

国立国語研究所学術情報リポジトリ

On the verbal concordance to analyze the usages
of a word in Japanese literature

メタデータ	言語: jpn 出版者: 公開日: 2017-03-31 キーワード (Ja): キーワード (En): 作成者: 鶴岡, 昭夫, TSURUOKA, Akio メールアドレス: 所属:
URL	https://doi.org/10.15084/00001044

国語研究のための索引作成システム

—新システムの基本思想を中心として—

霧岡 昭夫

1. まえがき

これまで、世の中には多くの総索引が作られている。その大半は、日本古典などの文学作品の語をカードに書き写して並びかえるという、手作業によって作られたものである。しかしこの十年来、電子計算機を利用して索引を作る方法が開発され発達して来ているので、われわれ国立国語研究所では、これを改良して言語研究に便利な索引や検索用言語データを作ること考えたのである。そしていくつかのシステムが作られてきたが、今回、田中章夫、土屋信一、中野洋と本稿筆者の四人の共同研究により新しい用語索引作成システムが完成した。そこで、この索引システムの内容を報告するとともに、ここに至るまでの経過と今後の見通しについて述べることにした。

2. これまでの索引システム

○国立国語研究所での索引作成システム

前節で書いたように、国立国語研究所では、これまでにさまざまな索引システムとそれによる総索引を作成してきた。われわれは、言語研究のためには、文脈付き索引（検索すべき語=key word が文脈中=in context にあるので、KWIC 索引と呼ぶ）が便利であると考え、この方式を持つ索引を多く作ってきた。今回新たに作成したものを含め、現在までの索引作成システムは、入出力の方法という点で、つぎのように分類される。

- ①漢字テレタイプ入力（漢字かな）→漢字テレタイプ出力（漢字かな）（注1）
- ②フレキシ入力（ローマ字）→ラインプリンタ出力（ローマ字）（注2）

- ③フレキシ入力（片かな）→ラインプリンタ出力（片かな）^{（注3）}
- ④漢字テレタイプ入力（漢字かな）→ラインプリンタ出力（片かな）^{（注4）}
- ⑤漢字テレタイプ入力（漢字かな）→高速漢字プリンタ（漢字かな）・ラインプリンタ出力（片かな・ローマ字）

電子計算機を利用しての索引作りは、まず①から始まった。日本語は漢字かな交り文が一般的であり、それを処理することを最初に考えたわけである。そして、このシステムを使って芥川龍之介の「くもの糸」の索引が作られた。しかし、このシステムは、出力機である漢字テレタイプ印字機の印字速度が遅い（120字/分）ので、短編作品の索引は作れても、中編・長編のものとなると非常に時間がかかって実用的でない。

そこで、出力を高速（90000～100000字/分）のラインプリンタ（電子計算機の印字機）による②③④のシステムが作られたのである。②③では入力機にフレキシ（英数字・片かなタイプライタに紙テープさん孔機の付いたもの）を用い、片かなで入力すれば片かな索引に（③）、ローマ字で入力すればローマ字索引（②）になるようになっている。また、④のシステムは、漢字（よみがな付き）とかなで漢字テレタイプによって入力し、ラインプリンタにより片かな索引を作るものである。これは、漢字かなで入力してあるので、漢字テレタイプで印字することもできるが、前に述べた理由からほとんど行われていない。③のシステムでは、主として日本古典の『遊子方言』『浮世風呂』『浮世床』『心中天の網島』『今昔物語集26巻・30巻』『当世書生気質』などの索引が作られた。また、④のシステムでは、昭和41年度から昭和48年度にかけて行われた「現代新聞の語彙調査」の長単位カナKWICが出力されたほか、昭和48年度から行われている「漱石・鷗外の用語研究」のために、『三四郎』『行人』『硝子戸の中』『夢十夜』『高瀬舟』『青年』のKWIC索引も作られた。②のシステムは入出力がローマ字ということから日本語への応用という点で利用度が低いが、現在『夢十夜』ローマ字索引が出来ているほか『星の王子さま』など外国語の索引がいくつか作られている。

