

万葉集の検索システム

蓼沼良一, 志村栄一 (山梨大学)

1. まえがき

原始日本語の母音は8通りで、古代日本語の文書、古事記、万葉集などには、まだ、その名残りがあるという。したがって、古代の文書を扱うには、この点を考えておく必要がある。

この研究は、万葉集の検索システムの開発を試みたもので、開発には、山梨大学の乙号機 F A C O M - S C R O - 30 を F O R T R A N によって行う。しかし、これはカナ文字が扱えないるので、適当な時期に I N T E L 8080 を用いた研究室の計算システムに切換えるつもりである。

システムの開発には、まず、その目的を明らかにすることが必要である。検索システムとしては「問い合わせ」と「索引書作り」とが対象になろう。索引書としては K W I C, K W O C などか、よく知られている。問い合わせでは、自然言語によるもの、記号化した人工言語によるものなどがあろう。そして、答には、とへようなものを作ろかと/or点を決めねばならない。

開発には、目的を明らかにした上で、とりかかるべきであるか、ここでは、これを後日にまわし、文字コードの問題と作者名の問題とを先に検討した。

2. 文字の種類

まず、F O R T R A N を用いて開発するという前提に立ち、万葉集をどのように表わしたらよいかという問題について考える。

万葉集を扱うには、当時の日本語の音を表わす文字と、中国式の音を表わす文字とか必要である。さらに、数字を利用しての表記の簡単化と、いくつかの記号とかいる。すなわち、表記に使う広い意味の字は、いくつも文字と数字・記号となる。

<字> → <文字> | <数字> | <記号>

<文字> → <日本音> | <中国音>

ここに、→は B a c K u l s 記号の "≡" と同じ意味に使ってある。

日本音は、母音だけか、子音母音とからめていて、子音だけのものはないと言える。母音は、古代には甲乙にわかれていた。これを甲母音、乙母音といふことにしよう。

<日本音> → <母音> | <子音> | <母音>

<母音> → <甲母音> | <乙母音>

甲母音アイウエオはローマ字表記法を用いて表わす。

<甲母音> → A | I | U | E | O

乙母音イエオについては、1文字では適當な文字がないので、:をつけて二字で表わした。

<乙母音> → I: | E: | O:

ローマ字表記法にはヘボン式、日本式などがあり、これに応じて、子音の表わし方も、いくつかがある。ここでは訓令式を用いる。

<子音> → K | G | S | Z | T | D | N | H | B | P | M | Y | R | W

訓令式の利点は、括弧、連濁などの扱いが簡単になることの他に、古代日本語の音は、訓令式に近い発音であったとすることもある。

中国式の音には日本語の音の他に、拗音、促音、撥音が加わる。

<中国音> → <日本音> | <拗音> | <促音> | <撥音>

= れらは、訓令式によれば、

<拗音> → <子音> + <母音>

<促音> → <子音> <子音> <母音>

<撥音> → N | N'

母音としては甲母音だけで充分で、乙母音を考える必要はない。

訓令式では、促音を表わすのに、同じ子音を重ねて用いる。これは訓令式の欠点の一つで、この改良案としてQを用いる方が提案されている。ここでは、Qの代りにCを用いることにする。

<促音> → C <子音> <母音>

撥音にNを用いると、母音の前ではナ行と混同する。これを区けるため、母音の前ではN'を用いる。したがってZ通りの表記が必要になる。この複雑さを区けるため、撥音は常に二字N:で表わすことにする。

<撥音> → N:

:を用いても、乙母音と混乱する恐れはない。

数字、記号はつきの通りとする。

<数字> → 0 | 1 | ... | 8 | 9

<記号> → , | . | □ (□は空白)

3. 文字コード

前節の検討の結果、才葉集を記すのに11万字の種類は、FORTRANの文字で英字20、数字10、記号4で、計34コである。したがって、これらを表わすには6ビット以上が必要になる。Z号機では16ビット、8080では8ビットを単位としているので、8セットを単位とした文字コードを作る。

