# 日本語情報検索システムにおけるキーワード自動抽出

## Automatic Indexing Methods for Information Retrieval System of Japanese Text

日本語文献の検索システムで、キーワード付与作業の省力化と品質の均一化を目的とした自動インデクシング方式として、日本語文構造解析方式と不要語除去方式を設定し、プロトタイプによる実験を行なった。前者は、キーワードの用語を限定した検索方式での利用を前提としており、各用語の文法・意味情報を用いた日本語文の構造解析により、キーワードとその文中での位置づけや意味をも合わせて抽出する方式で、抽出精度は80~85%であった。後者は、キーワードの用語を限定しない検索方式での利用を前提にしており、汎用性、高速性、保守性及びキーワードの抽出漏れの防止を主目的に考えられた方式で、抽出精度は約95%であった。

絹川博之\* Hiroshi Kinukawa 田中和明\* Kazuaki Tanaka 池上信男\*\* Nobuo Ikegami

#### 11 緒 言

日本語ワードプロセッサを含む漢字入出力機器の普及と日本語情報処理技術の進歩により、漢字や平仮名を含む日本語情報、文書情報を対象にした情報検索システム建設のニーズが高まっている。実際、日本科学技術情報センターでは、JOIS-II(JICST Online Information System II:科学技術情報のオンライン情報検索システム)<sup>1)</sup>を、日本特許情報センターでは、PATOLIS(Patent Online Information System:特許情報検索システム)<sup>2)</sup>を、日本経済新聞社では、NEEDS-IR(Nikkei Economic Electronic Data-Bank Service:新聞記事検索システム)<sup>3)</sup>をそれぞれ開発し、サービス業務を開始している。また、企業内では、これらの社外サービスとの連動及び社内OA(Office Automation)の一環として、特許情報、設計情報及び文書情報の検索システムの構築が各所で進められている。

このような状況を踏まえて、日立製作所では、汎用日本語 情報検索ソフトウェアとして、HITAC Mシリーズ計算機下 で動くORION(Online Retriever of Information)を開発し製 品提供を開始しているが、機能の充実が望まれている。例え ば、単に漢字をサポートしているだけでは十分でなく、使い やすく効率よく運用できる機能の具備が求められている。具 体的には、情報検索では大量情報の蓄積が前提であることか ら、情報入力蓄積作業の省力化が必要である。この作業のう ちインデックス付与作業は,入力情報を解釈加工する作業で あり、単なる情報入力作業とは質の違った作業である。この ため、情報入力作業の省力化という点から、このインデック ス付与作業の自動化がまず考えられ、このためのツール(自動 インデクシング)が求められている。また、情報検索のインデ ックス(例えば、キーワード)は、大量情報の中から所望情報 を短時間で探し出すための媒体であり、このインデックス付 与精度が検索精度を大きく左右している。しかし、従来では インデックス付与作業を人手に委ねていたため, 付与作業者 の個人差があってインデックス付与の品質の確保が難しいも のとなっていた。このため、インデックス付与の個人差の解 消と品質の確保という点からも, インデックス付与自動化ツ ールが求められている。更に、OAの進展に伴いオフィス各 所で発生する電子化された文書の検索を可能とすることが要 請されている。このためには、検索用のインデックス(例え

ば、キーワード)の付与作業の自動化ツールを構成機能として、文書検索システムに組み込むことが必要となってきている。日立製作所では、これらの要請に応ずるため、日本語文構造解析方式と不要語除去方式の自動インデクシングの研究開発を進めている。

以下, プロトタイプによる実験を踏まえて, 方式の特徴, 適用対象などについて述べる。

#### 2 日本語情報検索システムにおける位置づけ

自動インデクシングシステムは、キーワード自動抽出処理と、その校正処理を行なう各サブシステムから構成され、その方式設定に当たっては、検索系と蓄積系から成る検索システム全体にわたった連係を十分配慮する必要がある。

図1に、日本語情報処理システムの全体像を示す。 キーワード検索方式としては、次の2タイプ<sup>4)</sup>がある。

#### (1) 標準ターム方式

検索対象は限定されるが、キーワードとなる用語を標準化して、表記の違いや同義語による検索再現率\*1)(検索漏れ率の補数)の低下を防ぐ方式で、用語の文法・意味情報をもとに、キーワードの文中での位置づけや意味をもインデックスとして合わせて検索することで、検索適合率\*2)(検索ノイズの補数)の向上を図ることも可能である。

