

地学用かな漢字変換辞書の作成

佐藤 岱生・村田 泰章・青木 光子 (地質情報解析室)

Taisei SATO・Yasuaki MURATA・Mitsuko AOKI

はじめに

RIPS (Research Information Processing System: 工業技術院共同利用研究情報処理システム) における日本語処理の地学用かな漢字変換辞書を作成した。その手順や経緯・考え方・問題点を述べ、今後似たような問題が生じた時に参照していただくために、記録として残しておきたい。

昨今のコンピュータの発展は、すさまじいばかりである。地質調査所でコンピュータを使用した標本管理システム GEMS が動きだしたのは、つい12年ほど前の1975年である。その時は、漢字を主体とした日本語の使用など思いも寄らないことであった。日本語タイプライタとも言えるワープロが普及しだしたのも、ここ数年のことである。地質調査所で最初のワープロの報告書原稿が提出されたのは1982年のことで、地域地質研究報告書: 信濃池田地域の地質 (1983) であろう。

かな漢字変換方式の日本語処理をしていくには、熟語とその読み方を対応させた巨大なかな漢字変換辞書があってこそ可能となる。特に地球科学関係では、山や川の名前等を含む地名が、基本的に重要な用語である。標本産地、地層名、岩体名、構造帯名、異常帯名等々は地名を冠して命名される。また、岩石名や地質現象に関する用語も一般辞書に期待することはできない。標本管理や文献情報用のデータベースだけでなく、鉱物資源データベース・岩石分析値データベース・火山データベース・地震データベース等々の開発中あるいは開発の予想されるデータベースのためにも、地学用かな漢字変換辞書を作成することは、重要な課題であることがわかる。

地質情報解析室が、この仕事に取り掛かったのは次のような経緯による。地質調査所では、資料室が担当して地質文献目録を年単位で発行してきた。この地質文献目録を計算機にデータをためこむことによって作成することは、単に印刷物を作るだけでなくデータベースからキーワードによって検索が可能になるという点で非常に意義のあることである。地質情報解析室は、1985年7月の発足早々に資料室より申し入れがあって、地質文献データベース (GEOLIS; 村田泰章 1986: 日本地質文献

データベース GEOLIS のシステム開発について) の構築について協議した。入力件数・入力項目・検索のやりかた・インデックスタームの取り方などを検討するなかで、かな漢字変換効率によって文献入力スピードが大きく規制されることが認識された。地学用かな漢字変換辞書の作成については、「全所的に必要なものだから地質情報解析室が……」ということであった。地質情報解析室は、研究部門として発足したばかりだったし、予算も人手もなかったが、引き受けることにした。文献入力数が多くなり、ユーザーにとって便利になることを期待したのである。

この辞書は、当面は RIPS の日本語機能の中で使われ、次のような用途に必要である。

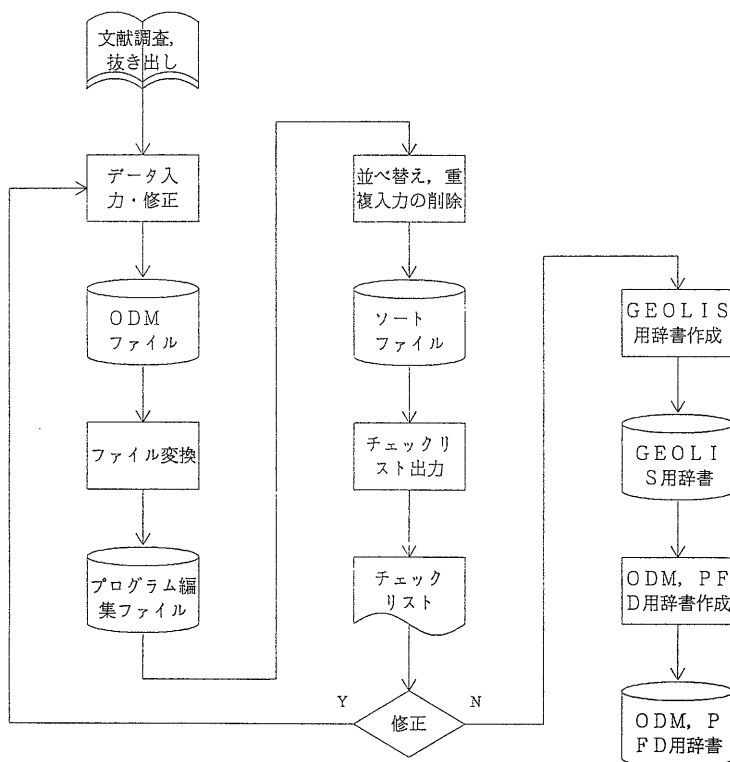
地球科学関係の文献 すなわち GEOLIS 入力
地質標本データベースへの標本産地等の入力
研究論文や地質図幅作成等の ODM 機能使用
PFD を利用した日本語データ作成