②③④のシステムではラインプリンタによる出力方式をとったために出力時間は短くなった。しかし、②はローマ字文しか適用できないし、③④は片かな出力のため読みにくく、同形見出しの識別もやっかいである。そこで、最近開発された高速漢字プリンタ（50000～100000字／分）を用いて漢字かなまじり文を高速出力するシステムの開発にかかったのである。

高速漢字プリンタを出力に利用するだけならば①や④のシステムの出力部分を変えればよい。しかし、どうせ言語データを作るならば、索引作りばかりでなく、さまざまな言語研究に使えるデータを作っておこうということと、さまざまな形のデータを扱える汎用性のある処理システムにしておこうということになって、⑤の新システムの開発が始められた。このシステムにおける、言語データ作成のための前処理（プレディット）については本稿（次節）で、またそのデータを処理するプログラムシステムについては中野洋論文（本集16ページ以下）において書かれている。

このシステムによって作成された（作成されつつある）KWIC索引をはじめとした言語データはつぎのとおりである（語数は概数。*印は第一次出力—mini KWIC—による校正用出力、**印は処理途中の未完成品）。

- 森 鷗外『寒山拾得』（4100語）
- // 『雁』（50000語）
- // 『山椒大夫』*（16000語）
- // 『渋江抽斎』**（150000語）
- 夏目漱石『草枕』*（58000語）
- // 『坊っちゃん』*（55000語）
- // 『こころ』**（150000語）

なお、昭和49年度から開始された「高校教科書の用語調査」のシステムもこのシステムを基礎としたものである（教科書調査のシステムについては別な機会に発表が行われる予定）。

○国立国語研究所外の索引作成システム

電子計算機を用いて索引を作る実験や作業は現在では非常に多くのところで行われている。それらの多くは特徴のある興味深いものであるが、言語研究を目的とした索引はそれほど多くはないように思われる。まとまった作品の索引では、『平家物語総索引』（金田一春彦・清水功・近藤正美編・1973. 学習研究社刊）があるが、これは従来の手作業方式の一部を電子計算機（HITAC 10）にやらせるようにしたものである。そのほか植村俊亮「漢字かなまじり文KWIC」（『情報処理 10-5』1969年）で発表されたシステムおよびそれによる新聞長単位60万語のKWIC索引がある。また、KWIC索引を利用して分析した結果が「放送用語のKWIC索引」（管野謙『文研月報』昭和46年6月）にあげられている。

（注1） 斎藤秀紀作成。「電子計算機と漢テレによる用語総索引の作成」（『電子計算機による国語研究』1968）

（注2） 斎藤秀紀作成。および江川清作成。

（注3） 土屋信一作成。「カナ入力による日本語文総索引の作成」（『電子計算機による国語研究Ⅳ』1972）

（注4） 石綿敏雄作成。「新聞用語調査の用例印字プログラム“COBOL-KWIC”」（『電子計算機による国語研究Ⅲ』1971）

3. 新しい索引作成システム

3.1 言語研究の方向

電子計算機は、入力するデータとそれを処理するシステムとがしっかりしていさえすれば、われわれの必要とする形式・内容をもった言語データ（索引を含めて）を、忠実に作ってくれる。したがって、われわれが言語データを作るには、どのような言語研究をするか、そのためにどのような形式・内容をもった言語データにするか、そしてそれにはどんな情報を付加する必要があるか、ということを検討することが必要になる。

今われわれが考えている研究は、音素・文字・語（自立語・付属語・語構成要素など、さまざまなレベルがある）・文節・句・節・文・文章といったいろ

いゝな言語単位について、用法（出現環境・出現率など）を考察して行く計画である。

文字の出現環境（例えば連続確率など）や出現率などは、一文字が一単位となっている原文データの磁気テープ（MT）ファイルがあればそれをもとにして電子計算機で自動的に調査することができる。音素については漢字かなまじり文においては漢字にふりがなを付けておけば、それをもとにして自動的に全文かなのデータが得られるから、あとは文字の場合と同じように調査をすることができる。語や文節の研究には、文章から語、文節の単位で用例が得られるように単位切りをする必要があるが、これについては次項で述べる。文や文章の研究は、音素・文字・語・文節などをもとにしたやり方、本文から直接データ（文の種類わけなど）を得るやり方といった、さまざまな調査方法を考へている。（注5）

