

## メールの変換道

Dehenkenでは、各種フォーマットの文書ファイルから安全にテキストデータを取り出す技術を大切にしています。さまざまな形式のフォーマットから、精度良く・安全に取り出す継続的品質・性能向上活動を「変換道」と呼んでいます。

### ●多様なメールファイルからのメール情報の取り出し

本ライブラリは、さまざまなメールで作成された1つのEMLファイルについて、ヘッダとボディとその添付ファイルを取り出します。添付ファイルはさらに添付展開し、そこに圧縮ファイルがあればそれを展開して、最終的に1つのファイルになるまで展開します(展開レベル指定による制限可)。対応している添付ファイルのエンコードは、uuencode(MIME形式/本文理込方式)、base64、quoted-printable、binhexです。

また、UNIXのmailコマンドでいうmboxファイルのような複合メール形式について、1つずつのメールを取り出し展開するAPIを用意しています。対応している復号メール形式は、mbox、thunderbird、Outlook、Outlook Express、Becky!等のソフトウェアで作成されたものです(Becky!のみ添付ファイルの展開は未対応)。

### ●さまざまな圧縮ファイルからのファイルの取り出し

本ライブラリは、さまざまな圧縮ソフトで作成された圧縮ファイルを展開します。展開は1つのファイルから再帰的にいき、最終的には1つのファイルのレベルまで掘り下げます。圧縮したアーカイブファイルは、zip(winzip/pkzip)、lha(自己解凍形式を含む)、gzip、tar+gz、rar、bzip2です。圧縮していないアーカイブファイルはtar/gnutarです。

※各種パスワード付は除きます。

### ●メールヘッダを変換して提供

メールのヘッダ部分を変換して取り出しをすることができます。メールヘッダのCc:などに多くの配信先を指定した場合には、複数行におよぶ場合があります。本ライブラリは、複数行に及んだ行は連結して1行にして結果の提供を行います。また、ヘッダにB-エンコード、あるいはQ-エンコードがあれば、デコードして結果の提供を行います。また、メールの取り出し時の文字コードの指定(後述のlanguageによる)ができますので、ユーザが得られるヘッダ情報は指定された漢字コードで得ることができます。

### ●メール本文を指定の文字コードに変換

メール本文のテキスト情報を指定の文字コードで取り出すことができます。メール本文が途中で破損し、バイナリデータが混入していた場合、破損部以外の文字情報を可能な限り取り出します。また、メールの途中の文字コードに異常データ(0x00-0x1fで改行とタブを除く)が含まれていてもそれを取り除く処理をします。



### ●メール添付情報と、付属するMIME属性を得られます

メールに添付されている複数の添付ファイルの1つ1つを順に取り出すAPIを用意しています。また、このときのメールの添付に付属する添付ファイル名の情報を取り出すことができます。

### ●ユニコードの対応

本ライブラリは、日本語文字コードの2つの方式 Windows-31J(CP932)と JIS X0213:2004に対応しています。これらの文字コード形式を指示した上で、JIS(ISO-2022-JP)/EUC-JP/Shift\_JIS/UTF-8/UTF-16の変換を行います。本ライブラリは、content-typeの指定によって、テキストの文字コードの判別を優先的に行うようにしています。Dehenken MFXライブラリ multi-language版では、ヨーロッパ圏の文字コード ISO-8859 や、アジアの中国語、韓国語についても対応しています。(GB2312/GB18030/KSC 5601/Big5等)。

### ●開発生産性を支援する情報の提供と処理の工夫

本ライブラリを使用して開発する技術者のために、開発生産性を高めるために様々な情報の提供と工夫をしています。

●サンプルソースを用意しており開発時の参考にいただけます。

●ファイルをAPIに与えるだけで、ライブラリが自動判別して動作します。

●ファイルの判定にはサフィックス(接尾文字)の情報を使いません。

●ファイル名指定方式とファイル内容メモリ渡しとの2つのインターフェースをご用意しています。

一般に、プログラム経験のある方であれば、容易に使い方をご理解いただけます。

### ●コールバック関数の指定

ユーザは、ライブラリにコールバック関数を登録することができます。本ライブラリは、大量の圧縮ファイル展開のために長時間処理を占有する場合がありますが、コールバック関数を用いて処理の途中に、中断やWaitの処理を入れることができ、CPUの占有を防ぐことができます。

### ●マルチスレッドに対応

本ライブラリは、CPU数やコア数に応じてテキスト抽出の分散処理による速度の向上ができるよう、マルチスレッドに対応しています(マルチスレッドセーフ)。

### ●破損ファイル対象時の安全性(Broken Files Safeness)

文書ファイルやテキストファイルは、保存時や複製時にデータの破損事故が発生する場合があります。本ソフトウェアが破損したファイルであっても、安全に抽出を継続したり、中断したりできることを確認するために、意図的に破損したファイルで網羅的に動作確認を実施しています。当社のこのような破損ファイルによる安全確認を Broken Files Safeness と呼んでいます。

