

# パソコン用地学かな漢字対応表の作成と ワープロ辞書としての利用について

野 呂 春 文・村 田 泰 章・佐 藤 岱 生・小松崎 峰 子 (地質情報解析室)  
Harufumi NORO・Yasuaki MURATA・Taisei SATO ・Mineko KOMATSUZAKI

## はじめに

パソコンワープロや 日本語フロントエンドプロセッサを用いたプログラムでの入力作業の能率は“かな”漢字変換辞書がどれだけ充実しているかで決まるといっても過言ではありません。購入したばかりのワープロソフトについている辞書では まったく不十分で 必要な単語を追加して はじめて 使いものになるのがふつうです。そのため パソコンワープロの上手な人とは 辞書登録をどんどん気楽にやる人のことだと言われるほどです。

とはいえ 辞書登録を系統的におこなうことは 大変めんどろで退屈な作業です。“かな”漢字変換辞書の充実が大事であることはわかっていても 個人では なかなか十分なものが作れないのが実状でしょう。

地質情報解析室では すでに工業技術院共同利用研究情報システム (RIPS) の FACOM メインフレーム用に “かな”漢字変換辞書を作成していただきましたので (佐藤ほか 1987) その成果をもとにして パソコンワープロの辞書を充実させるための作業に着手し このたび一応の完成をみました。ここに その作業過程使用法等を公開して みなさんの参考に供したいとおもいます。

あいのうらがた,	相浦型,	1
あいのうらそう,	相浦層,	1
あいのさわそう,	合ノ沢層,	1
あいのしまたい,	相ノ島帯,	1
あいのしまたい,	相島帯,	1
あえん,	亜円,	1
あえんがんたい,	亜沿岸帯,	1
あえんはく,	亜鉛白,	1
あおいし,	青石,	1
あおがしま,	青ガ島,	1

図1 パソコン用地学かな漢字対応表のファイル形式。

一行が一つの単語の よみがな 漢字属性を表現している。属性の項は「一太郎」特有のもので数字1は 普通名詞を意味している。他のワープロソフト 他のコンピュータで 地学かな漢字対応表を使う場合は ファイル形式の変更が必要である。

## 1. どのようなものを作ったか

作ったのは 図1のような形式の MS-DOS のテキストファイルです。一行が一つの単語の“かな”と漢字の対応を表しています。各行は まず 単語の読み のひらがな 全角のコンマ 漢字表記 全角のコンマ そして 全角の数字からできています。

このような形式にしたのは 実は ワープロソフト「一太郎 Ver. 3」の一括辞書登録機能を利用するためなのです。一太郎 Ver. 3 では 図1のような形式のファイルを用意しておけば 7千語までの単語が 一度に辞書登録できます。ただし 後でも述べますが この“かな”漢字対応表は わずかな変異で 他のワープロソフトに対応できます。決して「一太郎」専用ではありません。

“かな”漢字変換辞書の形式はワープロソフトごとにまったく異なっています。そのため すべてのワープロソフトごとの辞書を作るようなことは 一研究室の能力を越えています。しかも 特定のワープロソフト用の辞書を作った場合 そのコピーを他の人にあげることが 使用許諾権契約で禁止されていることが多いのです。せっかく辞書を作っても他の人が使えないのでは意味がありません。そこで 上のような 単純な“かな”と漢字の対応表を作成したわけです。各自 この対応表を 現在使っているワープロソフトの辞書に組み込んで 辞書を作成していただく という考えです。

この“かな”漢字対応表の単語の総数は 約2万8千語です。しかし 上で述べたように 一太郎の一括辞書登録機能を利用する場合は 7千語が限度です。一太郎用に 単語総数を7千語におさえた“かな”漢字対応表も作成しました。それは (1) 地学用語 (chigaku.dic) 5,683単語 (2) 地名 (chimei.dic) 4,141単語 (3) 国土地理院発行 図葉名 (zumei.dic) 1,858単語 (4) 山川等の自然地名 (yamakawa.dic) 5,689単語 の4つのファイルです。

この“かな”漢字対応表に収録している単語の範囲に

／あいのうらがた／相浦型／  
 ／あいのうらそう／相浦層／  
 ／あいのさわそう／合ノ沢層／  
 ／あいのしまたい／相ノ島帯／  
 ／あいのしまたい／相島帯／  
 ／あえん／亜円／  
 ／あえんがんたい／亜沿岸帯／  
 ／あえんはく／亜鉛白／  
 ／あおいし／青石／  
 ／あおがしま／青ガ島／