以上の研究は大きく分けると、MTファイルをもとにして電子計算機で検索・カウントなどを処理して行うものと、電子計算機で編集し、漢字プリンタでプリントアウトした KWIC 索引によるものとなる。そして従来、語彙調査は前者を中心に、索引作りは後者を中心に行われていたように感じられる。本システムは、どちらにも使える、というよりも文章を材料にした総合的言語研究ができるように開発されたものである。

（注 5） これまで、このシステムで作られたデータを用いた研究は次のものが発表になっている。

- 鶴岡昭夫「電子計算機による代表構文作成の試み」（『ことばの研究 5』1974）
- 鶴岡昭夫「文節タイプ連続の研究」（『電子計算機による国語研究 VII』1974）
- 米田正人「“文の長さ”の統計学的一考察」（『電子計算機による国語研究 VII』1974）

3・2 単位切り

今までの索引は、ほとんど一種類の単位を用いて作られている。カードに転

写する手作業方式では単位の種類をふやせば、それだけ作業が増大するから、複数の単位を用いることはまずなかった。電子計算機を用いた索引作りも、ほとんど単一の単位で行われてきた。国立国語研究所で作られた索引の単位はこれまでに行われた用語調査で用いられた単位がそのまま用いられることが多かった。それらは次のようなものである。

- α 単位 (国立国語研究所報告4 『婦人雑誌の用語』1953. 19ページ以降参照)
- β 単位 (国立国語研究所報告21 『現代雑誌九十種の用語用字』1962. 6ページ以降参照)
- 長単位 (国立国語研究所報告37 『電子計算機による新聞の語彙調査』1970. 13ページ以降参照)
- 短単位 (国立国語研究所報告37. 15ページ以降参照)

このうち、たとえば前節で示した作品のうち、『高瀬舟』『青年』などは短単位で単位切りがなされているが、それは次のようになっている(単位の切れ目は/で示す)。

／高瀬舟／は／京都／の／高瀬川／を／上下／する／小舟／で／ある／。／
徳川／時代／に／京都／の／罪人／が／遠島／を／申し／渡さ／れる／と／
……(『高瀬舟』より)

／純一／は／立ち留まつ／て／名前／を／読ん／で／見／た／。／自分／の
／捜す／大石／絹太郎／と／いふ／名／は／上／から／二三／人／目／に／
書い／て／ある／ので、／すぐ／に／見附かつ／た／。／…(『青年』より)

しかし、今までの単位は、もともと語数カウントを中心とする語彙調査のために考えられたものであるので、ことばを探す場合には都合の悪いことが少なくない。たとえば、上にあげた『高瀬舟』や『青年』の索引からは、「高瀬舟」の「舟(ブネ、またはフネ)」、「高瀬川」の「川(ガワ、またはカワ)」や「立ち留まつた」の「留まつ(ドマツ、またはトマル)」などの語を検索することが困難である。

また、本研究は前項で述べたとおり、さまざまな言語単位について研究をす

る計画なので、一種類の単位では不足である。

そこで、今回開発した新システムでは、C単位、L単位、S単位の三種類の単位を用いることができるようになっている。

[C単位]

C単位は、文節に相当するもの、および記号（連続した記号は全体で一単位）である。文節の認定は、大すじ次のように定めた（本稿では切れ目を/で表す）。

(1)スペース（一字分以上の空白）で切る。

(2)助詞・助動詞のあとで切る。助詞助動詞は『現代語の助詞・助動詞』（国立国語研究所報告3）によるほか、「じゃ」「にゃ」「ちまう」「ちやう」「てる」「とる」や、形容詞型活用語・形容動詞型活用語の連用形の直後の「ある」「ござる」を加える。「一の～」の形の体言句で、新潮国語辞典で一語扱いにされているもの（「板の間」「菜の花」など）は、「の」の後を切らない。