本地学用かな漢字変換辞書の開発は、富士通の機器を使用したので JEF コードで行われた。しかし JEF 拡張漢字を除けば JIS コードへの変換も容易なので、将来はパソコンやワープロでの地学用かな漢字変換辞書としての用途も開けるかもしれない。

RIPS における日本語処理

RIPS の富士通 M-380 を使用した日本語機能としては、種々のものがある。日本語の研究論文を書くなどのためのワープロ機能を持つ ODM (Office Document Manager: 文書処理システム) があるほか、日本語データの編集のために日本語 PFD (Programming Facility for Display users) がある。この PFD でデータベースに登録するためのデータを作成することもできる。また、ユーザプログラムと日本語辞書とのインタフェースも用意されていて、GEOLIS (日本地質文献データベースシステム) 等で使用されている。

これらの日本語処理におけるかな漢字変換や文法解析のために、RIPS にはいくつかの辞書が登録されている。



第1図 地学用かな漢字変換辞書作成の流れ

例えば ODM・PFD で使用する一般用語辞書 RIPS 専用のプログラムの上でのみ使用できる氏名・住所変換辞書そして漢字一文字の読み方を記憶した辞書などである。その辞書一つ一つは数万語という内容の膨大なものであるが一般的な用語が多く専門用語については私用語辞書を作成して各自で対応するという環境にある。専門用語を多用する文書作成には私用語辞書あるいは専門用語辞書が不可欠である。

RIPS では 都道府県名・郡名・市町村名までの地名はかな漢字変換されるが地学関係で頻繁に使われる大字小字名や山川の名前は辞書に入っていない。

作業の流れ

作業はまず かなと漢字が対応した適当な文献を探すところから始まった(第1図)。文献名や用語の選択については後で述べる。

データの入力作業は FACOM6658 (WDS) を RIPS の端末としてではなくローカルのワープロとして使用した。かなと対応する漢字を一対にした一覧表をワープロのファイルとして幾つか作成する。

こうして作成したファイルを RIPS の ODM のファイルに転送し結合してひとつのファイルとする。こ

の ODM ファイルをプログラムで編集可能なプログラム編集ファイルに変換しさらに読みの順番にソートして重複しているものを削除してソートファイルを作成する。このファイルからチェックリストを出力してデータのチェックと修正を行った。修正は元の ODM ファイルに対して行う。修正がすんだら再びファイル変換をしてチェックリストを出力して修正すべき所が無いことを確かめる。次にソートファイルを ODM・PFD の私用語辞書として使用する ODM・PFD 用辞書と GEOLIS 用辞書の2つの形式の辞書に変換する。

この方式で分県地図索引から山川湖名と平凡社地学事典から地質用語の入力をおこなった。この約5000語を入力した辞書が出来たところで地質調査所のおもだったところへ案内を出した。1986年1月末頃で開始から約3か月かかっている。時間はほとんどが文献からの用語の抜き出しと入力にかかっておりプログラムの作成に数日を要したほかはファイルの変換は1-2時間で終了する。

この段階で種々のアドバイスによって文部省制定の「学術用語集」からの入力の必要性を感じたとともに日常使っていないながら一般辞書には出てこない単語がたくさんあることに気が付いた。たとえば「三鉦学会・地質図幅・演旨」のたぐいである。このような用語は

ワープロ等の私用辞書から拾うことにした。

再び 地学用かな漢字変換辞書作成の流れ(第1図)の最初に戻る。次に入力した古今書院地学辞典全3巻以降は ODM の辞書参照コマンド「SHOW」を用いて既入力チェックを行いながらファイルへの書込みを行った。したがってデータ入力作業は ODM のワープロ機能を使用して行われた。

人名については 姓名とその読みが対になり しかも重複が無いような適当な文献がなかった。また 学会名簿等では重複が多く 大量のデータになることが予想された。そこで 工技院筑波研究センターの電話番号簿用のファイルから 姓と名前をそれぞれ一語とし 編集して重複をはぶいた。

