独自インデックス作成を指定 (デフォルト : all) : all そのプロパティのすべての利用可能な言語項億がインデックスされません。 subset XmpLang 要素で指定されている言語群だけがインデックスされます。

6.3 サンプル構成ファイル

以下に、TET PDF IFilter のための完全な XML 構成ファイルを示します：

```
<?xml version="1.0" encoding="UTF-8"?>
<!--
    TET PDF IFilter用XML構成ファイル
    (c) PDFlib GmbH 2008-2012 www.pdflib.com
    このファイルは下記レジストリキーに構成される必要があります：
    HKEY_LOCAL_MACHINE\SOFTWARE\PDFlib\TET PDF IFilter4\configfile
-->

<n:TetPdfIFilterConfig
  xmlns:n="http://www.pdflib.com/XML/TET_PDF_IFilter3/TET_PDF_IFilter_Config-3.0.xsd"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  xsi:schemaLocation="http://www.pdflib.com/XML/TET_PDF_IFilter3/TET_PDF_IFilter_
Config-3.0.xsd http://www.pdflib.com/XML/TET_PDF_IFilter3/TET_PDF_IFilter_Config-3.0.xsd"
  version="4.0">

  <n:Tet>
    <n:TetOptions></n:TetOptions>
    <n:DocOptions></n:DocOptions>
    <n:PageOptions></n:PageOptions>
  </n:Tet>

  <n:Filtering indexNestedPdf="true" metadataHandling="property" useIdentifier="true">
    <n:LocaleId
      detection="auto"
      useCatalogLang="true"
      default="0x0800"
      arabic="0x0401"
      chinese="0x0804"
      cyrillic="0x0419"
      latin="0x0409"/>
  </n:Filtering>

  <n:Metadata>
    <n:PropertySetCollection
      documentXmp="false"
      imageXmp="false"
      internal="true"
      pdf="false"
      shell="true"/>

    <n:PropertySet guid="b01ca440-1b9f-11dd-8b87-0002a5d5c51b">
      <n:Property identifier="2">
        <n:Source pdfObject="/Info/Producer"/>
      </n:Property>
    </n:PropertySet>
  </n:Metadata>

</n:TetPdfIFilterConfig>
```

A 定義済みメタデータプロパティ

表 A.1 に挙げるプロパティは、TET PDF IFilter に内蔵されているので、カスタム構成を一切必要としません。実際のプロパティ定義は TET PDF IFilter に内蔵されていますが、これらのプロパティを使うには、使いたいプロパティセットコレクション群を XML 構成ファイル内で有効にする必要があります (41 ページの 3.3 「定義済みメタデータプロパティ」を参照)。Windows Search の場合には、プロパティ群を *registerpropdesc.exe* ツールで登録する必要があります。インデックスサービスの場合には、対応する列定義ファイルを構成する必要があります。