(3)助詞・助動詞を伴わないで、主語、連用修飾語、連体修飾語、接続語、独立語となっている自立語は一単位。体言の連体修飾はみとめないが、時間、場所、肩書きなどで、下位・上位のものとの並列は切り離す（/5月/1日/正午に/ /東京都/北区/西ヶ丘の/）。

[L単位]

L単位は、C単位（文節および記号連続）を、次のように切り離したものである（本稿では切れ目を|で表す）。

(1)記号は一個につき一L単位。

(2)助詞・助動詞（それらが連続している場合はその先頭のもの）の前で切る。

助詞・助動詞はC単位の(2)で述べたとおりである。

(3)副詞語尾・形容動詞語尾は語幹から切り離さない（/はっきりと/ /すぐに/ /静かだ/ /立派です/ cf. /友達!と/ /人間!だ/ /うさぎ!です/）。

[S単位]

S単位は、L単位を、現代語で意味を担う最小の言語単位（最小単位）^(注6)の一～二回結合によって分割したものである。その切り方のあらまきは次のようになっている（本稿では切れ目を／で表す）。

- (1)助詞・助動詞は一語一単位とする。C単位で切られなかった「菜の花」「板の間」などの「の」も一S単位にする。しかし、「この」「その」「どの」などの「の」は切り離さずに全体で一単位とする。
- (2)和語の自立語は一要素を一単位とする。ただし、名詞・副詞・形容動詞語幹は二要素結合までを一単位とする^(注7)。（I走り／こむI I雨あがりI）
- (3)漢語^(注8)の自立語は、二要素（漢字二字）の結合までを一単位とする。（I社会／主義I I駐在／員I）
- (4)外来語の自立語は、一要素（原語で一単語）を一単位とする。ただし、「ネクタイ」「クーデター」など、日本語で一単語と考えられ、分割できないものは切らない。（Iウォーミング／アップI Iアンダー／ショットI）
- (5)和語要素と漢語要素の混種語である自立語は、(3)と同様に二要素結合を一単位とする（I重箱I I場所I）。それ以外の混種語は種によって切り離し、それぞれの中を(2)～(4)にしたがって処理する。ただし、「デモる」「タクる」などの活用語尾は切り離さない。（Iナイロン／ザイルI I水割り／ウィスキーI）
- (6)動詞型活用の接辞（「…がる」「…めく」など）や、形容詞型活用の接辞（「…がましい」「…っぽい」など）は一要素を一単位とする。また、形容詞・形容動詞の語幹に付く「さ」「み」「げ」も一単位とする。名詞性接辞は一要素として、(2)～(5)の中で処理されるので、外来語、および用言の中では一単位（Iプレ／オリンピックI Iお／美しいI）となるが、体言の中では一次結合は切らない（Iお足I Iお体I cf. Iお／父さんI）。
- (7)サ変動詞「する・ずる」は一S単位とする（Iびっくり／するI I心配／するI）。同じ位置に来る「できる」「いたす」なども同様に処置する。ただし、一字漢字に付いた「す・する・ずる・じる」などは切らない（I愛する