#### (2) フリーターム方式

キーワードとなる用語を制限しない方式で、処理方式が適 用対象に依存せず、新語の取込みも容易である。ホストコン ピュータの演算処理やファイルアクセス処理の高速化に伴い、 検索質問中の検索キーと同義語であるキーワードをもつすべ ての文献を検索できる機能の実用化により、最近見直されて きている。

これらの検索方式に対応する自動インデクシング方式として,次の2方式を設定した。

※1) 検索再現率:
(検索文献中の適切な文献数)
(文献中で、本来検索されるべき文献数)

※2) 検索適合率: (検索文献中の適切な文献数) (検索された全文献数)

<sup>\*</sup> 日立製作所システム開発研究所 \*\* 日立製作所ソフトウェア工場

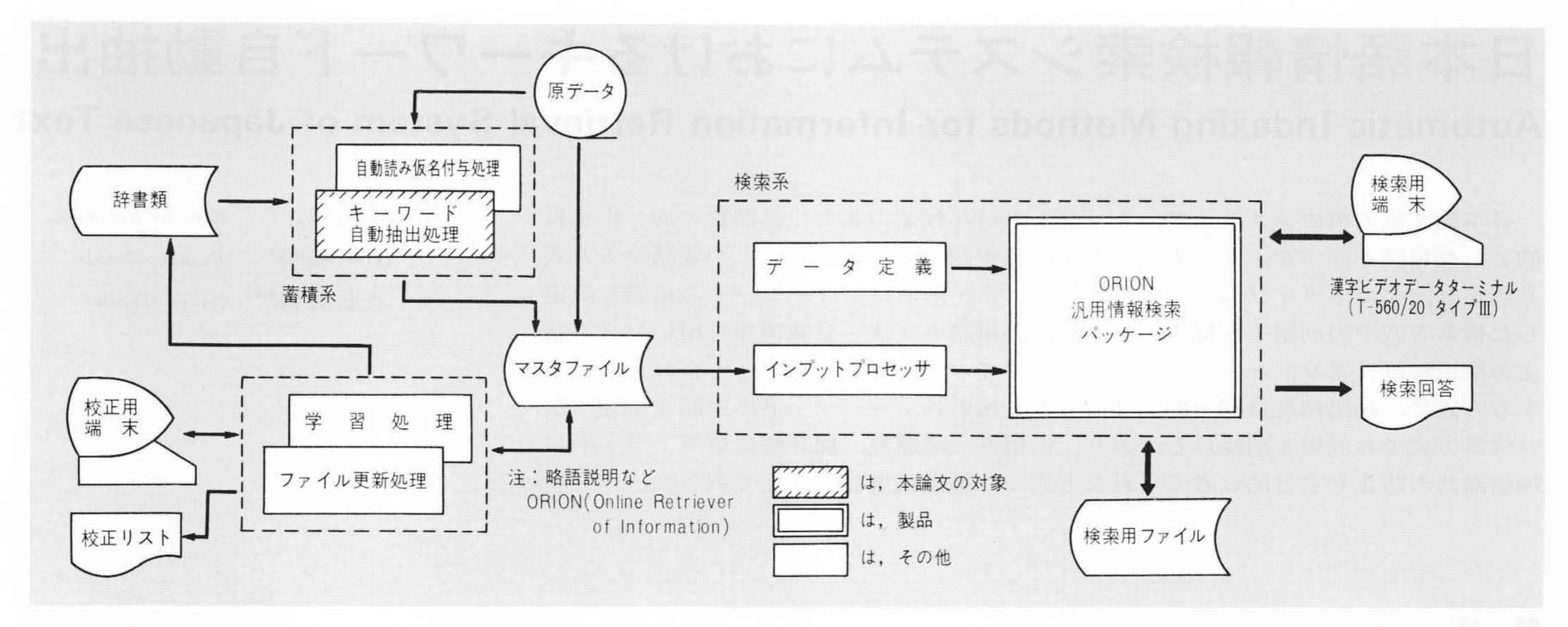


図 | 日本語情報検索システムの全体像 日本語情報検索システムは、ORIONを核とする検索系と、自動インデクシング処理やファイル更新処理を含む 蓄積系から構成される。