(1) 1字1バイト方式

最も簡単なものは、FORTRANの1字を1バイト(8ビット)に表わす方法である。これは簡単であるが、つきの欠点がある。

(ア) 記憶容量の無駄

扱うべき文字の種類は34で、5ビットで表わせる。これに8ビットを使うのであるから、約3ビット、40%の無駄となる。

(イ) 記事修正の困難

たとえば、甲母音を乙母音にかえること、あるいは乙母音を甲母音にかえるには:を加え、あるいは除かねばならない。記憶装置のなかで、このようにずらすことは計算機の苦手とする仕事の一つである。しかし、このような修正は、かなり生ずると思える。

(エ) カナ文字1バイト方式

前案の欠点を直す一方法は、カナ1字を1バイトに表わす方法である。日本音では母音8コで3ビット、子音14コで4ビットあればよい。したがって、8ビットをわり、5+3、あるいは4+4方式ができる。しかし、5+3方式では、母音のいらない撥音N:が扱えなくなる。そこで、4+4方式について考えを進めよう。

つぎに、中国音を考えると、拗音、促音、擦音が加わるから、4+4方式で扱えるかどうかをしらべよう。

まず、拗音については、

〈子音〉Y + 〈母音〉

〈子音〉+ Y 〈母音〉

とする2つの方法があろう。前の案は、〈子音〉Yを子音として扱う方法である。この案において、拗音記号Yに結びつく子音は、ヤ行Yとワ行W以外のすべてである。したがって、子音の種類は $14 + 12 = 26$ で4ビットでは表めせない。

これに対し、後の方式は、Y 〈母音〉を母音として扱うものである。いわば、拗母音を設定し、ローラ語の軟母音のように扱うものである。

〈母音〉 → 〈甲母音〉 | 〈乙母音〉 | 〈拗母音〉

拗音記号Yのつく母音はA, U, Oの3つだけであるから計11コで、4ビットで表めせる。

促音についてしらべよう。促音記号Cの扱いは、

C 〈子音〉 + 〈母音〉

とする方式と、前の母音とあすび

〈母音〉 C + 〈子音〉 + 〈母音〉

とする方式とがある。

促音記号Cは、すべての子音とあすびつくわけではなく、K, S, T, Pの4コに限られる。しかし、CK, CS, CT, CPという子音を設けると、子音は $14 + 4 = 18$ となり、4ビットでは表めせない。

オフの〈母音〉Cという方式では、Cはすべての母音と結びつくが $8 + 5 = 13$ で、4ビットで表めせる。

〈母音〉 → 〈甲母音〉 | 〈乙母音〉 | 〈促母音〉

しかし、拗母音3コは扱えるが、拗促母音3コまでは扱えない。この問題については、あとで再び考えよう。

擦音N: について考えよう。

N + :

N: + ④

とする2つの方式がある。前のものは、子音Nに母音:がついたとする考え方、後のものは子音N:に母音なしとする考え方である。

後の方式は、子音を増やすという方式であるが、4ビットを越えることになり4+4方式では扱えなくなる。したがって、前の母音を増やす方式によらざるをえない。

以上の音の他に、必要と考える子音として

KW, GW, TS, ジズ, ハ

がある。ハはHWとするとしても、子音の数は増やせないから、母音として扱う方式を考えねばならない。これは、擦音記号:が子音Nにしかつけないものと考えたように、特別な子音、あるいは子音にはつけない母音とし、不完全母音と名付けよう。

〈不完全母音〉 → W | S | Z | : | C

そうすると、つぎのようなコード案ができる。

I. 第1案

〈文字コード〉 → 〈子音部〉 〈母音部〉

〈母音部〉 → 〈母音〉 | 〈不完全母音〉

〈母音〉 → 〈甲母音〉 | 〈乙母音〉 | 〈拗母音〉

この方式によるコードの具体案の1例をあげよう。

コード	子音部	母音部(Ⅰ案)	母音部(Ⅱ案)	数字記号
0	母音	:	:	0
1	K	A	A	1
2	G	I	I	2
3	S	I:	I:	3
4	Z	U	U	4
5	T	E	E	5
6	D	E:	E:	6
7	N	O	O	7
8	H	O:	O:	8
9	B	YA	AC	9
10	P	YU	IC	>
11	M	YO	UC	・
12	Y	W	EC	□
13	R	S	OC	
14	W	Z	W	
15	記号	C	Y	