Ⅰ Ⅰ信じるⅠ)。

- (8)形容動詞語尾は語幹から切り離す。(Ⅰ静か／なⅠ Ⅰ立派／なⅠ)
- (9)副詞語尾「に」「と」は前と切り離さない(Ⅰ割にⅠ ⅠじっとⅠ)。ただし、「に」「と」を切り離して単独でも副詞として用いられるもの、三音節以上の擬音語・擬態語などに付く「に」「と」は切り離す(Ⅰ割合／にⅠ Ⅰすぐ／にⅠ Ⅰぐらり／とⅠ)。
- (10)同音・類音の反復形である副詞・擬音語・擬態語は、反復部分が三音節以上の場合切り離す。(Ⅰのっし／のっしⅠ Ⅰのらり／くらりⅠ)
- (11)独立して用いられないものを含むもの、切るべき位置のわからないもの(略称も含む)などは、二要素以上でも一S単位とする。(ⅠけだものⅠ Ⅰ都区内Ⅰ Ⅰ有頂天Ⅰ ⅠPTAⅠ Ⅰべ平連Ⅰ)
- (12)固有名詞(人名・地名)は一要素を一単位とする。(Ⅰ森／鷗外Ⅰ Ⅰ東京／都Ⅰ)
- (13)数(算用数字・漢数字・ローマ字のほか「幾度」「何人」の「幾」「何」などを含むは)一字一単位とする(Ⅰ一／回Ⅰ Ⅰ百／人Ⅰ)。ただし、一、二、……十……というように数え進むことのできないもの、位取りを表す十、百、千、万、億……などは前と切り離さない(Ⅰ一旦Ⅰ ⅠふたりⅠ Ⅰ五十万／円Ⅰ)。

単位切りの作業は、本文のC単位の切れ目に黒鉛筆で/ (スラッシュ)を入れて行き、それが済んだもののL単位の切れ目に青鉛筆で/を入れ、最後にS単位の切れ目に赤鉛筆で/を入れるという、三段階方式をとって行われる。なお、複合してできた接続詞(例えば「それに反して」「そのほか」など)や、連体詞(例えば「こういった」「そうした」など)のようなものはその認定が人によって、また時によってゆれるおそれがあるので、今までに述べたCLS単位で切れるものは切っておき、あとから連語コードを付けていくようにした。例えば、「それに反して」の場合は、/それⅠに/反しⅠて/と切っておいて、それぞれに接続詞の連語コードを付けておくのである(連語コードにつ

いては次項でのべる)。連語が多く集まれば、それを全部電子計算機に登録して自動的に連語コードを付けることも可能になる。

以上のような規則で処理したものは次のようになる。

／唐Ⅰの／貞観Ⅰの／頃Ⅰだ／と／言ふⅠから／、／西洋Ⅰは／七／世紀Ⅰの／初／日本Ⅰは／年号Ⅰと／云ふ／ものⅠの／やつと／出来／掛かつⅠた／時Ⅰで／ある／。／……(『寒山拾得』より)

／親譲りⅠの／無／鉄砲／で／小供Ⅰの／時Ⅰから／損Ⅰばかり／しⅠて／居る／。／小学／校Ⅰに／居る／時分／学校Ⅰの／二／階Ⅰから／飛び／降りⅠて／一／週間Ⅰ程／腰Ⅰを／抜かしⅠた／事Ⅰが／ある／。／……(『坊っちゃん』より)

上の単位切りで、／はC単位の切れ目であると同時にL単位、S単位の切れ目でもあり、ⅠはL単位、S単位の切れ目でもある。したがって／をCLS, ⅠをLS, ／をSと表わすと『寒山拾得』と『坊っちゃん』は、

〔寒山拾得〕				〔坊っちゃん〕		
C	L	S		C	L	S
CLS唐	Ⅰ	Ⅰ		CLS親譲り	Ⅰ	Ⅰ
LSの	Ⅰ	Ⅰ		LSの	Ⅰ	Ⅰ
CLS貞観	Ⅰ	Ⅰ		CLS無	Ⅰ	Ⅰ
LSの	Ⅰ	Ⅰ		S鉄砲	Ⅰ	Ⅰ
CLS頃	Ⅰ	Ⅰ		Sで	Ⅰ	Ⅰ
LSだ	Ⅰ	Ⅰ		CLS小供	Ⅰ	Ⅰ
Sと	Ⅰ	Ⅰ		LSの	Ⅰ	Ⅰ
CLS言ふ	Ⅰ	Ⅰ		CLS時	Ⅰ	Ⅰ
LSから	Ⅰ	Ⅰ		LSから	Ⅰ	Ⅰ

