

GEOLIS (日本地質文献データベース) 作り 舞台裏からの報告

本 庄 時 江・武 田 福 美・菅 原 義 明 (地質情報センター資料情報課)
Tokie HONSHO・Yoshimi TAKEDA・Yoshiaki SUGAWARA

1 年程前の本誌紙上 (No. 396) で 日本地質文献データベース: GeoLiS の作成に至る歴史的経過と製作工程を紹介した。

手作業の目録作りを卒業してコンピュータによるデータベース作りに移行してから早や2年以上が経過した。初年の1986年は 5,825件 2年目の1987年には 6,987件の地学関係論文の入力を終了し この2年分はデータベースから出力して出版し 大学等の関係機関にお届けした。今 3年目を迎えて今年の入力は果して何千件となるやらと思いを馳せるこの頃である。

約2年前 当所情報解析課 (元地質情報解析室) によって開発された入力ソフトを前にして さあ入力開始! となった時のことが昨日のように思い出される。

私達図書館員6名 (室長補佐も含む) は 英文タイプライターには日常馴れ親しんではいたものの コンピュータは初体験に等しい。うっかりキーを押してしまったら せっかく作って頂いたソフトを壊してしまうのではないかと恐れおのいたのである。しかし 幸いコンピュータの扱いは入力だけなら“馴れ”で通じることががすぐ解った。

およそ千件の入力終了毎に 入力ミスなどをチェックして 検索用データベースとしてローディングすることに決めた。数回のローディングを重ねて1年分の入力を終了し 再度入力ミス探しのチェックを始めた。我が目を疑うほどの驚きであった。何たる不統一! 著者名索引では同一著者が数カ所に分れている。収録資料名一覧をみると これまた 同一雑誌名が数カ所に ひどいのは7~8ヶ所にも分れている。

何故こんなに別れてしまったのか。コンピュータは恐ろしい。ヨミのたった1字が違ってても スペースの入れ方 コンマの打ち所1つが異っても 数百字の中のたった1字が異っても 全く別のものとしてコンピュータは認識してしまうのだ。例えば「長田 芳一」長田:ナガタ・オサダ 芳一:ヨシイチ・ヨシカズ・ホウイチ この姓と名のヨミを組み合わせると6通りのヨミとなる。また「地質調査所」は チンツ┘チュウサシヨ・チンツ┘チュウサジヨ・チンツ┘チュウサ┘シヨ・チンツ┘チュウサ┘ジヨと4通りにも分れてしまう。このような具合で数限りない不統一が生れてしまったの

である。

GeoLiSの入力プログラムによる日本語の入力は ヨミとして カナもしくはローマ字を付ける 一方 GeoLiSから出力して冊子体目録の著者名索引 キーワード索引を作るためには 入力の段階で著者名 論題名 資料名などすべての日本語はヨミを付けねばならない。このヨミで私達は大いに苦しんだのである。その経験を紹介して データベースを構築されている方々への参考になれば また 何らかのアドバイスを頂ければと願っている。

1. 地名の“ヨミ”

漢字には旧字体・新字体があることは良くご存知の通りなのだが この事が私達が入力する上での支障の一つとなった。地名の表記は 公の出版物は新字体で統一されている。しかし 地名が論題名に使われる場合には必ずしも新字体に統一されてはおらず しばしば旧字体で出てくる。入力は論題名の項目には止むを得ず記載通りとするが キーワードの項目では新字体に修正して入力することに統一している。

地名のヨミを調べる方法としては 古今に出版されている各種の地名辞典 毎年改訂し出版されている「日本分県地図」がある。「日本分県地図」は市町村名にふりがなを付けてあるので私達の入力作業にたいへん役立っている。しかし 論題名ではもっと細かい地名が出てくることもしばしばである。この時には文献での調査はお手上げで 町村役場に電話をして問い合わせる。役場の方に用件を話すと とても丁寧に対応して頂けるので助かるが たまには不正確なこともある。なぜなら 連音や濁音 促音、方言によって日常使われているうちに音読みが変化してしまうこともあるらしい。この傾向は山や河川等の自然地名には特に顕著である。例えば「樽前山」のヨミは 国土地理院の資料によると「タルマエサン」道立地下資源調査所の資料では「タルマイザン」とある。どちらも公の機関の資料であるので統一するには困ってしまうのだが 独断と偏見も含めて 同じ国の機関である国土地理院のヨミを採用ことにした。