表 A.1 TET PDF IFilter におけるプロパティ処理：定義済みプロパティセットコレクション一覧

Windows Search とインデックスサービス (IS) のためのプロパティ名。そのプロパティがテキストとしてインデックスされる場合には接頭辞を導出するためにも用いられます	データ型	多値	プロパティセット GUID / プロパティ ID	源： XMP プロパティまたは pCOS パス
シェルプロパティセットコレクション				
System.Document.Contributor	String	○	F334115E-DA1B-4509-9B3D-119504DC7ABB/100	dc:contributor
System.Document.DateCreated	DateTime	×	F29F85E0-4FF9-1068-AB91-08002B27B3D9/12	xmp:CreateDate, /Info/CreationDate
System.Document.DateSaved	DateTime	×	F29F85E0-4FF9-1068-AB91-08002B27B3D9/13	xmp:ModifyDate, /Info/ModDate
System.Document.DocumentID	String	×	E08805C8-E395-40DF-80D2-54F0D6C43154/100	dc:identifier
System.Document.PageCount (IS : DocPageCount)	Int32	×	F29F85E0-4FF9-1068-AB91-08002B27B3D9/14	length:pages
System.Document.Version	String	×	D5CDD502-2E9C-101B-9397-08002B2CF9AE/29	xmpMM:VersionID
System.Search.Contents (WDS 2 : なし、Indexing Service : Contents)	String	○	B725F130-47EF-101A-A5F1-02608C9EEBAC/19	PDF ページ群のテキスト内容
System.Title (WDS 2.6 ・ IS : DocTitle)	String	×	F29F85E0-4FF9-1068-AB91-08002B27B3D9/2	dc:title["x-default"], /Info/Title
System.Subject (WDS 2.6 ・ IS : DocSubject)	String	×	F29F85E0-4FF9-1068-AB91-08002B27B3D9/3	dc:description["x-default"], /Info/Subject
System.Author (WDS 2.6 ・ IS : DocAuthor)	String	○	F29F85E0-4FF9-1068-AB91-08002B27B3D9/4	dc:creator, pdf:Author, xmp:Author, /Info/Author
System.Keywords (WDS 2.6 ・ IS : DocKeywords)	String	○	F29F85E0-4FF9-1068-AB91-08002B27B3D9/5	pdf:Keywords, /Info/Keywords

表 A.1 TET PDF IFilter におけるプロパティ処理：定義済みプロパティセットコレクション一覧

Windows Search とインデックスサービス (IS) のためのプロパティ名。そのプロパティがテキストとしてインデックスされる場合には接頭辞を導出するためにも用いられます	データ型	多値	プロパティセット GUID / プロパティ ID	源： XMP プロパティまたは pCOS パス
System.MIMETYPE	String	×	0B63E350-9CCC-11D0-BCDB-00805FCCCE04/5	application/pdf (固定)
System.DateModified (IS : Write)	DateTime	×	B725F130-47EF-101A-A5F1-02608C9EEBAC/14	xmp:ModifyDate, /Info/ModDate
System.ApplicationName	String	×	F29F85E0-4FF9-1068-AB91-08002B27B3D9/18	xmp:CreatorTool, /Info/Creator
System.Kind	String	×	1E3EE840-BC2B-476C-8237-2ACD1A839B22/3	Document (固定)
PDF プロパティセットコレクション				
PDFlib.TETPDFIFilter.pdfversion (PDF バージョンに 10 を乗じた内容。例 : 「16」)	String	×	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/2	pdfversion
PDFlib.TETPDFIFilter.pdfa	String	×	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/3	pdfa
PDFlib.TETPDFIFilter.pdfx	String	×	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/4	pdfx
PDFlib.TETPDFIFilter.font	String	○	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/5	fonts[*]/name
PDFlib.TETPDFIFilter.bookmark	String	○	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/6	bookmarks[*]/Title
PDFlib.TETPDFIFilter.annotation	String	○	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/7	pages[*]/annots[*]/Contents
PDFlib.TETPDFIFilter.width	Double	○	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/8	pages[*]/width
PDFlib.TETPDFIFilter.height	Double	○	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/9	pages[*]/height
PDFlib.TETPDFIFilter.producer	String	×	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/10	/Info/Producer
PDFlib.TETPDFIFilter.trapped	String	×	E544AFE6-13E2-40F1-A702-DCEBE8FB7B02/11	/Info/Trapped
XMP 文書メタデータプロパティセットコレクション (XMP 2005 仕様から)				
Dublin Core				
PDFlib.TETPDFIFilter.dc.contributor	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/2	dc:contributor
PDFlib.TETPDFIFilter.dc.coverage	String	×	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/3	dc:coverage
PDFlib.TETPDFIFilter.dc.creator	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/4	dc:creator
PDFlib.TETPDFIFilter.dc.date	DateTime	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/5	dc:date

