

日本語のコンコーデンス

坂本義行(電子技術総合研究所)
岡本哲也(電気通信大学)

1. はじがき

自然言語の調査、研究を目的とし、その構文あるいは意味の解析を行なうための資料作成の自動化の1つとして、ALPS(Automatic Language Processing System)-KWDI C Index-Fileシステムをすでに開発し、これに関する報告は行なっている。⁽¹⁾⁽²⁾この既存のシステムの汎用化ならびに実用化を目的として、漢字かな混りテキストを探本としたコンコーダンスの入出力形式、蓄積、検索、統計、編集の各プロセッサを総括したシステムを作成し、Contextの自動処理の実験を行なった。以下に、本システムの概要、実験結果ならびに利用目的等について報告する。

2. 原デーティの入力方式

電子計算機で情報処理を行なうには、1次情報である原データをどのように形で計算機の認識系と結合させねば問題となる。

ニヤテは、

使用文字

符号变换

デーラの構造

の3兵はついで述べる。

使用文字 計算機へ日本諸を入出力させるには、種々の方法と問題点が報告されており。本システムでは、HITAC-8400 計算機システムで使用するといふことを方針の目的としたものである。字種として、かな、ローマ字書すによる翻字入力も可能であるが、原テータが漢字か混りである場合から、漢字を直接入力することにより、打鍵における視認性の良さ、同音異義語の識別、テータの検索(処理の中間結果)等に大きな利点をもつてゐる。当研究室には、漢字テレタイプライタ(漢字レシピ:西村氏により詳細な報告がなされてゐる)があり、漢字1字に対して、8単位2列の符号が紙テープ媒体に出力される。これは、計算機とは offline で、紙テープ読取機から入力される。この漢字の符号系は、計算機内での符号変換、行書きの処理、KWIC Index File の配列処理等で大きな特徴を有してゐる。漢字、カタカナ、ひらがな、ロシア文字、ギリシア文字、特殊記号等約2,300種が収められており、文字種が豊かな点、漢字を素人が打鍵するのに便利な音韻別配列ヒューリズムに考慮が払われている。

符号実操 翻字された漢テレコードは、オフ回
1.示すように、8単位2列で構成されており、文字
情報は、各列のオフ1～6単位で表現され、オフ単位
は補助用、オフ単位は奇数パリティに用いられて
る。HTTAと内へ1文字のオフ1列、オフ2列を各々
1バイトとして詰め込まれた場合、EBS DJK の内
部コードと一緒に致しなり。これは、COBOLによる
ログラム作成での煩雑さ、データの検査のとき、ア

第2圖 級子-7°エウ翻字

リニア上に現われない、後述する配列等の問題を解決するため、オフ表に示す符号変換を行なう。

第 1 表 漢テレ - HITAC コード変換表

漢テレコード	HITAC コード	漢テレコード	HITAC コード
0 128.	ホ 163	M 164	O 214
1 001	マ 164	N 037	P 215
2 002	ミ 165	O 038	Q 216
3 131	ム 166	P 167	R 217
4 004	エ 167	Q 168	S 226
5 133	エ 168	R ! 041	T 227
6 134	ヤ 169	! 042	U 228
7 007	ユ 170	半 171	V 229
8 008	ヨ 172	* 044	W 230
9 137	モ 173	* 176	X 231
ち 138	リ 174	/ 049	Y 232
= 011	ル 175	S 050	Z 233
！ 140	レ 186	T 179	0 240
+ 016	ロ 187	U 052	1 241
A 145	ヲ 188	V 181	2 242
B 146	ン 189	W 182	3 243
C 019	ア 193	X 055	4 244
D 148	ビ 194	Y 056	5 245
E 021	シ 195	Z 185	6 246
F 022	ド 196	← 186	7 247
G 151	エ 197	> 059	8 248
H 152	フ 198	(188	9 249
I 025	ギ 199	(問) 013	064
? 026	ヒ 200	064	064
. 155	イ 201	上記以外のコード	→ 095
) 028	ジ 209		
- 032	キ 210		
J 161	ル 211		
K 162	ム 212		
L 035	ヌ 213		