メールと圧縮ファイルの階層展開ライブラリ

# Dehenken MFX Library



株式会社 データ変換研究所 Dehenken Limited

本社 〒604-8155 京都市中京区錦小路通室町東入占出山町 308 ヤマチュウビル 1F

TEL 075-254-8780 FAX 075-254-8790

横浜営業所 〒231-0048 神奈川県横浜市中区蓬莱町 2-4-7 澤田聖徳ビル 204

URL : <http://www.dehenken.co.jp/> E-Mail : [info\\_ml@dehenken.co.jp](mailto:info_ml@dehenken.co.jp)

EST'D 1999 Dehenken Limited © Copyright Dehenken 2019. All rights reserved.

品質マネジメントシステム ISO 9001:2008 の認証取得

株式会社データ変換研究所は、2011年9月27日付で全社統一の品質マネジメントシステムとしてDNV GLよりISO-9001:2008の認証を取得しました。(現在は2015年版に移行)

認証の対象は「ソフトウェアプロダクトのデザイン・開発・製造」です。

Certificate No. : 02523-2011-AQ-KOB-RvA

Initial certification date : 27 September, 2011

Valid : 27 September, 2017 - 27 September, 2020





## 開発生産性を高める、メールの使用環境を提供

デ変研 MFX ライブラリ (以下、本ライブラリ) を用いてアプリケーションを開発する OEM ユーザを支援するために、開発しやすく工夫した様々な提供を行います。API 使用時のサンプルソース (再利用できる著作権付) もご提供しています。展開した 1 つのメールファイルのファイル名を指定して API に与え、ライブラリがメール形式や圧縮形式を自動判別し展開します。このときのファイルの判定には拡張子 (サフィックスまたは接尾文字) の情報を使いません。

## 本ライブラリで対応している E-Mail 形式の対応仕様について

複合メール形式	ヘッダ展開	本文展開	添付展開
Mbox 形式 ※1※2	○	○	○
PST (Outlook 32 / 64 bit) 形式 ※1※3	○	○	○
DBX (Outlook Express) 形式 ※4	○	○	○
Becky! 形式	○	○	×

- ※1…複合メール形式から展開した 1 通の E-Mail 形式のサイズは 2Gbyte まで、またメールの通数は 2147483647 (signed int の最大値) までとなります。
- ※2…MBOX 形式に対応しているアプリの 1 つの例として、Thunderbird があります。
- ※3…Outlook2002 は、2Gbyte 以上のサイズのファイルへは対応していません。
- ※4…2Gbyte 以上のサイズのファイルへは対応していません。

- **圧縮ファイル**
  - zip (winzip / pkzip : 圧縮形式 / 自己解凍形式)
  - lha (lh1 / lh5 / lh6 / lh7 : 圧縮形式 / 自己解凍形式)
  - tar+gzip / tgz / gzip
  - rar (圧縮形式 / 自己解凍形式)
  - bzip2
  - 7z (圧縮形式 / 自己解凍形式)
  - ※それぞれの圧縮形式において、パスワード付きのものを除きます。
- **アーカイブ形式**
  - tar / gnutar
- **メールヘッダのエンコード**
  - メールヘッダのエンコードは、MIME Q/B 及び RFC2231 に対応しています。

- **メール本文のエンコード**
  - メール本文のエンコードは、エンコードなし (text/plain) base64、quoted-printable に対応しています。
- **添付ファイルのエンコード**
  - 添付ファイルのエンコードは、uuencode (MIME 形式 / 本文埋め込み形式)、base64、quoted-printable、binhex に対応しています。
- **パート判別**
  - E-Mail 形式におけるパートの判別方式は、From と Date ヘッダがあり、改行が 2 つ連続して存在するまでを メールヘッダ、以降をメール本文として扱います。メール本文及びそれ以降が複数のパートで構成されている場合、後続のパート以降は添付ファイルとして扱います。
- **添付ファイル展開**
  - 添付ファイルが MS-Office (Word / Excel / PowerPoint) / PDF ファイルなどの場合、デ変研 TF ライブラリと連動して、テキスト抽出後のファイルも取り出すことができます。
- **圧縮ファイル展開**
  - 添付ファイルが圧縮ファイルであった場合に、圧縮ファイル内を展開して取り出すことができます。添付ファイル中にメールファイルの添付や、圧縮内にメールファイルがあった場合、もしくは、添付ファイルの中に zip 圧縮があり、そのなかに lha 圧縮があるような階層圧縮ファイルも、展開して順に取り出すことができます。
- **安全のための限界値設定**
  - メモリ使用制限 (limit\_total\_memory)
  - 1 つのファイルの制限 (limit\_one\_file)
  - ファイルの使用制限 (limit\_total\_file)
  - ヘッダの最大値の指定 (limit\_eml\_header)
  - メール本文の最大値の指定 (limit\_eml\_body)
  - メールの入れ子展開の階層指定 (limit\_level)
  - テキストとして取り出す文字コードの指定 (language)