II 第Ⅱ案

母音ヒ促母音を加え、拗母音を除くとオフの案ができる。

〈母音〉 → 〈甲母音〉 | 〈乙母音〉 | 〈促母音〉

この具体案をあげると、つきのようになる。子音部には変化がないから、母音部のみをあげる。

この案ではTS, オズの音は表わせないことになる。

この案による表記例をあげよう。子音部△, 母音部△からできる文字コードを(△, △)で表わす。

	Ⅰ	Ⅱ
花	(H,A)(N,A)	(H,A)(N,A)
神	(K,A)(M,I:)	(K,A)(M,I:)
絵画	(K,W)(-,I)(G,W)(-,A)	(K,W)(-,I)(G,W)(-,A)
京都	(K,YO)(-,O)(T,O)	(K,Y)(-,O)(-,O)(T,O)
国家	(K,O)(-,C)(K,A)	(K,OC)(K,A)
出張	(S,YU)(-,C)(T,YO)(-,O)	(S,Y)(-,UC)(T,Y)(-,O)(-,O)

ここに、一は母音の印を示す。

ここにあげた例だけでは、両者の優劣は決めていくので、音便の立場からくらべてみよう。口語には、動詞の活用に際し、音便がある。

書き立 → 書いた (イ音便)

咬み立 → 咬んだ (撥音便)

勝ち方 → 勝つ方 (促音便)

これを、西方式で表わしてみよう。

	I	II
(K,A)(K,I)(T,A)	(K,A)(-,I)(T,A)	(KA)(-,I)(T,A)
(K,A)(M,I)(T,A)	(K,A)(N,:)(D,A)	(K,A)(N,:)(D,A)
(K,A)(T,I)(T,A)	(K,A)(-,C)(T,A)	(K,A,C)(T,A)

この例のように、II案では、促音便に対して語長が変ると、語形の変化が面倒であるという欠点がある。また、TS, ハズの扱いなどもあり、I案の方がすぐれていいことわかる。

4. 作者名の扱い

万葉集のもつ情報の量をしらべよう。岩波版日本古典文学大系万葉集1~4巻をもとに概算すると

4巻 × 2000 頁 × 400 字 × 0.5 = 24万字
である。この式の系数 0.5 は、同書が原文と訓み下し文とから構成してあるためのものである。万葉集中は、約 4500 首の歌があるが、これを全部対象として

$$31 \text{字} \times 4500 \text{首} = 14 \text{万字}$$

となる。これは、すべてが短歌だけとしての計算であるから、実際はもっと多くなる。検索の効率をあがるためには、このような多量の情報を、如何に圧縮するかということになる。ここでは、まず作者名について検討してみよう。

作者名のような長さの変るものを探うには、(1) 可変長方式、(2) 固定長方式との二つの方法がある。記憶容量、検索時間を物指しとして、これらを評価比較を行う。

(1) 可変長方式

作者名を可変長のまま扱う方法について考えよう。いま、作者名の平均長を 10 字と仮定すれば、この総字数 K は

$$K = 10 \text{字} \times 4500 \text{首} = 4.5 \text{万字}$$

となる。これが、必要記憶容量の目安になる。

つぎに、ある作者の全作歌を取り出すに要する検索時間をしらべる。この時間は、作者名をしらべるに要する字数で評価できよう。名前の違うことを知るには平均長の半分、すなはち 5 字しらべればわかる。したがって、全部の検索には