データの構造 入力データの構造について、KWIC Index File の場合の入力形式と同じであります。上回の処理の対象となるデータをファイルと定義します。漢テレによるデータの打鍵のとき、オフ表に示す記号を pre editing として挿入

す。

実際の原文中には、つゞきのようだけ問題点がある。

1) 強調文 し例: ニュース述べる重力とは、哲學的なものである)。

2) 文の途中の圖、化学式表、数式等がその説明文。

3) 简易書き文。

4) ページ

5) 参考文献の指標

本来、原文が有するあらゆる情報を翻字しておこうが、実際にはその処理目的への必要性と翻字化の煩雑さとの兼ね合いで決定され3かであります。今回、本報告に使用したデータを中心以下に以下の処理を施した。

図方すべきの説明文は、すべて除外し、数式は可能な限り原データに忠実に入れる。だが、数式の終りには、これを文と見立て、文末の記号「。」を挿入した。(「。」がないと、数式と次の文がつながったものと計算機は解釈する。) ページ内添付とけなこ処理した。

参考文献の指標は、原データで「...である。」と「...で」の部分は「...である。」と変更した。文の処理の判定で、「？」をつづく文と解釈されたためであります。

だが、漢テレになり未登録の漢字については、カタカナ表に示すように、その読みをカタカナに直し(シグナ)ではさ表現方法を用いた。

例: 兩棟類 → 兩IIセイII類

翻字の手順は、

1) 原データを入力構造の約束に従って Pre Editing を行なう。

2) Pre Editing されたファイルを漢テレを使、テキスト上に穿孔する。

3) 穿孔データを計算機により磁気データ(MT)に高積する。

3. 分ち書きの手法

日本語には、英語文にみられるような語を単位に分ち書きせざれど記述されたという特徴がある。本システムの基本ルーチンである ALPS-KWIC Index File は、分ち書きされた語を単位とする処理方法が用いられておりため、Pre Editing として、人手による計算機による分ち書き処理を必要とする。

第2表 入力データ構造

<ファイル>	::= <データ> △
<データ>	::= <テキスト> <データ><テキスト>
<テキスト>	::= <テキスト名><注釈><著書名><テキスト内容>
<テキスト名>	::= ◇<文>
<注釈>	::= ◇<文>
<著書名>	::= ◇<文>
<テキスト内容>	::= <章> <テキスト内容><章>
<章>	::= <章名><章内容>
<章名>	::= §<文>
<章内容>	::= <節> <章内容><節>
<節>	::= <節名> <節内容>
<節名>	::= ◇<文>
<節内容>	::= <段> <節内容><段>
<段>	::= ◇<文> <段><文>
<文>	::= <文内容>。
<文内容>	::= <語> <文内容><語>
<語>	::= <文字> <語><文字>
<文字>	::= <外国文字> <カタカナ> <ひらがな> <漢字> <数字> <特殊文字> <漢テレにない字>
<外国文字>	::= A B C … a b c … A B B … α β γ …
<カタカナ>	::= アイウ…
<ひらがな>	::= あいう… つゅよ…
<漢字>	::= ハタヤ … 伴鶴腕
<数字>	::= 1 2 3 … 1 2 3 … りゅき
<特殊文字>	::= + - / △ * : ..
<漢テレにない字>	::= <カタカナ>

但し、上に用いられた記号は次のことを意味する。

< > ; 最小単位の集合

::; 左辺は右辺の構造で成り立つ

| ; 単位の和集合

<> | <> <> ; 繰り返しを許した和集合

… ; 同種類の文字によって左辺の構造を規定する

自動行分書きについては、多くの方法が紹介されています。こゝでは、植村氏が報告されていける单纯な手続手による方法を採用した。すなはち、文字種の異りによる分割による方法である。漢テレのコード配分は、字種に対して部分集合を形成するように設計されている。とくに「文字コード」(すなはち2列)や「オフ2列」のコードを判定するだけで、字種識別が可能であり、単純で高速な処理ができる。その処理手順は、入力された連続した文字列中、隣接した2文字を比較し、その字種の異なった部分にスペースを挿入する。たゞしこれを漢字種を判定し、分割を行なわなければ。