榮・榮 櫻・桜 海・海 瀧・滝 勤・勤
 廣・広 國・国 齋・斎 澤・沢 壽・寿
 淺・浅 莊・庄 總・総 藏・蔵 濱・浜

同一姓名の統一を図るために私達は旧字体・新字体のある文字は原則として新字体に統一することにした。

原則としてと書いたのは 例えば龍・竜 嶽・岳 などでは割り切るには迷いがあるって「全てを新字体とする」と断言出来ない所があるからで あと1年位は迷ってみようと考えている。

新しい著者が毎年殖え続けるであろうが 著者名に関してはヨミも表記も統一化の方向で大部分が固まりつつあると言えるまでになった。

3. 最大の悩みはキーワード付け

1) 何故 論題名からキーワードを選ぶのか？

所内に「文献データベース委員会」(委員長以下15名で構成)が昭和61年に発足 委員のうち研究者8名が作業部会委員を兼ねて各人の(学問)専門分野で GeoLiS に入力する論文の採・否の選択に当たっている。

作業部会長は 資料室の受入れ手続きを終了した資料を各委員の専門分野に基づいて振り分ける。各委員は資料室に来てこの各専門分野ごとに振り分けられた資料に随時目を通し入力論文を選択することになっている。

当所資料室は年間約13,000冊の資料を受け入れており8名の研究者が平均に資料を扱ったとして年間1人1,625冊 1カ月135冊の資料に目を通さねばならない。雑誌は1冊に少なくとも7~8論文 学会講演要旨だと数百論文が掲載されているので 1カ月135冊とはいえ 論文数にしてみると この十数倍もの論文を扱って採録論文の選択を行っていることになる。

作業部会の研究者は自分の研究のかたわら この採録作業に当たっているのでその労力は大変なもので 研究者の協力があるからこそ GeoLiS の信頼性が高められているのである。

研究者は論文選択だけでもかなりの負担となっている上に キーワード付与もということになると作業量は更にふえ 現在の数倍もの協力が必要となる。現状ではそのような体制は到底取り得ないため 次善の策として私達図書館員である入力者がキーワードを付与することになった。しかし 論題名が論文内容を必ずしも的確に表現していない場合も多く せめて内容を示している最小限のキーワードを補いたいのだが私達は地質学の専門分野に入ると いわゆる“門前の小僧”である。門前の小僧たちが出来得ることは

論題名に使われている語からキーワード(いわゆる自然語*)を抜き出すこと位がせいぜいである。

2) キーワード選定のむづかしさ

このような事情から思いがけない作業の1つとしてキーワード付けが私達入力者の肩にのしかかってきたのである。

入力開始以前の私達の打ち合わせでは 日常 地学用語には多少とも耳馴れているのでソースが作られていない現状からは“自然語”でいくしかないということになった。自然語という言葉の持つ曖昧な感じに 何となく安心感を持ってしまっただけでかなり気楽に考えてしまった。いざキーワード付けが始まったとき 私達は一様に立ち往生してしまったのである。キーワードとは何か? が皆正確には理解出来ていなかったからである。私達各自がそれぞれに曖昧なイメージを持ってはいたものの 6人のキーワードをつき合わせてみると 各人の選んだキーワードは語の区切りや集合のさせ方が各人各様であり無残な結果となった。

例えば「中新統砂子坂層産潮間帯性貝類化石群」からキーワードを採ると 二通り以上の採り方がある。

- a : 中新統・砂子坂層・潮間帯(性)・貝(類)・化石(群)
- b : 中新統・砂子坂層・潮間帯性貝類化石(群)・貝類化石(群)・化石(群)

aの採り方は 意味のある最小限の語まで細分するやり方で 各語はかなり大まかなとらえ方になる。bはかなり正確な採り方と上位概念を含めた採り方ができる。この他 a, b 組み合わせた採り方も出来る。