ということになる。Cから次のCの前までがC単位、Lから次のLの前までがL単位ということになる。

(注 6) 国立国語研究所報告21参照

(注 7) 体言類を二要素の結合を一単位としたのは、「屋根」「場合」「割合」「仕

事」など、二要素の結合全体で一語のように使われるものが少なくないという理由からである。しかし、そのために「食べすぎる」は「食べ／すぎる」と切れるのに「食べすぎ」は切れないので「食べすぎ」から「すぎ」が検索できないという問題がある。

(注 8) 本稿にいう漢語、外来語には、それぞれ和製漢語、和製英語などを含める。

3.3 付加情報

このシステムは、多種の情報を処理する能力が備えられている。それらは、つぎのようになっている。

C	親譲り	〔おやゆずり〕	(S1)		L	の	(WR)
⋮	⋮	⋮	⋮		⋮	⋮	⋮
単位 情報	出現 語形 (S 単位)	読み が な	語 種 ・ 品 詞		単位 情報	出現 語形 (S 単位)	語 種 ・ 品 詞

単位情報は、C、L、S単位の先頭、すなわち黒い／の直後の語にはCと書き、L、S単位の先頭、すなわち青い／の直後の語にはLと書き、S単位の先頭すなわち赤い／の直後の語にはSと書いておく。この情報をもとに、C単位のデータ、L単位のデータ、S単位のデータが自動的に出来る。

出現語形のあとに、読み仮名が必要な場合は〔 〕の中にひらがなで書き入れる^(注9)。そしてすべての出現語形のあと()内に語種、品詞の情報を書き入れ、またその品詞が動詞や動詞型活用の接辞(「がる」「めく」など)である場合は活用型、活用行の情報も付け、複数のS単位が複合して他の品詞となっているものには、それぞれに連語コードも付ける(連語コードの付け方は14ページ参照)。付加情報は、表1に示したようなものがある^(注10)。

単位切りの済んだ原文に書き込む空白はほとんどないから、単位切り作業と情報付加作業の間に原稿用紙に清書する作業がある。清書は出現語形一S単位ごとに一行ずつ書き、13ページの例のような形で情報を付ける。先に単位切りの例としてあげた『寒山拾得』と『坊っちゃん』の漢字テレタイプさん孔

表 1 付加情報コード表

1 ケタ目	2 ケ タ 目	3 ケタ目	4 ケタ目	3 あるいは 5 ケタ目
語種コード	品 詞 コード	活用コード	(活用形) (活用行)	連語コード
S 和語	1 純名詞	F 四段・五段 活用	ワわあ行	ササ変動詞の一部
T 漢語	2 (空番)	G 上一段活用	ああ行	ケ形容動詞の一部
U 外来語	3 サ変動詞語幹	H 上二段活用	かか行	コ固有名詞の一部
V 混種語	4 形容動詞語幹 (派生形を含む)	I 下一段活用	がが行	セ接続詞の一部
W 語種不要	5 形容詞語幹 (派生形を含む)	J 下二段活用	ささ行	カ感動詞の一部
X 数字	6 名詞性接辞, 助数詞	K 変格活用	ざざ行	フ副詞の一部
Y 記号	7 数詞	V 融合形, 変 則活用	たた行	レ連体詞の一部
Z 語種不明	8 固有名詞	Z 活用不明	だだ行	ド動詞の一部
%情報無視	9 代名詞		なな行	ヨ形容詞の一部
	A 接続詞		はは行	
	B 感動詞		ばば行	
	C 副詞		まま行	
	D 連体詞		やや行	
	E 動詞		らら行	
	+ 動詞性接辞		わわ行	
	- 形容詞性接辞		ん	
	M 口語形容詞		3 ケタ	
	N 文語形容詞		目 V Z	
	P 助動詞, 形容動詞語尾			
	Q (空番)			
	R 助詞			
	X 算用数字, ローマ数字			
	Y 記号			
	Z 品詞不明			
	%情報無視			