処理結果 日本語における正書法の規則がヨーロッパの諸言語に比較して厳密に守られており、そのため問題が生じる。第2回の例で示すように、用言の語尾、名詞、用言の漢字かひ振り表現、数の表現が分離する例と副詞、接続詞等にみられる分離が行なわれない例が多く見られる。後者については、量的で誤書を用いて認定が可能であるが、前者の場合には、多くの問題を含んでいる。

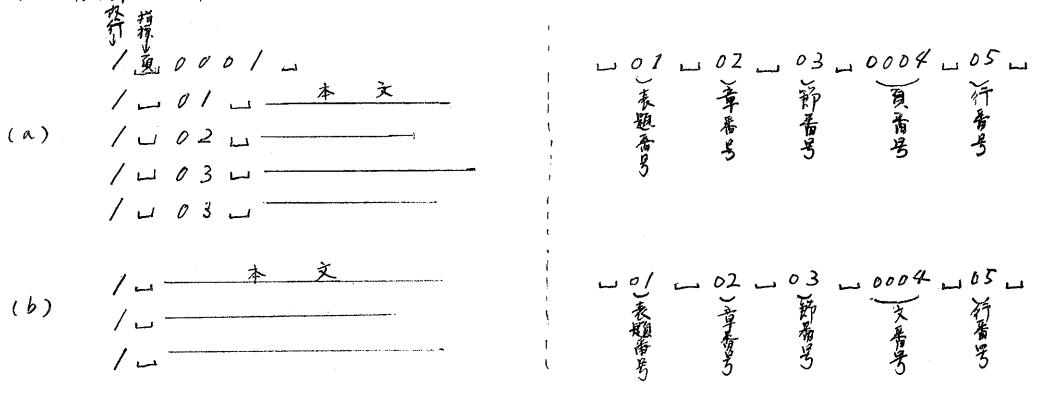
(a)	見立て → 見立て かけ離れた → かけ離れた は虫類 → は虫類 市工図 → 市工図	目覚ましい → 目覚まし 抽出出す → 抽出する 20億年 → 20億年
(b)	一見脇が → 一見脇が 一方抑制 → 一方抑制 諭じたもつとも → 諭じたもつとも	方2回 分離書きの不成功例

4. 記事の作成

作成された「コンコード」システムの利用にあたっては、その Keyword を含む文脈の原データとの参照が必要となる。その記事の書式を2種類とした。

- 1) 第3回(a)で示すように、原データでのページ、行を人手により、漢テレで翻字することで挿入する。
- 2) (b)で示すように、計算機で自動的に、文番号と同文中での行番号を挿入する。

2つの方法が選択できる。



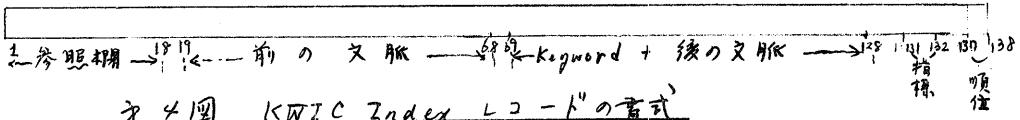
第3回 入力データの書式

文の1行の長さは、出力装置によって調整が可能であり、語間のスペースは、1行に標準化が行なわれている。

作成の手順は、入力データより、指標判別を行ない、参照番号をセットする。1行分の文脈を移した後、MTに蓄積される。なお、KWIC Index の対称から除外した表題、注等については、参照番号欄の行数、文数を0とする。

5. KWIC Index File の作成手順

KWIC Index の作成は、記事のMTを入力として、文を単位にコードに蓄積される。この文から Keyword の設定を行ない、参照番号のセット、Keyword とその Keyword に続く前後の文脈を移し、次回の書式で MT に出力する。これを1コードとする。