私達は a, b どちらの採り方が良いのか この2年間迷い続けてきたのである。

採り方の問題の他に 必要語 不要語の問題がある。例えば 地域・区域を表わす語: ~地方・~北部・~周辺など 集合を表わす語: ~類・~群など 年代や期間を表わす語: 中新世後期・上部デボン系など。

著者名の項で触れた表記の統一の問題はここにも共通している。旧字体・新字体 漢字・カナ表記の不統一 それに外来語のカナ付けの不統一 など問題点は山積し

* 自然語

非統制語・Free Term ともいう。ソースによる統制語に対する語 日常用いる語を索引語とするので利用者は思いつく限りの同義語・外来語・略語を並べて検索しなければ探し洩れを生む恐れがある。

一方 統制語は多義語の意味を1つに限定し 同義語などは1つの語のみを生かして その他は生かした語を参照するので索引語としてはすぐれている。しかしそのためにはソースの作成が前提となる。

ている。

キーワード検索をする場合 例えば“花こう岩”という表記の場合 花崗岩・花コウ岩の表記では検索の網からまれてしまって出て来ない。 シミュレーション・シミュレーション フォッサ・マグナ フォッサマグナなどの例でも同じことが起きる。 これでは検索者はあらゆる表記の種類を駆使しなければ完全な検索が出来ないので誠に不便である。 この部分に関しては 1年間の入力終了時にキーワードリストを出力してキーワードファイルだけは出来る限り統一することに努めた。

3) 問題点の解決方法は？

数限りない疑問と不安を抱えての2年間の入力であった。 1年目は入力するだけで精一杯 2年目に入って著者名・資料名の統一の土台がほぼ出来た。 ここまで

は図書館担当者の努力で解決できる範囲のものであった。

しかしキーワードに関わる問題は研究者の協力なしには解決できない事がかなりはっきりしてきた。

そこで 本年5月 「文献データベース委員会」の研究者委員にキーワードに関するアンケートをお願いした。 アンケートの内容は ①1987年に入力した論文の中から116論文題名を抽出して 私達入力者が付与したキーワードの適・否を問う問題 ②「日本地質文献目録1987」のキーワード索引約6ページ分を抽出して 語の要・不要 表記の統一化の必要性などを問う問題であった。

このアンケートに対して研究者12名から解答を頂きその結果を集約して「文献データベース委員会」で検討を行った結果をここに簡単にご報告したい。

第2表

論 題 名	入力者が付与した キーワード	アンケートでの研究者の指摘		数字は回答者数
		追 加	削 除	部分削除(下線) 分割(～と～) 修正(→)
小笠原弧火山フロントの海形海山カルデラ中に発見された熱水性硫化物	小笠原弧：火山フロント： 海形海山カルデラ： 熱水性硫化物	カルデラ1 硫化物1	火山フロント1	海形海山カルデラ1 小笠原弧1 熱水と硫化物1
人工知能的手法を用いた黒鉱生成モデルの検討(演旨)	黒鉱生成モデル	人工知能6 黒鉱1		黒鉱と生成モデル1
潜在的レアアース資源としての玄武岩類—佐賀県東松浦玄武岩類について—	レアアース資源：玄武岩類： 佐賀県：東松浦(佐賀)：		レアアース1	レアアース資源1： <u>(佐賀)</u> 2 玄武岩類1： レアアースと資源1
愛知県幡豆郡一色町佐久島の化石(海岸の転石を調査して)	愛知県：幡豆郡(愛知)： 一色町(愛知)：転石： 佐久島(愛知)：化石		転石3	～(愛知)2
津軽地方の歴史地震津波—湖沼底堆積物による歴史津波の研究	津軽地方：歴史地震津波： 地震津波：湖沼底堆積物： 堆積物：歴史津波：青森県	津波3	歴史地震津波3 堆積物3：地震津波1 歴史津波2	歴史地震津波1：地震津波1 歴史地震津波2 湖沼底堆積物1
東海地方の歴史津波—安政東海津波，下田での挙動から，ディアナ号の感じた海震(Sea Shock)	東海地方：歴史津波：津波： 安政東海津波：東海津波： 下田：ディアナ号：海震： Sea Shock		ディアナ号3：津波1 東海津波2 歴史津波1	安政東海津波1
仙台市西方，安達火山噴出物のストロンチウム同位体比	仙台市西方：安達火山： 噴出物：宮城県： ストロンチウム同位体比：			噴出物→火山噴出物2 ストロンチウムと同位体1 仙台市西方2
富士火山1707年(宝永4年)噴出物の層序にそった組成変化	富士火山：火山噴出物： 噴出物：層序：組成変化	1707(宝永4年)1 宝永1：県1名	噴出物1：層序1 組成変化1	火山噴出物2 火山噴出物又は噴出物のどちらか1つ 1
中部地方における後期更新世の古気候	中部地方：更新世後期： 古気候：			更新世後期1