## メールと圧縮ファイルの階層展開ライブラリ

# デ変研 MFX ライブラリ

本ライブラリは、メールからヘッダ・本文・添付ファイルの情報を取り出す【メール展開 (MX)】機能と、圧縮ファイルの展開をする【圧縮展開 (FX)】機能を統合したライブラリです。メールと圧縮ファイルを内部領域に展開し、1 つ 1 つのメールの内部情報を取り出すことができます。さらに mbox や PST (Outlook) といった複合メール形式にも対応しており、複合形式メールの最終個数を返した後、任意の 1 つのメールを取出し、展開することができます。

## 1 つのメールの展開後の結果の表示「result.txt」について

### 1 つのメール (eml 形式、RFC822) の場合

メールヘッダ

メールボディ

1 つ目の添付ファイル (dehenken\_word.doc)

2 つ目の添付ファイル (dehenken\_xlsx.zip)

この zip ファイル内には、dehenken\_xlsx.xlsx を内包

### 上記 1 つのメールを本ライブラリが展開した場合、展開後の result.txt の内容の参照方法

番号	TYPE	LEVEL	L_STRING	備考
No.0	E-MAIL	0	0	EML ファイル全体
No.1	EML HEADER	1	1	ヘッダ
No.2	EML BODY	1	2	本文
No.3	MS OFFICE	1	3	1 つ目の添付
No.4	PKZIP	1	4	2 つ目の添付
No.5	MS OFFICE VISTA XLSX	2	4.1	2 つ目の添付の展開後の内容

## ■対応 OS

**Red Hat Linux**  
AS3 / ES3 / WS3 / AS4 / ES4 / WS4 / EL5 / EL6 / EL7 / EL8

**Windows**  
2000 / XP / Vista / 7 / 8 / 8.1 / 10

**Windows Server**  
2000 / 2003 / 2008 / 2008R2 / 2012 / 2012R2 / 2016 / 2019

**Windows Storage Server**  
2012R2 / 2016

## ■対応 C コンパイラ

**Windows**  
Microsoft Visual Studio 2008 以上

**Linux**  
Gnu C Compiler (gcc)

## ■構成

**メモリ** 1GB 以上  
**HDD 利用量** 500MB 以上

※Windows は、x86 または x64 を対応に含めます。  
※Linux は、32bit 版と 64bit 版の両方を対応に含みます。  
※他の OS・コンパイラ・開発環境下で不明な点は、お問い合わせください。  
※ハードウェアの搭載メモリは推奨 2GB 以上で、メモリ量が多い方が大きな文書に対応できます。

## result.txt の全内容

NUM : 6

```
*** No.0 ***
TYPE :E-MAIL
DATA_SIZE :43557
TEXT_SIZE :0
LEVEL :0
L_STRING :0
NAME : (null)
STATUS :0
```

```
*** No.2 ***
TYPE :EML BODY
DATA_SIZE :184
DATA_ON_M :
out_dir.2228_1456884341/MEMORY/data003
TEXT_SIZE :168
TEXT_ON_M :
out_dir.2228_1456884341/MEMORY/text003
LEVEL :1
L_STRING :2
NAME : (null)
STATUS :0
```

```
*** No.4 ***
TYPE :PKZIP
DATA_SIZE :6383
DATA_ON_M :
out_dir.2228_1456884341/MEMORY/data005
TEXT_SIZE :0
LEVEL :1
L_STRING :4
NAME :dehenken_xlsx.zip
STATUS :0
```

```
*** No.1 ***
TYPE :EML HEADER
DATA_SIZE :897
DATA_ON_M :
out_dir.2228_1456884341/MEMORY/data002
TEXT_SIZE :897
TEXT_ON_M :
out_dir.2228_1456884341/MEMORY/text002
LEVEL :1
L_STRING :1
NAME : (null)
STATUS :0
```

```
*** No.3 ***
TYPE :MS OFFICE
DATA_SIZE :24064
DATA_ON_M :
out_dir.2228_1456884341/MEMORY/data004
TEXT_SIZE :186
TEXT_ON_M :
out_dir.2228_1456884341/MEMORY/text004
LEVEL :1
L_STRING :3
NAME :dehenken_word.doc
STATUS :0
```

```
*** No.5 ***
TYPE :MS OFFICE VISTA XLSX
DATA_SIZE :9363
DATA_ON_M :
out_dir.2228_1456884341/MEMORY/data006
TEXT_SIZE :162
TEXT_ON_M :
out_dir.2228_1456884341/MEMORY/text006
LEVEL :2
L_STRING :4.1
NAME :dehenken_xlsx.xlsx
STATUS :0
```