1986年11月に 合計25,000語の入力を終わったので第2図の様な案内を所内各部署とおもだった使用者に出した。1987年2月には 国土地理院発行の「図名索引」からの入力も終了した。蛇足ではあるが ODM および PFD での本かな漢字変換辞書の使用方法を第3図および第4図に示す。

準拠文献

辞書には どのような用語が必要かを 広く衆知を集めて検討を行うほうが良かったのであろうが そのような時間が無かった。GEOLISの入力開始に間に合わせるために 緊急に作成する必要があった。

始めに用語を選ぶ元となった文献を 入力に使用した順に挙げておく。

分県地図索引(人文社 昭和60年度版)

平凡社地学事典(増補改訂版)

古今書院地学辞典(全3巻)

文部省学術用語集地学編(丸善)

ワープロの私用辞書

工技院筑波研究センター電話番号簿

(昭和61年1月版)

図名索引(国土地理院 昭和52・56年度版)

などである。ただし 工技院筑波研究センター電話番号簿は RIPS 上のデータファイルを編集して使用したので 従来の印刷物として連想される文献とは少々異なるかもしれない。

地名について

すでに述べたように 地球科学では地名が

1986. 11. 14.

地質情報解析室

PFD/ODM使用の皆さんへ

地学用語辞書の試験公開について

日本語処理のために、昨年から試験的に公開しております地学用語辞書を5000語から約25000語に拡充しました。内訳は、地球科学用語約11000語、山川名約8500語、地名約5500語です。下記によりご利用ください。

地学用語辞書ファイル名

'G0216.CHIGAKU.JISHO'

PFDの場合

日本語EDITに入り、コマンド行でALTQ

出てきたメニューで、上記辞書名を指定する。PF3

(END)でメニューを抜け、EDIT画面に戻る。

ODMの場合

EDIT画面直前のメニューで、私用辞書名の所に上

記辞書ファイル名を指定する。

この辞書は、資料室で入力中の地学文献データベース(GEOLIS)にも使用されております。必要な熟語やお気付きの点は、情報解析室 露(3642) 佐藤岱生まで ご連絡ください。一応の締め切りを12月15日とさせていただきます。

第2図 1986年11月中旬の案内

入力開始から約1年で登録語数は約25,000語である。ファイル名とPFD・ODMでの辞書指定方法を示している。入力開始3カ月目の最初に作成した辞書に比べると登録語数は約5倍に増えた。

-----< 辞書メニュー >-----

コマンド ==>

辞書のデータセット名を指定して下さい。

単語辞書

データセット名 ==> 'G0216.CHIGAKU.JISHO'

データセット名 ==>

標準辞書を連結 ==> 1 (1-する,2-しない)

漢字辞書

データセット名 ==>

データセット名 ==>

標準辞書を連結 ==> 1 (1-する,2-しない)

第3図 PFDでの辞書指定メニュー

データセット名の所に第2図の辞書ファイル名を書く。

ページ	本文	訂正
5	圧力形受信器	圧力型受信器
22	外型雌型	外形雌型
27	genti-taisaisakibutu	genti-taiseikibutu
33	hanō-kankei	hannō-kankei
56	kai[sai]	kai[sei]
57	kaitei-taiseikiati	kaitei-taiseikitai
58	kakeki-gun[syū]	kaseki-gun[syū]
64	kasieigan	kaseigan
69	keisitu-keugan	keisitu-ketugan
70	keityōsekisitu	keityōsitsu
71	Kenryū-rui	Kenryū-rui
88	kutyū-hōsyanō-tankō	kūtyū-hōsyanō-tankō
94	menso-kōsi	mensin-kōsi
112	汙過器	汙がJ I Sコードに無い。
113	汙速	汙がJ I Sコードに無い。
117	ryūsutu-dosyaryō	ryūsyutu-dosyaryō
118	saikkessyō[-sayō]	saikessyō[-sayō]
119	Saiyū-rui	Sairyū-rui
123	青石綿	青石綿
135	自然砂汙井	汙がJ I Sコードに無い。
139	層層注入	層々注入
147	syō-sekkaingan	syō-sekkaigan
155	tatuki	tauki
159	tihyō-tyoryoryū	tihyō-tyoryū
162	地質判読	地質判読
168	Tyōzigai-zui	Tyōzigai-rui
175	yonkutu-ōkō	shikutu-ōkō
176	yūkō-ōryokyu	yūkō-ōryoku
180	ziki-kiben	ziki-kiban