図2 東芝 AS3000 シリーズ unix ワークステーションで辞書登録を行う場合の“かな”漢字対応表の形式。一行が‘/’で始まりよみがな‘/’漢字そして‘/’で終わる。‘/’は半角でも全角でも良い。

については佐藤ら(1987)にくわしく述べられていますのでここでは簡単に説明するにとどめます。

まず地名です。地名には自然地名と行政地名とがあります。行政地名は大方のワープロの辞書では村のレベルまで収録されていることが多いのでこの“かな”漢字対応表では山川湖岬などのいわゆる自然地名を主にして収録しています。また国土地理院発行の地形図の図葉名も収録しています。

地球科学用語は事典の見出し語から選択しています。よほど特殊な単語は別にして国内発行の学会誌や書物にあらわれる単語のほとんどを網羅していると思います。

## 2. ワープロ辞書として利用する

パソコンワープロの辞書として“かな”漢字対応表を使う方法の一例を紹介します。ここでは一太郎の例しか紹介できませんが他のワープロソフトでも同じようなことができるかもしれません。使っているワープロソフトのマニュアルをよく見てください。

一太郎の場合一つの辞書に1で紹介した4つの“かな”漢字対応表全部を登録することはできませんので必要に応じてひとつかあるいは総数7千語を越えない範囲でいくつかを選んで辞書登録をおこなうことになります。

辞書登録の手順は以下のとおりです。

- (1) 初期化の済んだフロッピーディスクを2枚用意します。
- (2) その一枚に現在使用中の辞書(ATOK. DIC)と一太郎用“かな”漢字対応表(chigaku. dic 其他)をコピーします。

(3) もう一枚のフロッピーディスクは新辞書用ディスクとしてあとで使いますのでそのままにしておきます。

(4) Aドライブに一太郎についてくる「ユーティリティーディスク」をセットしてUT(またはut)と入力します。あとは指示にしたがって作業を進めれば新しい辞書ができあがります。地学用語(chigaku. dic)約5千7百単語を登録した場合だいたい1時間かかります。

(5) できあがった新しい辞書を現在使用しているフロッピーディスクまたはハードディスクにコピーして作業終了です。

繰り返しになりますがこの“かな”漢字対応表の利用は一太郎のみに限られるわけではありません。

「新松」等他のパソコンワープロソフト多くのunixワークステーションでは一括辞書登録機能を持っていますのでその仕様に合わせて“かな”漢字対応表の形式を変更すればそれらの計算機で利用できるでしょう。

東芝AS3000シリーズワークステーションの場合漢字コードを拡張unixコード(後述)に変換して図2のような形式の“かな”漢字対応表を用意すれば簡単に辞書の拡張ができます。

SONYのNEWS京都大学数理解析研究所から公開されているGMW(萩谷1987)のフロントエンドであるWnn(桜川1987)などでも多少変更すればこの“かな”漢字対応表が使えると思います。

“かな”漢字対応表はMS-DOSのテキストファイルですから通常のテキストエディターで単語“読み”の修正追加削除ができます。ワープロソフトでもこの作業ができないことはありませんがファイルが大変大きいのでエディターを使ったほうが良いでしょう。

“かな”漢字対応表を他のワープロソフトやシステムに対応させるには仕様に合わせた簡単なフィルターに通すだけでできます。フィルターとはファイルを読んでなんらかの加工を施してからふたたびファイルに書き出す仕事をするプログラムのことです。