注 ○すべての語に、語種コード・品詞コードがつく。活用コードがつくのは動詞と動詞性接辞だけである。連語コードは必要な時に、人間又は機械によってつけられる。

○語種コードのW (語種不要) は、品詞が、助詞・助動詞・固有名詞である場合に付ける。

○品詞コードの4, 5に語幹とあるが、もともとは活用語の語幹であるが転成名詞の一部となっているものも含む。

用原稿はつぎのようになっている（実際にはこの他にページ情報，行情報，段落先頭情報，本文外の語の情報といったシフトコードが書かれるが，これについては中野論文 25 ページ参照）。

『寒山拾得』

C唐〔とう〕（T1）
Lの（WR）
C貞観〔じょうがん〕（T1）
Lの（WR）
C頃〔ころ〕（S1）
Lだ（WP）
Sと（WR）
C言ふ〔いふ〕（SEFは）
Lから（WR）
C，（YY）
C西洋〔せいよう〕（T1）
Lは（WR）
C七〔なな〕（X7）
S世紀〔せいき〕（T1）
Lの（WR）
C初〔はじめ〕（S1）
C日本〔にっぽん〕（W8）
Lは（WR）
C年号〔ねんごう〕（T1）
Lと（WR）
C云ふ〔いふ〕（SEFは）
Cもの（S1）
Lの（WR）

『坊っちゃん』

C親譲り〔おやゆずり〕（S1）
Lの（WR）
C無〔む〕（T6）
S鉄砲〔てっぽう〕（T1）
Sで（WR）
C小供〔こども〕（S1）
Lの（WR）
C時〔とき〕（S1）
Lから（WR）
C損〔そん〕（T1）
Lばかり（WR）
Cし（SEKさ）
Lて（WR）
C居る〔いる〕（SEIあ）
C。（YY）
C小学〔しょうがっ〕（T1）
S校〔こう〕（T6）
Cに（WR）
C居る〔いる〕（SEIあ）
C時分〔じぶん〕（T1）
C学校〔がっこう〕（T1）
Lの（WR）
C二〔に〕（X7）

C やっと (SC)	S階 [かい] (T6)
C 出来 [でき] (SEIか)	Lから (WR)
S 掛かっ [かかっ] (SEFら)	C 飛び [とび] (SEFば)
L た (WP)	S 降り [おり] (SEIら)

S単位がいくつか集まって別な品詞となっている連語の処理は、たとえば、「とはいうものの」という接続詞や「非常識だ」という形容動詞の場合、つぎの左側の入力原稿のように情報を付けておけば、右側のように処理されるようになっている。

《入力原稿》	《S単位》	《L単位》	《C単位》	《連語》
Cと (WRセ)	と	と	とは	とはいうものの
Lは (WRセ)	は	は		
Cいう (SEFアセ)	いう	いう	いうものの	ものもの
Lものの (WRセ)	ものもの	ものもの		
C非 [ひ] (T6ケ)	非	非常識だ	非常識だ	非常識だ
S常識 (じょうしき) (T1ケ)	常識			
Sだ (WPケ)	だ			

以上のようにプレエディットした原稿を漢字テレタイプで紙テープにさん孔して電子計算機処理にまわすのである。電子計算機の処理システムについては中野論文で述べられている。

(注9) 作業では読み仮名を現代仮名づかいで入れた。それは、ルビの付いていない本文の場合、作業者への負担がより軽いと考えたからである。

(注10) これらの情報のうち、語種・品詞・活用コードは、『新聞の語彙調査Ⅱ』(国立国語研究所報告38, 11ページ)に発表されたものを基本にした。

4. あとがき

今まで述べてきたことにより、新しく開発されたシステムと、それによって作られた言語データが、他のシステム、データとどう違うかということが明らか

かになった。

これだけのシステムができたのは、電子計算機・漢字プリンタなど機械類の発達と、電子計算機による言語処理の研究の成果があったからである。

今後望まれるのは、これらのデータを蓄積し、それをもとに単位切りや情報付加の作業を自動的に行う方式（「一貫処理システム」と名付けている）の実用化によって、人手作業を軽減することと、入力のスピード化・簡単化——たとえば光学文字読取装置の実用化など——といったことである。これらも現代の技術開発の状況を見るとそう遠い将来ではないようである。