第7図 文部省学術用語集地学編の疑問点

これは文部省にお知らせしたリストの写である。地学愛好者に早く情報が伝わることを願ってあえて掲載する。日本語がかなのみで書かれた用語はチェックしていない。

読みを対応させた一覧表をワープロ(WDS)で作成する。出来た一覧表はRIPSへ転送して辞書に変換する。山川名と平凡社地学事典の入力が終わった段階で第1図の手順に従って辞書に変換した。ここで辞書に変換してしまうのはとりあえず作成した小さな辞書でその後の入力における重複をチェックするために必要だからである。このようにするとRIPSのワープロ機能であるODM中のSHOWコマンドによって辞書の内容を見ることが出来て既に登録した熟語を再び入力しないで済む。また短い熟語は長い熟語の一部となっていることが多いためにより長い熟語を入力する場合に変換効率が上がることもねらいのひとつである。

古今書院の地学辞典では当用漢字に無い漢字はかな書きするという原則がある。また鉱物名はカタカナ

で書いている。しかしこのかな漢字変換辞書ではカナはほとんど漢字に置き換えた。たとえばカクセン岩相を角閃岩相 カコウ岩は花崗岩 角レキ岩は角礫岩とした。

文部省学術用語集地学編は「複雑・難解な学術用語を整理・統一する」目的を持って昭和24年から検討を始め昭和57年に制定された地学の学術用語を収録したものである。ここでは標準的に使用される用語として特に注意深く入力した。「花こう岩・斑れい岩・隆起さんご礁」等のかな文字を含む熟語はもらさず入力した。したがってこのかな漢字変換辞書には「花崗岩」も「花こう岩」も登録されている。

私達が普通に使っている言葉で標準的な辞書にないものがかなりあると考えられる。前述の「三鉱学会・地質図幅・演旨」のほかにも「岩鉱学会・地調月報・暗灰色・淡黄色・査読」等々限りがない。これらの用語は良く使いこなされているワープロの私用辞書から選択するのが最も良いと考えられた。そこで曾屋龍典氏のODM私用辞書および地質部のワープロ私用辞書から用語を選択させていただいた。

人名について

人名は前述のように適当な文献が見当たらない。一時は地質学会年会のプログラムから講演者をリストアップしてみたが読みを付けるのもたいへんな作業で断念せざるを得なかった。次に思い付いたのが工技院筑波研究センターの電話帳でRIPS上のファイルになっているものを編集して使用することにした。名前の姓と名を切り離してJEFコード順に一列に並べて重複を取る。これで何人もいた佐藤さんや鈴木さんが1個の「佐藤・鈴木」という漢字の熟語になる。これに適当に読みがなを付けたがどうしても分からない読み方十数件については各所の庶務係に問い合わせた。地球科学のうえで有名な人名を選択すべきかもしれないが現時点ではむずかしい。上記電話帳ファイルによって地質調査所の本所の職員が1985年の時点で網羅されているにすぎない。

原典の疑問点

このかな漢字変換用学術辞書は過度の厳密性は迫及していない。すなわち常識的であると言うことで前述のように名前の読み方も常識に従って付けたにすぎない。しかし刊行されている辞典類のことになると針の先ほどのことでも取り上げたくなる。

第8図

「地名地学用語読み方辞典」の一部

地学用かな漢字変換辞書に登録された熟語の漢字からその読みを探す逆引き辞典である。冊子体で利用する。

分県地図の自然地名索引についてはすでに述べたので省略するとして平凡社の地学事典には気が付くようなミスは見付けられなかった。ただ「熔岩・熔結凝灰岩」のように「熔」の字を使用している。ここでは「溶岩・溶結凝灰岩」のように「溶」の字を使うものと両方を登録した。

古今書院の地学辞典にはごくわずかではあるがミスプリントと思われるものがある。例えば第II巻237ページの「自動酸化」が「自動酵化」に394ページの「同生団塊」が「同生団魂」になっている。また「しんせいこうしょう」の様に「ょ」が「よ」になっているものが幾つかある。

文部省学術用語集地学編で気付いた疑問点は第7図のとおりである。この図はすでに丸善経由で文部省にお知らせしたものの控えて文部省がしかるべき処置を行うものと思う。にもかかわらずあえてここに紹介したのは地学に興味を持っている方々にできるだけ早く情報をお伝えしたいからである。