一例として付録にJEF漢字コードをシフトJIS漢字コードに変換するためのフィルタープログラムを載せておきますので参考にしてください。

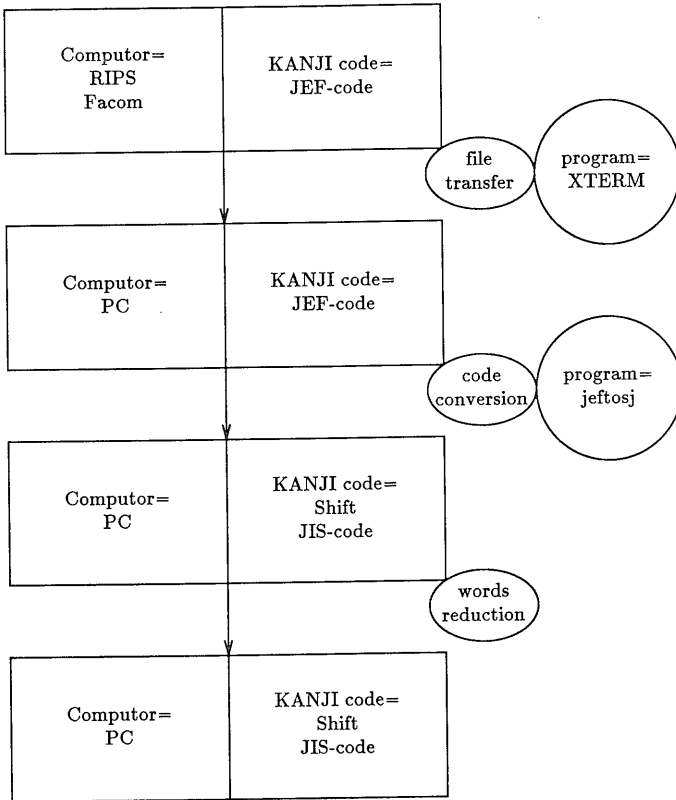


図3 パソコン用地学かな漢字対応表の作成過程。四角の箱の中は ファイルの所在と漢字コードを示す。楕円の中は作業を 円の中は そのために用いられたプログラム名を それぞれ示す。作業は それぞれ (1) RIPS上のJEFコードのファイルをPCに転送する (2) PC上でJEFコードからシフトJISコードに変換する (3) PC上で単語を削除するである。

### 3. どのようにして作ったか

“かな”漢字対応表を どのようにして作ったか 簡単に説明しましょう(図3)。

工業技術院筑波研究センターの共同利用計算機(RIPS)のFACOM M380で利用できる形の“かな”漢字変換辞書が 情報解析室の佐藤らによって 約2年がかりで作られました(佐藤ほか 1987)。

この“かな”漢字変換辞書(以降 これをオリジナル辞書と呼びます)が 今回作成した“かな”漢字対応表の母体です。オリジナル辞書の詳細については佐藤ら(1987)を参照していただくとして ここでの要点はそれが (1) 共同利用メインフレームの上にあること (2) 漢字コード体系が JEF コードと呼ばれる FACOM 独自のものであること (3) 登録単語の総数が 地学用語約1万3千 地名関連約1万5千 計2万8千という膨大なものであることです。

というわけで オリジナル辞書をパソコンに移植するには いくつかクリアすべき関門があります。どのような関門をどうクリアしたか順に述べてみましょう。

まず 最初の関門は メインフレームからパソコンへのファイル転送です。

オリジナル辞書の大きさは 約5百キロバイトあります。それだけの大きさのファイルをメインフレームからパソコン端末に誤り無しに転送するのは 簡単ではありません。通常のターミナルエミュレーション等に用いられる公衆電話回線経由の無手順垂れ流し通信ではほとんど不可能です。しかし さいわいなことに 最近 物理探査部の中塚によって ファイル転送プログラムが開発され(中塚 1987) それを利用して まったく誤りなしのファイル転送ができました。

ファイルがパソコンに移ったところで 次は 漢字コード体系の違いが重要な問題になります。

FACOM メインフレームは JEF コードという呼び名の独自の漢字コードを採用しています。一方 パソコンの MS-DOS は シフト JIS コードと呼ばれる漢字コードを採用しています。また 近ごろ さかんに導入されている unix ワークステーションでは 拡張 unix 漢字コード(16進表現は JEF 漢字コードと同じです。

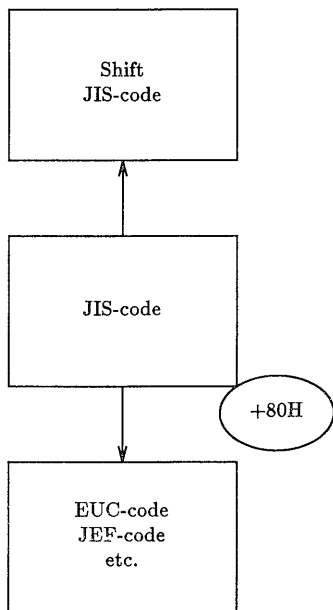


図4 各種漢字コードの間の関係。 シフト JIS 漢字コード JEF 漢字コード 拡張 unix 漢字コード (EUC コード) 等が JIS 漢字コードから派生していることをしめす。