①キーワードの付け方は個人差が大きい

私達入力者が付けたキーワードに対して研究者に適否を指摘して頂いた例のごく一部を 表2 に紹介する。この例からも キーワード付けがいかに個人差があるかわかるが 次の1例は最もよくその特徴を示している。(第3表参照)。

第3表 <論題名> 津軽地方の歴史地震津波—湖沼底堆積物による歴史津波の研究

付与者 キーワード	資料室	研究者(回答者)					
		a	b	c	d	e	f
津軽地方	○	○	○	○	○	○	○
歴史地震津波	○	○	○	○	○	○	○
地震津波	○	○	○	○	○	○	○
湖沼底堆積物	○	○	○	○	○	○	○
堆積物	○	○	○	○	○	○	○
歴史津波	○	○	○	○	○	○	○
青森県	○	○	○	○	○	○	○
回答者の追加		津波 (3名)		歴史地震津波	湖沼	歴史津波	

委員は各学問分野毎に選ばれていることもあって全問題に回答してもらうには無理があり 各自の分野だけを回答した委員が少なくなかった。このため 回答の集約は同一意見が多い場合でも4～5名 殆んどが1～2名の指摘であった。委員会で一堂に会して改めて検討した結果も 結論を得た事項はごく一部で 殆んどが今後の課題として持ち越されることとなった。

②ほぼ合意に達した事項

a 論文中对象となっている地域があるときは その地域の所属する県名とキーワードとして付与する。

例 論題名：岩手火山周辺の岩屑流

キーワード：岩手火山・岩屑流・岩手県

(県名の追加はデータベース構築開始時から合意を得ていたので付与されている)

b 地域・区域の部分を表わす語は原則として不要

例 ～南部 ～西方 ～地方 ～地域 ～地区 ～南縁部 ～周辺

・河川の周辺を表わす語 例例えば ～流域 ～下流 ～上流域 なども不要

・海域を表わす語は省略しない。 例例えば～沖 ～海域 ただし 海域を表す語が重複する場合はどちらかを省略する。

・例外として 東北地方 関東地方 中部地方などの表記はそのままとする。 東北 関東などの表記は東北地方 関東地方に直す。

また 日本の区域を表す中部日本 東北日本など

の表記は必要とする。

c 集合を表す語は省略する。

例 火山岩類 遺跡群

・～層群はそのままとする。

d 地震・噴火など現象の起きた年は省略する。

例 伊豆火山噴火(1986) 御岳くずれ(1984)