$$H = 5 \text{字} \times 4500 \text{首} = 2.3 \text{万字}$$

の比較がいる。この K と H を比較の物指しとする。

(2) 固定長方式

作者名を固定長にするにはいろいろな方法があるが、辞書を利用するのが簡単である。名前の辞書に要する記憶容量は

$$10 \text{字} \times 500 \text{人} = 0.5 \text{万字}$$

となる。ここでは、作者は 500 人として計算した。

作者は 500 人であるから、この番号は 2 バイトで表わせる。したがって歌集中の作者名は

$$2 \text{字} \times 4500 \text{首} = 0.9 \text{万字}$$

となる。よって、総字数Kは、これらの和となる

$$K = 1.4 \text{ 万字}$$

名前の検索には、辞書の検索が加わる

$$\text{辞書検索: } 5 \text{ 字} \times 250 \text{ 人} = 0.13 \text{ 万字}$$

$$\text{歌集検索: } 1 \text{ 字} \times 4500 \text{ 首} = 0.45 \text{ 万字}$$

$$\text{計 } H = 0.6 \text{ 万字}$$

(3) 語方式

作者名を見ると、同じ語が何人かの名前に現われる。たとえば、大伴家持、大伴東人、大伴旅人の如くである

これから、語辞書を用い、語番号の列として名前を表わす中間の方式を考えられる。この方式について述べよう。

作者名は3語、1語は4字でできると仮定しよう。この仮定で、作者名の述べ語数は

$$3 \text{ 語} \times 500 = 1500 \text{ 語}$$

となる。異なり語は、この半分とすれば 750 語となるこれにより字数は

$$\text{語辞書: } 4 \text{ 字} \times 750 \text{ 語} = 0.3 \text{ 万字}$$

$$\text{歌集: } 2 \text{ 字} \times 3 \text{ 語} \times 4500 \text{ 首} = 2.7 \text{ 万字}$$

$$\text{計 } K = 3 \text{ 万字}$$

となる。

検索時間は、語辞書を3回までして作者名が決まり、歌集の検索には、語番号が2字であるから、つきのようになる。

$$\text{辞書検索: } 2 \text{ 字} \times 375 \text{ 語} \times 3 \text{ 回} = 0.225 \text{ 万字}$$

$$\text{歌集検索: } 3 \text{ 字} \times 4500 \text{ 首} = 1.35 \text{ 万字}$$

$$\text{計 } H = 1.6 \text{ 万字}$$

となる。

以上の計算結果を表下すると、つきのようになる

	可変	固定	語
容量	4.5	1.4	3.0
時間	2.3	0.6	1.6

この表によれば、固定方式が最もすぐれ、語方式、可変長方式の順になつていて。しかし、データの作成誤りの発見などの立場からは、語方式がすぐれてい。たとえば、乙母音を用いる ABE: (阿倍) を、甲母音にして ABE の誤りは、語辞書方式により、容易に発見できる。

5. 作者名表の作成

前節の検討結果にもとづき、語辞書を利用して作者表を作る作業を行つた。使用した計算機は乙号機 FACOM 270-30 で、FORTRAN を用いた。作業は、次のよう順序で行つた。

(1) 作者カードの作成

資料乙の時代順作者人名録をもとにし、前記資料 1 を参照して作者カードを作成。作者カードの構成は、作者番号、作者名、必要なときは注釈を加えてある。

〈作者カード〉 → 〈作者番号〉 〈作者名〉 [* <注釈>]