JIS 漢字コード (2 バイト) の各 1 バイトに 80H (16 進) を加えると JEF 漢字コードその他になる。

JIS 漢字コードから シフト JIS 漢字コードへの変換は 少々こみいっているので 文献を参考にしていただきたい。

DEC 漢字コード ATT 漢字コード等も同様です) が 多く採用されています (石田 1986)。 いずれにしろ これらのいろいろな漢字コードは JIS 漢字コードをもとにして それを一定の手続きにしたがって変換して得ることができます。 したがって “かな” 漢字対応表を広い範囲の計算機に移植するには JEF コードから JIS コードあるいはシフト JIS コードへの変換が必要です (漢字コードの変換については 図4 および 付録をご覧ください

い)。 JIS コードとシフト JIS コードとの変換については 井上ほか (1987) を ご覧ください。

JEF コードから JIS コードへの変換は 各バイトから 80H (16 進の 80 10 進では 128 です) を引くだけの単純な作業です (図4)。 この 作業は パソコンでおこないました。 これで JIS コードによる “かな” 漢字対応表ができました。

JEF コードからシフト JIS コードへの変換は かなり複雑ですので 手続きの説明は省いて 付録に変換プログラムを載せておきます。 図5は “かな” 漢字対応表の一部を 16 進数の JEF コード JIS コード シフト JIS コードで表現したものです。

このようにして オリジナル辞書がパソコンの MS-DOS で利用できるかたちになりました。 最後に 辞書の大きさの問題を解決しなければなりません。

1. でもふれましたが オリジナル辞書は 地名関係 15,003 単語 地学用語 12,707 単語 合計 27,710 単語を含む膨大なものです。 一方 代表的なワープロソフトである「一太郎」では 辞書の一括登録は 7 千語までに制限されています。 単語の削除が必要です。

オリジナル辞書では 変換効率を重視して かなり長い複合単語が含まれています。 最長の単語は 「八方尾根超塩基性岩体」で 漢字 10 文字になります。 これらの 単語は 削除が可能です。 また「一太郎」付属の辞書を使って簡単に “かな” 漢字変換できる単語も削除します (“簡単に” というのがクセモノですが 変換キーを 2 3 回押せば変換できるのを “簡単” ということにしています)。

この方針で 単語の削除を行った結果 地学用語は 5 千 7 百程度まで減らすことができました。 しかし 残念ながら 地名関連用語は削除の余地が少なく 最終的に約 1 万 2 千単語となりました。 そのため 「一太郎 Ver. 3」用の “かな” 漢字対応表としては 1 で述べたとおり (1) 地学用語 (chigaku.dic) 5,683 単語 (2) 地名 (chimei.dic) 4,141 単語 (3) 国土地理院発行

	べ	っ	と	ぶ	,	別	飛	,
J I S コード :	2459	2443	2448	2456	2124	4A4C	4874	2124
J E F コード :	A4D9	A4C3	A4C8	A4D6	A1A4	CACC	C8F4	A1A4
シフト J I S コード	82D7	82C1	82C6	82D4	8143	95CA	94F2	8143

図5 パソコン用地学かな漢字対応表の一行を JEF 漢字コード JIS 漢字コード シフト JIS 漢字コードのそれぞれで どう表現しているか示す。 表記は 16 進数である。

図葉名 (zumei. dic) 1,858単語 (4) 山川等の自然地名 (yamakawa. dic) 5,689単語 の4つを作成しました。

連絡先は 下記のとおりです。

郵便番号305 茨城県つくば市東1-1-3  
工業技術院地質調査所地質情報解析室  
野呂春文 村田泰章または佐藤岱生

#### 4. 地学かな漢字対応表の提供について

以上 紹介してきました地学かな漢字対応表を 広く公開して 多くの地学関係者に使っていただきたいと 思います。 提供できるのは (1)メインフレーム用オリジナル辞書を図1の形式に変換して得られたファイル (2)「一太郎」で そのまま一括辞書登録できるように単語の削除を行ったファイル です。 ファイルの形式は MS-DOS のテキストファイルで 1メガバイトフロッピーディスクが2枚になります。 MS-DOS システムワープロソフト ワープロソフト用辞書等は 当然ながら含みません。

提供の方法 提供に際してのルール等については 所内で検討中ですので 詳しく述べられませんが この地学かな漢字対応表の入手を希望される方は 担当者までご連絡くだされば お答えします。