#### (1) 日本語文構造解析方式

あらかじめ登録されたキーワードとなる言葉と、その文法・ 意味情報をもとに、日本語文の構造解析を行なってキーワー ドを選定する方式である。

#### (2) 不要語除去方式

汎用性, 高速性, 保守性, 検索漏れの防止を主目的に考え, 不要語が多少含まれていても, キーワードとなる言葉を漏れ なく抽出する方式である。

次章では,この両方式について述べる。

### 自動インデクシング方式

#### 3.1 日本語文構造解析方式による自動インデクシング5)

漢字仮名まじり日本語文からのキーワード抽出と,各キーワードのロール(文中での意味的役割:①主体,②客体,③時,④場所,⑤活動,⑥その他の主題,の6種)付与を自動化し,それによって,省力化と検索適合率の向上をねらうものである。以下,その処理方式について説明する。

#### (1) 文節構成語の認定

文字種の変化点で「文節」分割し、その「文節」と自立語辞書、 付属語表を照合させ、名詞、動詞、付属語の認定を行なう。

- (2) 日本語文構造解析によるロール付与
  - (a) 日本語文型表を参照して、文構造を認定する。
    - (i) 複文パターン表参照による複文構造の認定
    - (ii) 受身態から能動態への変換及び格の標準化
    - (iii) 日本語文型表参照による格支配関係の認定
    - (iv) 名詞同士, 連体形用言の修飾関係の認定
  - (b) 次の条件により、ロールを付与する。
    - (i) 当該文節の支配する動詞
    - (ii) 当該文節を構成する名詞の意味
    - (iii) 当該文節に付く付属語
  - 〔例〕「支配する」

(3) 自動処理結果の確認と修正のために、漢字ビデオ端末を用いて対話形で校正を行なう。

当方式の処理手順を図2に示す。

#### 3.2 不要語除去方式による自動インデクシング方式

英文情報を対象とした情報検索システムでは,英文中から不要語を除去し,残りの単語をキーワードとする方式が実用化されているが,日本語文では,単語単位で分かち書きされておらず,また,明確な分割規則がないことから,独自の方式を検討する必要があった。以下,処理方式について説明する。(1) 漢字仮名まじり日本語文の自動分かち書き