方針わち、

(1) 参照番号：記事中の Keyword を含む文脈に与えられてる参照番号がそのまま与えられる。

(2) 前の文脈：記事中の Keyword の前の文脈を規定数だけ出力させる。ただし次の文の先頭が規定字数内にある場合は、それ以前をスペースとする。

(3) 後の文脈：記事中の Keyword を含む後の文脈を規定字数だけ出力する。ただし、規定字数内で、文が終つてる場合には、残余をスペースとする。

(4) 指標：記事中、標題、章名、節名の指標、Keyword が含まれてる行にコメントが付されてる場合には、右端にこれを示す記号とする。

(5) 出現順位：記事中に、その Keyword が出現した順位を示す。

なお、逆配列(語尾から語頭へ文字配列)、アルファベット順(五十音順)は、ソートする。KWIC Index File の作成は、1レコードを中心て反転したコードについて、ソートを行なう。その結果を逆転して出力する。

1ファイル分の KWIC Index レコードが MT に蓄積されると、COBOL のソート、アソート用の、ディスク・ソートを行なう。すなはち、Keyword とそれに関連する文脈について、指定の字数についてソートされる。結果は、記事の後に統合して、MT に KWIC Index File として蓄積される。

6. Keyword の統計処理

KWIC Index File で得られた結果を用ひ、自然言語の分析に必要な以下のいくつかの統計データを計算し出力するプログラムを作成した。

(1) 正規語数：記事中に出現した Keyword の总数。

(2) 異なり語数：記事中に出現した異なる Keyword の总数。

(3) 異なり語の出現頻度数。

(4) 正規語数に対する各異なり語の出現頻度数の割合。

(5) (4) の累積頻度。

処理手順は、出力形式により、Keyword 配列と頻度順配列の 2つがある。

Keyword 配列

(1) Keyword の総数を記憶。

(2) KWIC Index File からレコード"単位で読み込み Keyword の抽出を行なう。

(3) 前の Keyword と比較、同じ場合は出現頻度の計算、累積していき場合は、前の Keyword についての情報を MT に出力し、新しい Keyword をセットする。

(4) (2), (3) を繰り返し、全レコードの処理を終ると、次に図に示す統計情報を出力する。

*****	統計情報	*****
* 文献名	脳とオートマトン。	*
* 見出し語の取り出し方法	逆順	*
* 見出し語の配列方法	出現頻度順	*
* 異なり語数	1116語	*
* 認語数	3425語	*
*****	*****	*****

頻度順配列

次に図 統計情報の出力例

- (1) Keyword 配列の MT を入力とし、出現頻度数を方 1 Key, Keyword を方 2 Key として、ソートを行う。
(2) (1) の結果について、累積頻度を計算し出力する。
(3) (1), (2) を繰り返し、全レコードを処理後、統計情報を出力する。
以上の手順を繰り返して、Keyword の取出し方法(正, 逆), 配列順序(keyword, 頻度) の 4 種のうち、必要なものを選択出力である。

7. 出力形式(たて横変換)

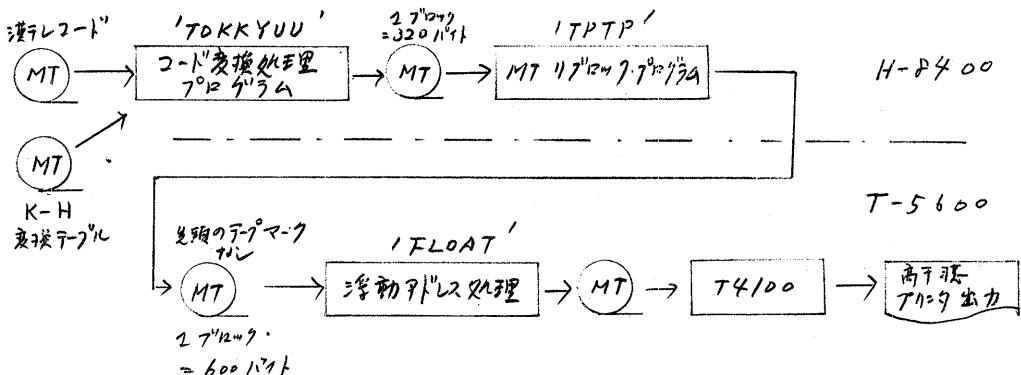
KWIC Index を印刷するには、漢テレの場合も、後述する高画質アリニアで、紙の横幅に物理的な制限があり、MT 上に蓄積されていけるレコードを全部印字することができない。そこで以下に示すいくつかの方法がある。