表 A.1 TET PDF IFilter におけるプロパティ処理：定義済みプロパティセットコレクション一覧

Windows Search とインデックスサービス (IS) のためのプロパティ名。そのプロパティがテキストとしてインデックスされる場合には接頭辞を導出するためにも用いられます	データ型	多値	プロパティセット GUID / プロパティ ID	源： XMP プロパティまたは pCOS パス
PDFlib.TETPDFIFilter.dc.description	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/6	dc:description
PDFlib.TETPDFIFilter.dc.format	String	×	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/7	dc:format
PDFlib.TETPDFIFilter.dc.identifier	String	×	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/8	dc:identifier
PDFlib.TETPDFIFilter.dc.language	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/9	dc:language
PDFlib.TETPDFIFilter.dc.publisher	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/10	dc:publisher
PDFlib.TETPDFIFilter.dc.relation	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/11	dc:relation
PDFlib.TETPDFIFilter.dc.rights	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/12	dc:rights
PDFlib.TETPDFIFilter.dc.source	String	×	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/13	dc:source
PDFlib.TETPDFIFilter.dc.subject	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/14	dc:subject
PDFlib.TETPDFIFilter.dc.title	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/15	dc:title
PDFlib.TETPDFIFilter.dc.type	String	○	D92BB3CA-CE2B-4B9B-972A-5BF54B468171/16	dc:type
XMP Basic				
PDFlib.TETPDFIFilter.xmp.Advisory	String	○	C60E822A-074F-4BD5-9889-6EBD372F2000/2	xmp:Advisory
PDFlib.TETPDFIFilter.xmp.BaseURL	String	×	C60E822A-074F-4BD5-9889-6EBD372F2000/3	xmp:BaseURL
PDFlib.TETPDFIFilter.xmp.CreateDate	DateTime	×	C60E822A-074F-4BD5-9889-6EBD372F2000/4	xmp:CreateDate
PDFlib.TETPDFIFilter.xmp.CreatorTool	String	×	C60E822A-074F-4BD5-9889-6EBD372F2000/5	xmp:CreatorTool
PDFlib.TETPDFIFilter.xmp.Identifier	String	○	C60E822A-074F-4BD5-9889-6EBD372F2000/6	xmp:Identifier
PDFlib.TETPDFIFilter.xmp.Label	String	×	C60E822A-074F-4BD5-9889-6EBD372F2000/7	xmp:Label
PDFlib.TETPDFIFilter.xmp.MetadataDate	DateTime	×	C60E822A-074F-4BD5-9889-6EBD372F2000/8	xmp:MetadataDate
PDFlib.TETPDFIFilter.xmp.ModifyDate	DateTime	×	C60E822A-074F-4BD5-9889-6EBD372F2000/9	xmp:ModifyDate