この学術用語集は読み方が訓令式のローマ字で示されているのでミスプリントを発見しにくく校正ミス

凡 口 刀

14 鳳来湖	ホウライコ
鳳凰山	ホウオウサン
鳳凰山花崗岩	ホウオウサンカコウカン
鳳凰山阿部川断層	ホウオウサンアベカワタマッコウ

口 部

5 凹入角	オウニュウカク
凹地	オウチ
凹部	オウフ
出芦	テト
出芦川	テトカワ
出水山地	イスミツチ
出石川	イスシカワ
出羽	イスハ
出羽	テウ
出羽山地	テウサンチ
出羽川	イスハカワ
出来山	テキサン
出屋敷	テヤシキ
出屋敷峠	テヤシキトウケ
出砂	イテスナ
出砂島	イテスナシマ
出前沢	テマエザワ
出島	テシマ
出島	イスシマ
出雲岳	イスモタケ
出雲盆地	イスモホチ
出雲層群	イスモロクン
出森	イテモリ
出森山	イテモリヤマ
凸凹割れ口	テコホコワレクチ
凸斜面	トツシヤメン
凸部	トツフ
8 函岳	ハコタケ
函渕層群	ハコフチソクン
函館港	ハコタテコウ
函館湾	ハコタテ湾

刀 部

2 刀形類	トウケイルイ
3 刃切	ツルキリ
刃切山	ツルキリヤマ
刃石器	シンセツキ
刃状板位	ハショウテンイ
刃物ヶ崎	ハモンカサキ
刃物ヶ崎山	ハモンカサキヤマ
4 刈込	カリコミ
刈込湖	カリコミコ
刈田	カッタ
刈田	カリタ
刈田岳	カリタタケ
刈田岳	カッタタケ
刈安	カリヤス

刈安山	カリヤサヤマ
刈屋	カリヤ
刈屋川	カリヤカワ
刈寄	カリヨセ
刈寄山	カリヨセヤマ
刈場	カリハ
刈場山	カリハヤマ
切甲類	セッコウレイ
切込	キリコミ
切込砂利	キリコミシトリ
切込炭	キリコミタン
切込湖	キリコミコ
切出し	キリタシ
切目	キリメ
切目川	キリメカワ
切目崎	キリメサキ
切羽	キリハ
切取り	キトリ
切取り斜面	キトリシヤメン
切面	セツメン
切峰面	セツホウメン
切峰面図	セツホウメンズ
切断	セツタン
切断山脚	セツタンサンキョク
切雲	キリモ
切雲谷	キリモタニ
切頭	セツトウ
分子	フンシ
分子化合物	フンシカゴウフツ
分子容	フンシヨク
分子結合	フンシケツゴウ
分子磁化率	フンシシカリツ
分子熱	フンシネツ
分化	フンカ
分化作用	フンカサヨウ
分化指数	フンカシスウ
分化紋床	フンカコウシヨウ
分化説	フンカセツ
分水界	フンスイカイ
分庄	フンツ
分布	フンフ
分布様式	フンフヨウシキ
分光	フンコウ
分光分析	フンコウフンシセキ
分光光度計	フンコウコウツケイ
分光器	フンコウキ
分岐	フンキ
分岐砂嘴	フンキサシ
分岐脈	フンキミツク
分岐進化	フンキシンカ
分別	フンベツ
分別作用	フンベツサヨウ
分別結晶	フンベツツケツシヨウ
分別結晶作用	フンベツツケツシヨウサヨウ
分別晶出作用	フンベツツシヨウシュツサヨウ

32

と思われる。また和英が第1部 英和が第2部となっているがミスは両方の対応する所でまったく同じである。これは活字を拾った時に生じるミスではなくおそらくどちらか一方を電算機に入力して他方を計算機処理によって作成したものと考えられる。はじめに一方を作成した時のミスが他方に引き継がれたものと想像される。とはいえこのような電算機処理の試みは今後ますます推し進められるべきであろう。

漢字 JIS コードについて

本州最東端の岬はどこですか？ 小中学生むけのクイズだが 答えは岩手県宮古市の「鮎ヶ崎（とどがさき）」である。このことは 広辞苑にも記載されているがこの「鮎」の字が JIS コードに無い。

文部省学術用語集では 当用漢字以外の漢字の使用をごくわずかに限って認めている。その中に「汭」の字が含まれており「汭過器・汭速・自然砂汭井」に使われている。ただし 文部省学術用語集の化学編では「ろ過（汭過）」のように書かれている。