【東山旅行】

みー	465	06	C	SD	----	0007120	うーなー、ー	あー、ー	あー、ー	あー、ー
ほーじー	462	01	C	WR	----	0001690	一、二、三、の、	四、五、六、	七、八、九、	十、
生動ー	462	06	C	SC	----	0002490	しーなー時、	、	生動、	こらへら
深うーをー	466	12	C	WR	----	0012640	出、な、の、に、	深、う、を、	時、	、
深うーをー	462	11	C	WR	----	0003150	、	、	、	、
真ー白いー	468	08	C	SM	----	0012090	一、二、三、の、	真、白、い、	日、が、	
肩ーつくー寝れる	463	04	C	WP	----	0003940	て、あ、り、な、	、	、	、
明末ーじー	470	08	C	WR	----	0015240	、	、	、	、
明末ーじー	470	10	C	WR	----	0015610	手、の、あ、な、	、	、	、
明末ーをー	471	02	C	WR	----	0016190	な、に、し、け、に、	、	、	、
濃念めー	467	12	C	SEIC	---	0010890	な、人、だ、と、	、	、	、
上げろーなりー	471	13	C	WP	----	0017340	、	、	、	、
上げろーなりー	471	13	C	WP	----	0017370	、	、	、	、
免罪ーをー	470	12	C	WR	----	0015830	上、の、一、定、の、	、	、	、
定ーじーばー	473	09	C	WR	----	0019970	一、出、し、て、	、	、	、
定ーじーばー	473	10	C	WR	----	0020000	一、被、つ、て、	、	、	、
はうーてー	462	02	C	WR	----	0001910	さ、い、こ、と、を、	、	、	、
与へるーのーてある	468	10	C	WP	----	0012460	に、一、満、足、を、	、	、	、
顔ーにー	465	03	C	WR	----	0006940	つ、と、一、顔、の、	、	、	、
顔ーをー	473	09	C	WR	----	0019840	一、一、排、け、な、	、	、	、
顔ーをー	473	12	C	WR	----	0020180	け、な、一、時、に、	、	、	、
当りてー	473	07	C	WR	----	0019760	こ、で、一、火、に、	、	、	、
当りてー	471	05	C	WR	----	0016930	の、と、一、火、に、	、	、	、
あからーのー	471	04	C	WR	----	0016460	は、す、	、	、	、
鈴ーからー	471	02	C	WR	----	0016230	な、	、	、	、
鈴ーはー	470	07	C	WR	----	0015170	一、其、一、然、説、の、	、	、	、
鈴ーをー	474	06	C	WR	----	0021100	、	、	、	、
あななーはー	463	09	C	WR	----	0004390	つ、な、	、	、	、
あななーはー	466	02	C	WR	----	0008120	よ、く、一、い、い、	、	、	、
深うーのーどー	463	03	C	WR	----	0003020	と、一、一、聞、い、て、	、	、	、
深うーたいどー	463	07	C	WR	----	0004300	、	、	、	、
深うーにー	466	09	C	WR	----	0008720	一、台、用、に、	、	、	、
胸んるーどー	473	09	C	WP	----	0019060	水、の、一、突、で、	、	、	、
胸裡裏てりー	470	07	C	WR	----	0015130	水、一、一、流、さ、な、	、	、	、
寝めー	468	06	C	SEIC	---	0011690	一、一、一、一、一、一、	、	、	、
洗ふーせーてー	472	04	C	WR	----	0017750	の、一、一、一、一、	、	、	、
洗ふーまーりー	472	08	C	WR	----	0018140	一、一、一、一、	、	、	、
洗るー	472	01	C	SD	----	0017460	一、一、の、一、一、	、	、	、
洗るー	470	01	C	SD	----	0014560	一、一、一、一、	、	、	、
洗るー	470	04	C	SD	----	0014740	こ、い、い、し、	、	、	、
ありー	461	05	C	SEFE	---	0000740	ト、一、一、一、一、	、	、	、
ありー	462	05	C	SEFE	---	0002350	一、一、一、一、	、	、	、