#### 引用文献

B. W. カーニハン D. M. リッチー (1978) プログラミング言語C (石田晴久 訳) 共立出版。  
井上尚司 大野浩之 柳染直樹 民田雅人 池田けんしろう (1987): UNIX ワークステーション NEWS, アスキー。  
石田晴久 (1986): 日本語機能も定まった最近の UNIX 事情 bit, vol. 18, no. 6, 4-9。  
萩谷昌己 (1987): GMW ウィンドウシステムについて bit, vol. 19, no.3, 4-19。  
中塚 正 (1987): ファイル転送ソフトウェア XTERM, RIPS (工業技術院共同利用研究情報処理システム) 共用ソフト登録中。  
桜川貴志 (1987): 開かれた日本語入力システム Wnn, bit, vol.19, no. 9, 13-23。  
佐藤岱生・村田泰章・青木光子 (1987): 地学かな漢字変換辞書の作成 地質ニュース No. 394, 42-49。

#### 付 録

```
/* jeftosj.c -- J E F 漢字コードのデータファイルを          *
*                  シフトJ I S 漢字コードのファイルに        *
*                  変換するためのプログラム                  *
*                  注意 2バイト漢字コード以外の処理は      *
*                  考慮していない                            */
#include <stdio.h>

main()
{
    int k, sk, d, c, cjis;
    long count=0L;
    FILE *from,*fto,*fopen();
    char buffer[20];

    /* ファイルの準備 */
    printf(" J E F 漢字コードのファイル名 : ");
    scanf("%s", buffer);
    if((from=fopen(buffer, "rb")) ==0) /* バイナリー */
        exit(printf("file %s not found\n", buffer)); /* オープン */

    printf(" 出力ファイル名 : ");
    scanf("%s", buffer);
    if((fto=fopen(buffer, "wb")) == 0) /* 同上 */
    {
        fclose(from);
        exit(printf("file %s can't open\n", buffer));
    }

    while((c = getc(from)) != EOF) /* ファイルの終わりまで */
    { /* 1バイトずつ処理する */
        count++;
    }
}
```

```

if( c == 0x40 ) ;                               /* J E F の埋め草 を */
                                                /* 捨てる           */

else if( c == 0x0a || c == 0x0d )              /* CR/LF はそのまま */
    putchar(c, fto);                            /* ファイルに出力   */

else
{
    cjis = c - 0x80;                            /* J E F => J I S の変換 */

    if( count % 2 == 1 )                       /* 第一バイトの処理 */
    {
        k = cjis - 0x20;                       /* 区数を求める */

        if( k <= 62 )                          /* 一区から六十二区なら */
        {
            if( k % 2 == 1 )                   /* 奇数区なら */
                sk = (k+1)/2 + 0x80;
            else
                sk = k/2 + 0x80;               /* 偶数区なら */
        }
        else                                    /* 六十三区以上なら */
        {
            if( k % 2 == 1 )                   /* 奇数区 */
                sk = 0xe0 + (k-62)/2;
            else                                 /* 偶数区 */
                sk = 0xe0 + (k-63)/2;
        }
        putchar(sk, fto);                      /* 第一バイトをファイルに */
                                                /* 出力する               */
    }

    else                                        /* 第二バイトの処理 */
    {
        if( k % 2 == 1 )                       /* 第一バイトが */
        {                                       /* 奇数区なら */
            st = cjis + 0x1f;
            if( st >= 0x80 )                   /* 未使用領域の */
                st = st + 1;                  /* 補正           */
        }
        else                                    /* 第一バイトが */
            st = cjis + 0x7e;                 /* 偶数区なら */

        putchar(st, fto);                      /* 第二バイトをファイルに */
                                                /* 出力する               */
    }
}

fclose(from);                                  /* ファイルを閉じて */
fclose(fto);                                  /* 後始末           */
}

```

JEF 漢字コードをシフト JIS 漢字コードに変換するためのフィルタープログラム

JEF 漢字コードで表現されたテキストファイルを シフト JIS 漢字コードのテキストファイルに変換するための フィルタープログラムです。

このプログラムは JEF 漢字コードの処理のみ行います。通常の EBCDIC コードは 正常に処理されません。

標準的な C (カーニハン リッチー 1978) で記述してありますので 現在 市場に出回っている MS-DOS 用の C コンパイラの大部分で コンパイルできると思います。unix マシンの C でも正常に動きます。