表 A.1 TET PDF IFilter におけるプロパティ処理：定義済みプロパティセットコレクション一覧

Windows Search とインデックスサービス (IS) のためのプロパティ名。そのプロパティがテキストとしてインデックスされる場合には接頭辞を導出するためにも用いられます	データ型	多値	プロパティセット GUID / プロパティ ID	源： XMP プロパティまたは pCOS パス
PDFlib.TETPDFIFilter.xmp.Nickname	String	×	C60E822A-074F-4BD5-9889-6EBD372F2000/10	xmp:Nickname
PDFlib.TETPDFIFilter.xmp.Rating	Int32	×	C60E822A-074F-4BD5-9889-6EBD372F2000/11	xmp:Rating
XMP Rights Management				
PDFlib.TETPDFIFilter.xmpRights.Certificate	String	×	0DE7A11C-E2C5-4EFA-8017-BECD888E7EC9/2	xmpRights:Certificate
PDFlib.TETPDFIFilter.xmpRights.Marked	Boolean	×	0DE7A11C-E2C5-4EFA-8017-BECD888E7EC9/3	xmpRights:Marked
PDFlib.TETPDFIFilter.xmpRights.Owner	String	○	0DE7A11C-E2C5-4EFA-8017-BECD888E7EC9/4	xmpRights:Owner
PDFlib.TETPDFIFilter.xmpRights.UsageTerms	String	○	0DE7A11C-E2C5-4EFA-8017-BECD888E7EC9/5	xmpRights:UsageTerms
PDFlib.TETPDFIFilter.xmpRights.WebStatement	String	×	0DE7A11C-E2C5-4EFA-8017-BECD888E7EC9/6	xmpRights:WebStatement
XMP Basic Job Ticket				
PDFlib.TETPDFIFilter.xmpBJ.JobRef	String	○	EBC983EF-C1CF-45C8-A29E-993543A0ECFB/2	xmpBJ:JobRef
XMP Paged-Text				
PDFlib.TETPDFIFilter.xmpTPg.NPages	Int32	×	7A9EB492-35AB-49FE-B364-A21FC9575C28/2	xmpTPg:NPages
PDFlib.TETPDFIFilter.xmpTPg.PlateNames	String	○	7A9EB492-35AB-49FE-B364-A21FC9575C28/3	xmpTPg:PlateNames
Adobe PDF				
PDFlib.TETPDFIFilter.pdf.Keywords	String	×	17EB8447-FC9B-4D4D-81DF-31E9AA770CBF/2	pdf:Keywords
PDFlib.TETPDFIFilter.pdf.PDFVersion	String	×	17EB8447-FC9B-4D4D-81DF-31E9AA770CBF/3	pdf:PDFVersion
PDFlib.TETPDFIFilter.pdf.Producer	String	×	17EB8447-FC9B-4D4D-81DF-31E9AA770CBF/4	pdf:Producer
XMP 画像メタデータプロパティセットコレクション (XMP 2005 仕様から)				
Photoshop				
PDFlib.TETPDFIFilter.images.photoshop.AuthorsPosition	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/2	photoshop:AuthorsPosition
PDFlib.TETPDFIFilter.images.photoshop.CaptionWriter	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/3	photoshop:CaptionWriter
PDFlib.TETPDFIFilter.images.photoshop.Category	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/4	photoshop:Category
PDFlib.TETPDFIFilter.images.photoshop.City	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/5	photoshop:City

表 A.1 TET PDF IFilter におけるプロパティ処理：定義済みプロパティセットコレクション一覧

Windows Search とインデックスサービス (IS) のためのプロパティ名。そのプロパティがテキストとしてインデックスされる場合には接頭辞を導出するためにも用いられます	データ型	多値	プロパティセット GUID / プロパティ ID	源： XMP プロパティまたは pCOS パス
PDFlib.TETPDFIFilter.images.photoshop.Country	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/6	photoshop:Country
PDFlib.TETPDFIFilter.images.photoshop.Credit	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/7	photoshop:Credit
PDFlib.TETPDFIFilter.images.photoshop.DateCreated	DateTime	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/8	photoshop:DateCreated
PDFlib.TETPDFIFilter.images.photoshop.Headline	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/9	photoshop:Headline
PDFlib.TETPDFIFilter.images.photoshop.Instructions	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/10	photoshop:Instructions
PDFlib.TETPDFIFilter.images.photoshop.Source	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/11	photoshop:Source
PDFlib.TETPDFIFilter.images.photoshop.State	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/12	photoshop:State
PDFlib.TETPDFIFilter.images.photoshop.SupplementalCategories	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/13	photoshop:SupplementalCategories
PDFlib.TETPDFIFilter.images.photoshop.TransmissionReference	String	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/14	photoshop:TransmissionReference
PDFlib.TETPDFIFilter.images.photoshop.Urgency	Int32	○	C9F08C60-189D-11DD-8441-0002A5D5C51B/15	photoshop:Urgency
内部プロパティセットコレクション				
PDFlib.TETPDFIFilter.version	String	×	007867F0-C59B-43FC-AB1E-8EEE77057254/2	4.1 (固定)
PDFlib.TETPDFIFilter.tetversion	String	×	007867F0-C59B-43FC-AB1E-8EEE77057254/3	4.0 (固定)
PDFlib.TETPDFIFilter.indextime	DateTime	×	007867F0-C59B-43FC-AB1E-8EEE77057254/4	インデックスランの日付・時刻
PDFlib.TETPDFIFilter.eval	Int32	×	007867F0-C59B-43FC-AB1E-8EEE77057254/5	IFilter が評価モードで走っているなら例外番号