(1987年入力分は 論題名に年号の表記のあるものは上記のように表記したが 1988年入力分からは省略する)

e 期間・年代を表わす語は省略する。

例 新期御岳テクラ層 中新世後期

f 量・規模等を表す語は原則として省略する。

例 大型化石 大規模崩壊 含マンガン鉱物 大蛇紋岩体 微量水銀

・1つの用語として成立している語 例例えば“大地震” “微小地震” は省略しない。

③ 結論を持ち越した事項

a キーワードの採り方

キーワードの採り方はかなり個人差があることは前述の通りであるが アンケートの回答と私達入力者の経験から二つのパターンで採る傾向があると言える。(第4表参照)。

第4表

例	Aタイプ	Bタイプ
大隅花崗閃緑岩	大隅花崗閃緑岩 花崗閃緑岩 閃緑岩	大隅 花崗岩 閃緑岩
地附山地すべり	地附山地すべり 地すべり	地附山 地すべり
大山系角閃石安山岩	大山系角閃石安山岩 角閃石安山岩 安山岩	大山 角閃石 安山岩

Aタイプ：細分から上位概念までを取り入れる

Bタイプ：最小区分の語に分け並列に採る。

Aタイプは 目的が的確に絞れる 上位概念が組み込める利点がある反面 キーワードが長くなる 年間分を纏めたときに大きくなりとならないのでキーワード数が非常に多くなってしまふ。 また冊子体目録にしたときに検索には有効であるがページ数が増えてしまうという欠点を持つ。

Bタイプは 大枠で検索するには便利 キーワードが短かくて簡潔 キーワード付与が簡単 などの利点が目立つが 利用者にはデータベースレベルでの検索では組合せが自由なのできわめて有効であるが 冊子体では大枠での検索しか出来ない不便さがある。

このように利点・欠点を合わせ持っているため 委員会の検討では結論が出ず 現場である私達に当面は任せられることになった。 私達もこの2年間 両タイプを往きつ戻りつ連れ動いてはいるが 利用者のことを考えると冊子体での利用が圧倒的に多いことから 最近ではAタイプの採り方の傾向が強い。

b 町・村名に県名を付記することの是・否

数名の委員から 細かい地名は同一名が多いので県名が付記されると便利 との要望があったため1987年入力分から町・村名には原則として県名を付記した。 例えば 佐之島(愛知) のように。

この件では委員の意見は半々で当面は現行通り付記することとなった。

c その他 形容詞 前置詞的な意味を表す語 及び普通名詞の扱い

例えば 熱的構造 地質学的評価 統計的手法 塩基性火山岩 クリーブ性崩壊 火成活動 地震活動 材料評価 物理的性質 減水特性 地質調査 土壌ガス調査

この類いの語は非常に多いのだが委員の殆んどが要・不要の判断に苦しむとのことであった。 例えば“地質調査”の調査は必要 “土壌ガス調査”の調査は不要との意見もあり 要・不要の基準を見出すことが出来な

かった。 また 英単語複数の～s 例えば Rocks Terranes などの要・不要も結論が出ず 当面は今まで通り表記通りの入力となった。

d 専門用語の統一はむづかしい

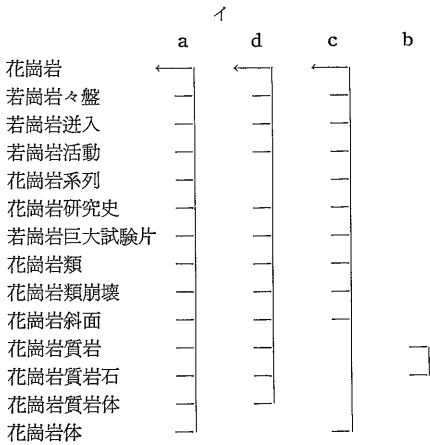
キーワードの数をなるべく少なくするには同義語を纏めるのが一番である。 アンケート調査②の結果からもキーワード統一の必要性では一致しているものの 具体的な語になると意見が大変分れた。

表5イは「日本地質文献目録1987」のキーワード索引中で“花崗岩”という語を頭に持つキーワードを並べたものである。 回答者のうち3名が“花崗岩”でのくくりを主張したが 統一すべき語の指定は必ずしも一致していない。 表5ロも“葛根田”を頭に持つ語での質問で これも3人3様の回答であり 表5ハも同じ結果であった。