ここで、〔……〕は、序くてもよい項目を意味する。

資料2の人名録は、ほぼ50音順に配列してあるの下、作者番号も、その順になつてある。

(1)	101 AUZI NO: OKISIMA
(2)	102 AGATA NO: INUKAHI NO: WOTOME
(3)	103 AGATA NO: INUKAHI NO: HITOKAMI
(4)	104 AGATA NO: INUKAHI NO: MIIYO:
(5)	105 AGATA NO: INUKAHI NO: MOTIWO
(6)	106 AGATA NO: INUKAHI NO: YOSIWO
(7)	107 AKI: NO: OHOKIMI
(8)	108 AKI: NO: WOSA NO: OBITOMARO: * SURUGA NO: SAKIMORI
(9)	109 ASAKURA NO: MASUSHITO * KAMITUKE: NO: SAKIMORI
(10)	110 ASADA NO: YASU
(11)	111 ASUKA NO: OHOKIMI
(12)	112 ASUKA NO: NADO:MARO:
(13)	113 ATUMI NO: OHOKIMI
(14)	114 ATO NO: TUSITARI
(15)	115 ATO NO: TOBIRA NO: WOTOME
(16)	116 ABE: NO: OKINA
(17)	117 ABE: NO: OKIMITI
(18)	118 ABE: NO: WOMINA
(19)	119 ABE: NO: KOOHODI
(20)	120 ABE: NO: SAMIMARO:
(21)	121 ABE: NO: TUGIMARO:
(22)	122 ABE: NO: TOYOTUGU
(23)	123 ABE: NO: BOU
(24)	124 ABE: NO: BOU * ABE: NO: TUGIMARO: NO: ZINAN:
(25)	125 ABE: NO: HIRONIHA
(26)	126 ABE: NO: MUSIMARO:
(27)	127 AMA NO: BOU
(28)	128 AMA NO: INUKAHI NO: WOKAMARO:
(29)	129 AMU NO: MOROTATTI
(30)	130 AHATA NO: BOU
(31)	131 AHATA NO: OHOKIMI
(32)	132 AHATAME NO: WOTOME
(33)	133 ARAUDI INASIKI
(34)	134 ARTMA NO: MIKO
(35)	201 IHOMARO:
(36)	202 IKUSA NO: OHOKIMI
(37)	203 IKUTAMABE NO: TARUKUNI * TOROTUAHUMI NO: SAKIMORI
(38)	204 IKEDA NO: BOU
(39)	205 IKEBE NO: OHOKIMI
(40)	206 ISIKAHANA NO: OKINA
(41)	207 ISIKAHANA NO: IRATUME
(42)	208 ISIKAHANA NO: HUHITO
(43)	209 ISIKAHANA NO: OHOBIA
(44)	210 ISIKAHANA NO: WOMINA * HUJIHARA NO: SUKUNAMARO: NO: IUMA
(45)	211 ISIKAHANA NO: WOMINA * OHOTU NO: MIKO NO: MIYA NO: MAKATATI
(46)	212 ISIKAHANA NO: KAKE NO: WOMINA
(47)	213 ISIKAHANA NO: KIMIKO
(48)	214 ISIKAHANA NO: TARUHITO
(49)	215 ISIKAHANA NO: TOSTTART
(50)	216 ISIKAHANA NO: BOU * MIYAMARO: KA KIMIKO KA
(51)	217 ISIKAHANA NO: BOU
(52)	218 ISIKAHANA NO: HIRONARI
(53)	219 ISIKAHANA NO: MIMITI
(54)	220 ISOUDI NO: NORIMARO:
(55)	222 TSONOKAMI NO: KATUWO
(56)	223 ISONO:KAMI NO: BOU * UTOMARO: KA

(2) 語辞書の作成

作者カードから、作者名を構成する要素より語をとり出し、これを並べかえる。そして、語番号をつける。こうしてできたのがつきの表である。

この時点では、文字コードの編成はしていないので下のORTTRANの文字順に並べてある。また、語番号が100以下のは、数字記号とした。ニラして得た要素より語の总数は、ほぼ600語である。

KOTONARTSU NO ME LIBIKI (GOSU= 597 GU)	
1 *	810 HATA
2 1	811 HATIMARO:
3 2	812 HATORIBE
4 3	813 HATUSE
5 4	814 HEGURI
6 5	815 HEGURIUDI
7 5NEN:	816 HIBITO
101 ABE:	817 HIMATURI
102 ADUMABITU	818 HIMEMIKO
103 ADUMAMARO:	819 HINO:KAMI
104 AGATA	820 HINO:KUMA
105 AHAMARO	821 HIOKI
106 AHATA	822 HIROKAHA
107 AHATAME	823 HIROKATA
108 AHUMI	824 HIROME
109 AKAHITO	825 HIROMIMI
110 AKAMARO:	826 HIRONAHA
111 AKIMOTI	827 HIRONARI
112 AKINIHA	828 HIRONIHA
113 AKI:	829 HIROSE
114 AMA	830 HIROSIMA
115 AMU	831 HIROTARI
116 ANAGA	832 HIROTU
117 ARAMIMI	833 HIROTUGU
118 ARAUDI	834 HITATTI
119 ARIMA	835 HITO
120 ASADA	836 HITOKAMI
121 ASAKURA	837 HITOMARO:
122 ASAME	838 HITONA
123 ASUKA	839 HITONUSI
124 AIO	840 HITOTARI
125 ATUMI	841 HITOWOSA
126 AUZI	842 HITUGIMIKO
127 AYUMARO:	843 HITUZI
201 BEN:KI	844 HIZEN:
202 BOU	845 HOHUSI
203 BUNIN:	846 HOUZI
401 DOUYOU	847 HOZUMI
501 ENITATI	848 HUBITTOUDI
502 ENO:WI	849 HUBUKI
503 EUZI	850 HUDIHARA
701 GEN:MEI	851 HUDIHARABE
702 GEN:SYOU	852 HUDIWITI
703 GITUU	853 HUHIITO
704 GIYAUMONI	854 HUKUSI
705 GO:NU:TANIWOTTI	855 HUMI
706 GUWANIGOZI	856 HUMIMOTI
707 GUWANINI	857 HUMINARI
801 HADA	858 HUMUYA'
802 HAHA	859 HUNADU
803 HAKURI	860 HUNE
804 HAKUTUU	861 HURU
805 HANISI	862 HURUHITO
806 HANIUDI	863 HUSASAKI
807 HARIMA	864 HUSTMARO:
808 HASETUKABE	865 HUUSI
809 HASIHITO	901 IHANO:HIME
	902 IHE