日本語文で重要な言葉は,非平仮名で書き表わされている という日本語の特性を利用して,文字列の分割を行なう。

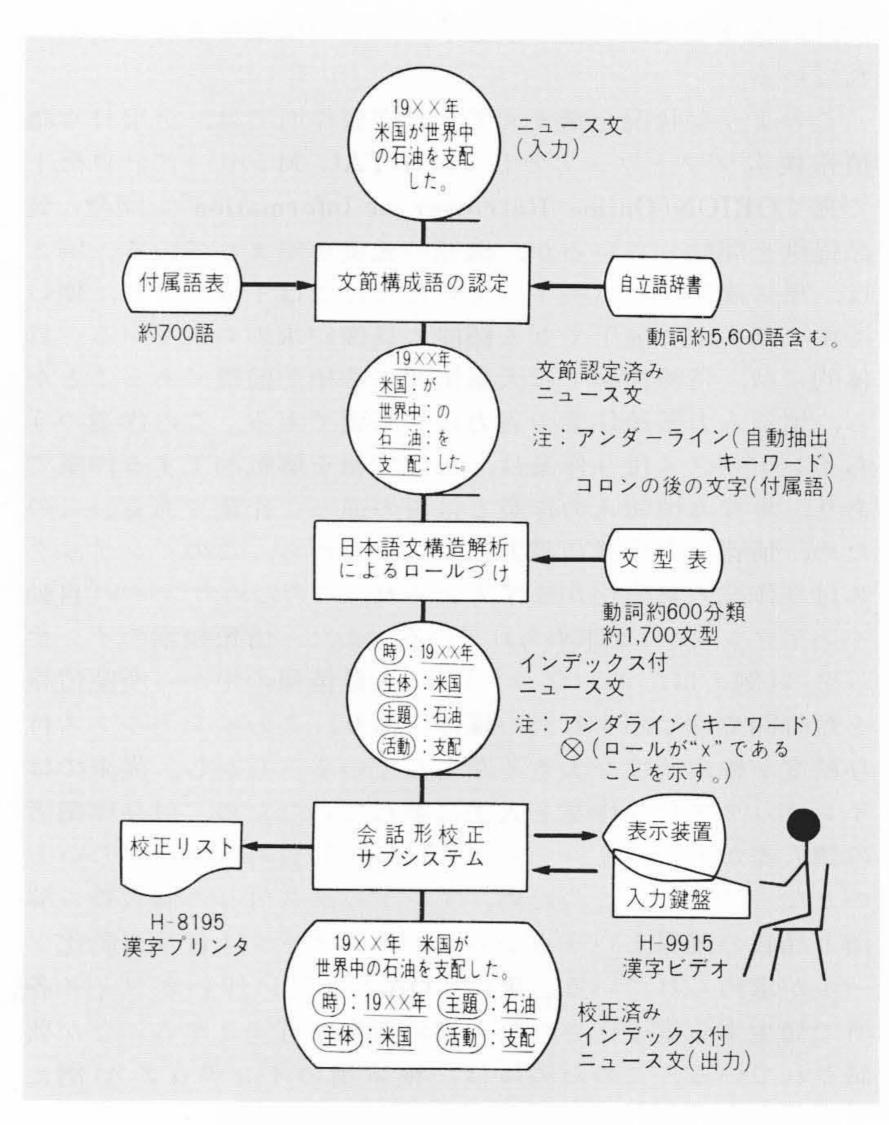


図 2 日本語文構造解析方式による自動インデクシング この方式の処理手順は、次の三つに大別される。(I) 文節構成語の認定、(2) 日本語文構造解析によるロールづけ、(3) 会話形式による自動抽出処理結果の校正