- (1) KWIC Index レコードを左右に 2 分し、別々に印刷し、後に貼付して横書きにする。
(2) 前文脈の左端から以後文脈の右端を切りとり中心部分 30 字程度を印字し、横書きにする。
(3) 横書きをたて書きに変換して、レコードの全文字を出力する。
(1) の方法は、漢テレの行送り機構の精度が低いため貼付の左右ですれか生じる。また(2) の方法は、文脈が短くなり、KWIC の特徴を損なうといった欠点があり、たての一行の長さが自由に選べる(漢テレ) (3) の方法をニコでは採用した。この方法の欠点は、特殊記号、数字等でたて書き用の記号がない場合、吹き出しがうまくできない完全な変換が行えないのがある。

処理手順は、1 行(40 行)分づつコア内に読み込み、参照番号を漢数字に変換して下段にたて書きとし、前の文脈を上段、Keyword と後の文脈を下段に漢テレコードに変換して移す。また真二に、製本に複利かように奇数は左、偶数は右側の下方に負数がつけられる。これを 1 単位として MT に出力される。

漢テレで印字する場合には、HITAC 4400 のユーティリティ T012 や T012 は紙テープに出力し、漢テレで打鍵される。(打鍵速度 = 130 字/分)

高画質アリニア、ミステム(T4/100)を用いた高速印刷を行なう場合、次に図に示すように、コード変換処理と MT ブロック、フォマットの調整が必要とする。



第6回 T4100 处理手順

3.

コード変換処理プログラム ($\text{J}^{\circ}\text{E}^{\circ}\text{G}^{\circ}\text{U}^{\circ}$ 名 'TOKKYUU') は植村氏が作成された漢字コード向の相互変換プログラムで、これは、漢字からT4100コードへの変換処理を行なつた。

浮動アドレス処理サポート $\text{J}^{\circ}\text{E}^{\circ}\text{G}^{\circ}\text{U}^{\circ}$ 'FLOAT' は、T4100システムの浮動アドレス方式 $\text{J}^{\circ}\text{E}^{\circ}\text{I}^{\circ}\text{N}^{\circ}\text{T}^{\circ}$ の入力 MT を作製する目的で当所の推論機構研究室で開発され、TOSBAC-5600 を使用して一ガ月間実験され、113 $\text{J}^{\circ}\text{E}^{\circ}\text{G}^{\circ}\text{U}^{\circ}$ である。

得られた MT を T4100 32 テムへ割り当てる。記事、KWIC Index File、統計処理結果等が高速印字出力される。(1 文字 = 32×32 ピット; 字種 = 8192)

8. 处理結果と応用

現在、数千語からなる日本語文書 6 種を出力結果が得られており、オランダ語 10 万字からなる特許公報のコード化が作成中である。また、T4100 1:12 の文字 $\text{J}^{\circ}\text{E}^{\circ}\text{G}^{\circ}\text{U}^{\circ}$ (記憶させ) = 1:1 となり、数万語からなる古代 12 三字のコンコードансを作成中である。

すて、KWIC Index File に記憶させてある Keyword の出現順位と、異なり語配列表を利用して、高速付箋検索システム、逆配列表による活用実験、頻度配列表による重要 Keyword の抽出等の応用を考えている。

最後に、本システムの $\text{J}^{\circ}\text{E}^{\circ}\text{G}^{\circ}\text{U}^{\circ}$ 作成を手伝つて下さつた、電気通信大学の岩津直和、佐藤雅之、石川富士夫の諸君に感謝いたします。

参考文献 (1) 坂本義行: Automatic Language Processing System, KWIC Index File, CL研究委員会, 69-5, 1969.3.

(2) 坂本義行、岡本哲也、加藤知巳: ALPS-KWIC Index File の応用 - Concordance & Syntax KWIC の作成、第7回国情報科学研究集会発表論文集, JICST, 1970.

(3) 西村赳彦: 電試漢字仕様案、機械翻訳研究委員会, 1965.9.

(4) 植村俊亮: 電子計算機による自動索引の研究(上), 電子技術総合研究所研究報告 第743号, P60, 1974.2.