JIS コードあるいは その富士通企画版である JEF コードには「炉・爐・瀘」とも在るのに「汭」はないが「爐」を「炉」と書くなら「瀘」を「汭」と書いても良いような気がする。ちなみに「瀘」の字の「盧」は「ろ」という音を表わしているが「戸」には「ろ」の音はなく単なる記号としての役割しか無いようにおもう。「炉」は JIS 第 1 水準であり JIS 第 2 水準には「戸・舂・鈔」の文字がある。どのような理由で「汭」が採用されなかったのであろうか。

もっと困ったことは 鉱山や鉱床学で使われる「鑛（ひ）」と鉱物学で使われる「榑（せつ）」の字が無いことである。「鑛」は「鑛押・鑛先・鑛の内・鑛肌」等に「榑」は「榑石・榑面・両榑体」のように使われる。JIS を決めた国語学者の視野が狭いのか 変な文字を使う地質屋がへそ曲りなのか いずれにしても興味深いことである。

また 斑瀛岩の瀛の字も「析・励」が許されるならば「析」と書いて良いような気がする。漢字は 画数の少ない方が覚え易く 書き易く 読み易い。さらに活字や計算機のプリンタ用のフォントも作り易い。JIS コードは なぜ画数の多い字が多くて 画数の少ない文字が少ないのだろうか。不思議に思う。「花崗岩」の「崗」も意味が同じで発音も同じなだから「岡」で良いのではないかと言う人もいる。画数の少ない方が何かにつけて良いと思うのだが「花崗岩」では感じが出ないだろうか。

日本語に起因する問題

日常生活で遭遇する熟語で読み方の分からないものはそれほど多くはないだろう。また 辞書を手に入れることもむずかしくはない。ところが 地名の類はそうではなく 同じ漢字がまったく予想もつかない読み方をされることもまれではない。

本地学用かな漢字変換辞書は 熟語の読み方が分から

なければ使いようがない。使う人の思い浮かぶ読み方で たまたま目的の漢字に変換されたとき それが正しい読み方だと誤解される可能性は高いであろう。また正しい読み方が入っていたとしても間違いは起こり得る。本辞書は このような誤解から生じたり 辞書の間違いによる日本語の変化に対して責任を負うものではない。使う人の責任にかかわる部分の方が大きいように思われる。しかし 本かな漢字変換辞書を使用して間違いを見付られた方は すぐにお知らせ頂ければ幸いである。

また 漢字から熟語の読みを探せるような「地名地学用語読み方辞典」を作成した（地質調査所研究資料集登録 No.31 第 8 図）。これは 地学用かな漢字変換辞書のために作成したファイルの並べ替えを行って 熟語を漢字の部首および画数順に配列して 漢字からその読みを引くようにしたものである。同じ漢字列に対して複数の読み方があるかどうか確かめることができる。地名・人名等の 難読熟語の読みを探すのにも有効であろう。GEOLIS 入力中に地名の読み方が分からないので苦勞することが多いのである。また 冊子体となっているので一覧性があり この地学用かな漢字変換辞書にどのような熟語が登録されているかを読んで確かめることもできるだろう。熟語をどの程度の長さで切ればよいのかも分かって 漢字変換効率を上げることもできよう。ちなみに 本地学用かな漢字変換辞書に登録されている最も長い熟語は「八方尾根超塩基性岩体」で 漢字10文字である。

おわりに

地学用かな漢字変換辞書を作成した途中経過とその間に感じた問題点を書いた。この辞書に対して お気づきの点や要望などがあたららどんどん知らせていただきたい。実現の努力を惜しまない。また 当面 GEOLIS に組込んで使用するほか PFD・ODM の私用辞書として RIPS の利用者は 断わりなしに使用することができる。また ODM では私用辞書の一つしか使えないので 時期をみて RIPS の標準辞書に組み込む方を考えたい。

謝辞：この辞書の作成にあたって 環境地質部の曾屋龍典課長・地質部山田直利課長ほかから有益なアドバイスをいただいた。データ入力は青木の他に 鈴木みゆきさん・広瀬真知子さんの手をわずらわせた。地質情報解析室の花岡尚之室長および野呂春文氏には 辞書作成過程で種々の協力を頂いたが 本報告作成にあたっては 粗稿を読んで有益なアドバイスをいただいた。以上の方々に深く感謝いたします。