B 更新履歴

本マニュアルの更新履歴

日付	変更
February 13, 2012	▶ <i>TET PDF IFilter 4.1</i> (<i>TET 4.1</i> をベースとする) のための更新
September 22, 2010	▶ <i>TET PDF IFilter 4.op2</i> (<i>TET 4.op2</i> をベースとする) のための更新
July 22, 2010	▶ <i>TET PDF IFilter 4.o</i> (<i>TET 4.o</i> をベースとする) のための更新
August 06, 2008	▶ <i>Search Server</i> のための更新
June 16, 2008	▶ <i>TET PDF IFilter 3.o</i> (<i>TET 3.opre2</i> をベースとする)
June 6, 2008	▶ <i>TET PDF IFilter 3.o beta3</i> (<i>TET 3.opre2</i> をベースとする)
2008 年 5 月 09 日	▶ <i>TET PDF IFilter 3.o beta2</i> (<i>TET 3.opre1</i> をベースとする) のための初版

索引

D

DateTime プロパティ型 42

DocOptions 要素 76

F

Filtering 要素 76

FiltReg.exe 65

G

GUID (グローバル一意識別子) 40

H

HRESULT エラー値 50

I

IFilter Explorer 65

ISO 32000 27

L

LocaleId 要素 77

M

Metadata 要素 77

P

PageOptions 要素 77

PDF バージョン 27

PrefixDeclarations 要素 77

PrefixDeclaration 要素 77

PropertySetCollection 要素 78

PropertySet 要素 78

Property 要素 78

prop コマンドラインツール 50

R

registerpropdesc.exe 50

S

Search Server 12

SharePoint 10

メタデータ 55

Source 要素 79

SQL Server 13

メタデータ 60

T

TetOptions 要素 79

TetPdfFilterConfig 要素 79

Tet 要素 79

U

Unicode

正規化 33

分解 30

Unicode 字形統合 28

Unicode マッピングテーブル 35

UUID (ユニバーサル一意識別子) 40

W

Windows Search 7

メタデータ 49

X

XML 構成ファイル 73

XML 要素・属性 75

XmpLangSelector 要素 79

XmpLang 要素 79

XMP メタデータ 21, 37

あ

暗号化 *PDF* 27

イベントログ 70

インデックス

プロパティをテキストとして 45

インデックスサービス 16

メタデータ 61

か

カスタムプロパティ 42

画像メタデータ 37

言語検知 24

互換分解 30

さ

しおり 22

修復モード
破損 PDF のための 27
字形統合 28
正規化 33
正準分解 30

た

多値プロパティ 44
単値プロパティ 44
注釈 22
定義済みプロパティ 41
トラブルシューティング 65

は

破損 PDF 27
パスワード保護された PDF 27
パッケージ 23
評価版 5
ファイル添付 22
フォームフィールド 21
分解 30
文書情報項目 20, 37
プロパティセットコレクション 41
ベクトル処理 44
保護された PDF 27
ポートフォリオ 23

ま

メタデータ 37
 SharePoint における 55
 SQL Server 60
 Windows Search における 49
 インデックスサービスにおける 61
 カスタムプロパティ 42
 多値プロパティ 44
 単値プロパティ 44
 定義済みプロパティ 41
 プロパティ識別 40
 プロパティセットコレクション 41
 プロパティのベクトル処理 44
 プロパティをテキストとして 45
メタデータを採りページ内容を無視 47

ら

ライセンスキー 5, 66
ログ記録 70
ロケール識別子 (LCID) 24

PDFlib GmbH

Franziska-Bilek-Weg 9
80339 München, Germany
www.pdflib.com
電話 +49・89・452 33 84-0
fax +49・89・452 33 84-99

疑問があたりの際は、PDF メーカーリストと、
tech.groups.yahoo.com/group/pdflib のアーカイブをチェックしてください

ライセンスに関するお問い合わせ
jp.sales@pdflib.com

サポート
jp.support@pdflib.com (お使いのライセンス番号をお書きください)