付録 I 記事と KWIC Index の出力例

付 錄 五 正配列の統計出力例

脳とオートマトン。

順位	頻度表 (配列順) 見出し語	頻度数	相対頻度	18頁 累積頻度
851	上段	1	0.0003	0.82425
852	成功	1	0.0003	0.82451
853	条件	2	0.0006	0.82514
854	状態	1	0.0003	0.82566
855	状態化	4	0.0012	0.82666
856	走式化	1	0.0003	0.82666
857	情報交換	1	0.0003	0.82666
858	情報処理	1	0.0005	0.82666
859	情報処理能力	5	0.0015	0.82666
860	場合	5	0.0006	0.82666
861	場合所	2	0.0009	0.82666
862	越	3	0.0006	0.82666
863	信号	2	0.0006	0.82666
864	神経回路	2	0.0003	0.82666
865	神経回路群	5	0.0015	0.82666
866	神経回路網	1	0.0003	0.82666
867	神経回路網内部	7	0.0023	0.82666
868	神経系	1	0.0003	0.82666
869	神経細胞	3	0.0009	0.82666
870	神経生理学者	1	0.0006	0.82666
871	神経信号	2	0.0012	0.82666
872	神経網	4	0.0012	0.82666
873	神経	0	0.0003	0.82666
874	神経振動	1	0.0003	0.82666
875	振動	1	0.0003	0.82666
876	進歩	1	0.0003	0.82666
877	進化	1	0.0003	0.82666
878	歩	1	0.0003	0.82666
879	新間	0	0.0003	0.82666
880	人工知能	1	0.0003	0.82666
881	人工知能学	2	0.0006	0.82666
882	水生生物学	1	0.0003	0.82666
883	推測	2	0.0006	0.82666
884	誰	1	0.0003	0.82666
885	…体細胞	5	0.0015	0.82666
886	数学	6	0.0023	0.82666
887	数学学者	8	0.0023	0.82666
888	数学	1	0.0003	0.82666
889	数学者	1	0.0003	0.82666
8890	数学	1	0.0003	0.82666
8891	数学	1	0.0003	0.82666
8892	数学	3	0.0009	0.82666
8893	数学	1	0.0003	0.82666
8894	数学	1	0.0003	0.82666
8895	数学	2	0.0006	0.82666
8896	数学	1	0.0003	0.82666
8897	数学	1	0.0003	0.82666
8898	数学	1	0.0003	0.82666
8899	数学	1	0.0003	0.82666
900	設計	1	0.0001	0.82666

付録 Ⅱ 逆配列の頻度順統計出力例

脳とオートマトン。

順位	頻度表 (頻度順) 見出し語	1 頁		
		頻度数	相対頻度	累積頻度
1	め	274	0.08000	0.08000
2	ニューロン	237	0.06922	0.14922
3	に	184	0.03042	0.20363
4	が	183	0.02422	0.22688
5	を	775	0.02212	0.24888
6	は	644	0.01816	0.26888
7	(648	0.01400	0.28888
8)	410	0.01200	0.30445
9	で	324	0.00870	0.32295
10	シナプス	323	0.00644	0.33953
11	1脳	221	0.00618	0.35617
12	と2脳	221	0.00558	0.36617
13	そのも	220	0.00500	0.37222
14	神経回路網	220	0.00477	0.37816
15	・な考	207	0.00477	0.38607
16	小脳	166	0.00477	0.39537
17	こする	166	0.00447	0.40423
18	オートマト	155	0.00447	0.41231
19	電位	144	0.00447	0.41916
20	として	143	0.00411	0.42261
21	には	122	0.00388	0.43000
22	蓄る	122	0.00355	0.43338
23	ある	111	0.00355	0.43680
24	大脳	111	0.00322	0.44002
25	細胞	111	0.00322	0.44346
26	とい	100	0.00299	0.44966
27	して	100	0.00299	0.45584
28	このよ	100	0.00299	0.46121
29	な進化	100	0.00299	0.46707
30	興奮性	100	0.00299	0.46753
31	多	100	0.00266	0.47800
32	神経網	100	0.00266	0.48006
33	もや	99	0.00266	
34	した	99	0.00266	
35	リババレー	99	0.00266	
36	ーション	99	0.00266	