(3) 作者名表

再び作者カードを読み直し、語辞書の番号によって作者名を表わしたのが、次の表である。

この表からわることは、駿河の防人、上野の防人などという注釈を除けば、殆んどが了語で表わせる。

SAKUSYAMEI NO GOKUBANGO LIST

101	= 126,1422,1533,
102	= 104,1422, 921,1422,2329,
103	= 104,1422, 921,1422, 836,
104	= 104,1422, 921,1422,1337,
105	= 104,1422, 921,1422,1354,
106	= 104,1422, 921,1422,2526,
107	= 113,1422,1507,
108	= 113,1422,2323,1422,1501, 1,1956,
* *	1422,1910,
109	= 121,1422,1312, 1,1119,1422,1910,
110	= 120,1422,2515,
111	= 123,1422,1507,
112	= 123,1422,1401,
113	= 125,1422,1507,
114	= 124,1422,2060,
115	= 124,1422,2048,1422,2329,
116	= 101,1422,1532,
117	= 101,1422,1530,
118	= 101,1422,2320,
119	= 101,1422,1162,
120	= 101,1422,1914,
121	= 101,1422,2068,
122	= 101,1422,2069,
123	= 101,1422, 202,
124	= 101,1422, 202, 1, 101,1422,2068,
* *	1422,2603,
125	= 101,1422, 828,
126	= 101,1422,1360,
127	= 114,1422, 202,
128	= 114,1422, 921,1422,2318,
129	= 115,1422,1352,
130	= 106,1422, 202,
131	= 106,1422,1507,
132	= 107,1422,2329,
133	= 118, 919,
134	= 119,1422,1324,
201	= 903,
202	= 908,1422,1507,
203	= 909,1422,2030, 1,2049,1422,1910,
204	= 906,1422, 202,
205	= 905,1422,1507,
206	= 923,1422,1532,
207	= 923,1422, 922,
208	= 923,1422, 853,
209	= 923,1422,1503,
210	= 923,1422,2320, 1, 850,1422,1955,
* *	1422,2073,
211	= 923,1422,2320, 1,1522,1422,1324,
* *	1422,1338,1422,1303,
212	= 923,1422,1112,1422,2320,
213	= 923,1422,1142,
214	= 923,1422,2029,
215	= 923,1422,2060,
216	= 923,1422, 202, 1,1339,1101,1142,
* *	1101,
217	= 923,1422, 202,
218	= 923,1422, 827,
219	= 923,1422,1326,
220	= 927,1422,1419,
222	= 925,1422,1137,
223	= 925,1422, 202, 1,1539,1101,

あとがき

語辞書を利用するものは、かなり有利な方法と思うが、駿河の防人などという句
辞書の用意も有効かも知れない。これからは更に、統計的な検討が必要であろう。

ありがとうございました。研究上種々助けて下さった本講座伊藤誠助教授、アロケラムやテ
ータの作成をして頂いた武藤厚子技官に感謝します。

参考文献

- (1) 岩浪版 万葉集1-4 日本古典文学大系
- (2) 武田祐吉校註 万葉集上、下 角川文庫