4. 今後の課題

入力開始3年目を迎えたこの4月に 当所情報解析課 村田泰章技官によって 入力ソフトの改良が行われこれまでは欧文か和文のどちらか一方しか入力できなかった 著者名・論題名・資料名・発行者名が これらを含めた

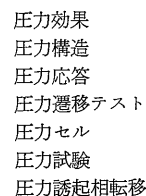
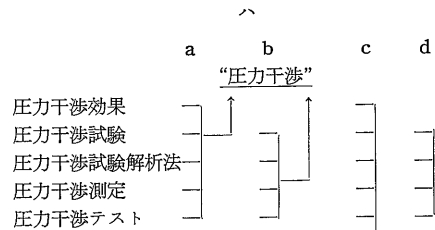
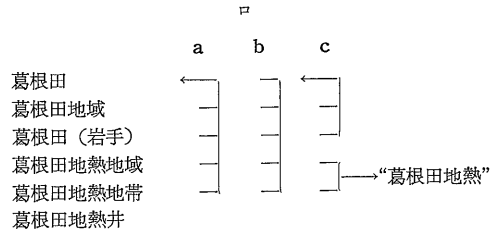
第5表



a, b, c……: 回答者

: 矢印のキーワードに統一

: 統一の必要性あり 但し
: キーワードの指定は別途検討



第6表

論文題名	GEOLIS のキーワード	GEOREF のキーワード	
Strike-slip fault (Nayorogawa Fault) in northern Hokkaido	STRIKE-SLIP FAULT NAYOROGAWA FAULT HOKKAIDO NORTHERN HOKKAIDO	Japan structural geology faults displacements strike-slip faults Far East Asia Nayorogawa Fault Hokkaido	Hidaka Belt folds shear deformation right-lateral faults subduction oceanic crust crust Eurasian Plate
Conical Folds in the Hitoyoshi Bending, South Kyushu, Formed by the Clockwise Rotation of the Southwest Japan Arc	CONICAL FOLDS HITOYOSHI BENDING SOUTH KYUSHU CLOCKWISE ROTATION SOUTHWEST JAPAN ARC	Japan structural geology tectonics Far East Asia Hitoyoshi Bending Kyoshu plate rotation Nansei-Shoto Arc	Hokusatsu Bending Nojiri Bending Shimanto Supergroup Hyuga Group folds Miocene Neogene Tertiary

全項目で欧和両語での入力が可能となった。また冊子体目録作成のための出力ソフトも少しづつ改良されている。

著者名・資料名の扱い・統一もほぼ見通しがついたが最後に残るのが キーワード付与の問題である。

ちなみに 米国 AGI (The American Geological Institute) 製作の地球科学文献データベース: Georef と当所の GeoLiS の双方に掲載された同一論文のキーワードを比較してみると 表6に示す通り非常に大きな差がある。

Georef は “Georef Theasaurus and Guide to Indexing” というシソーラスが作成されており 現在 4th Edition まで版を重ね A4判 512 ページという独自の辞書からキーワードが付与されている。

当所の GeoLiS のキーワードの充実を図るには 当所独自のシソーラスを作るか 研究者によるインデクサー・グループを組んで キーワード付与の作業に当る必

要があると考え、しかし 当面はこれら改善策の実行は不可能なことから 現在のやり方を継続せざるを得ない。

「文献データベース委員会」では データベース構築に関わる問題点の認識が やっと一致した段階で 改善策を練るのは今後に持ち越すことになったが 当面 改善の策として1988年分の入力終了した時点で キーワード索引を出力してキーワード統一化への検討の手がかりを考えることとなった。

このように 多くの研究者の協力を得ながら構築しているデータベース: GeoLiS は まだ所内利用だけに限定されている。 所外の方々には冊子目録での利用しかできないのが 精魂込めて製作に取り組んでいる私達にとって残念でならない。 地学関係の仕事に携わる多くの方々に 一日も早く利用して頂ける日が訪れるのを 私達担当者は願って止まない。