#### (2) 自立語の認定

付属語表を用いて,分かち書きされた文節から付属語を除 去することにより自立語を認定する。

#### (3) 不要語の除去

不要語テーブルの登録語や,キーワードの文字列構成基準 を満足しない自立語を除去する。

処理手順を図3に示す。

ここでは、キーワードを、**図4**の斜線で囲まれた領域内の 文字列と定義している。

#### 4 プロトタイプによる実験・評価

前3章で述べた両方式の実現性を検討するため、表1に示す文献を対象としたプロトタイプを作成し、実験・評価を行なった。

各プログラムはアセンブリ言語で記述されており、その構

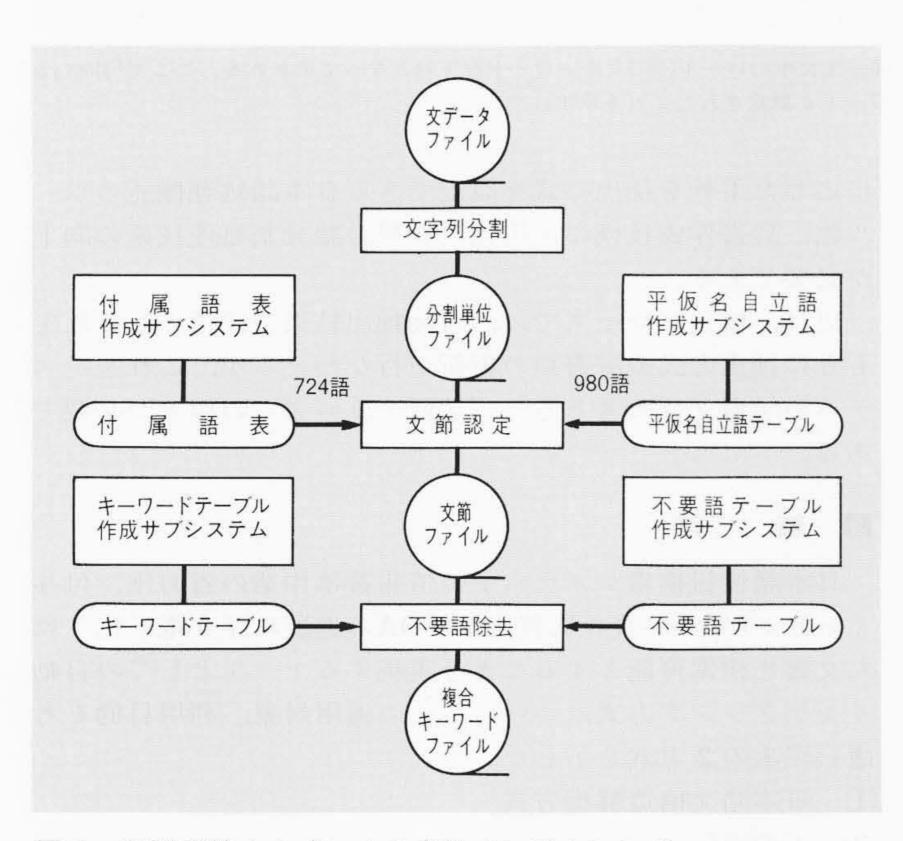


図3 不要語除去方式による自動インデクシング この方式の処理 手順は、次の三つに大別される。(1)漢字仮名まじり日本語文の自動分かち書き (2) 文節の中から自立語を認定する。(3) 不要語の除去

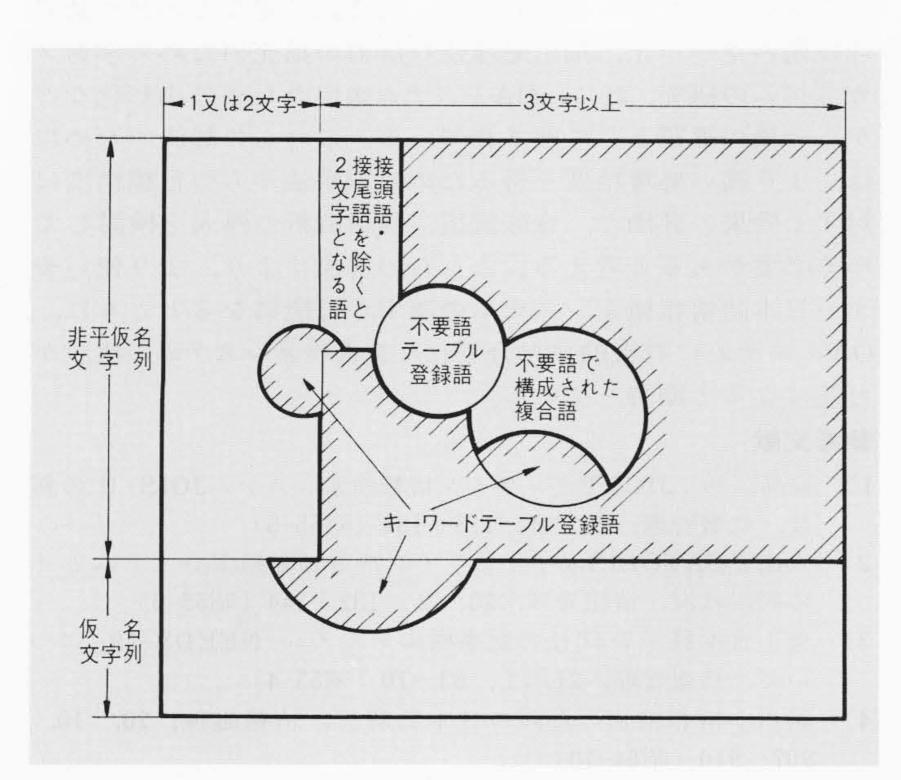


図4 不要語除去方式におけるキーワードの定義例 斜線で囲まれた部分をキーワードとする。

表 | 自動インデクシングの適用データ | 自動インデクシング方式の 適用データとして、日本語文構造解析方式については外電記事文を、不要語除 去方式については「日立評論」のタイトルと要旨を用いた。

自動インデクシング方式	実験対象データ	報告書数 センテンス数
日本語文構文解析方式	外 電 記 事 文	281 1,225
不要語除去方式	「日立評論」のタイトルと要旨	78 361

表 2 自動インデクシングのプロトタイプの構成と性能 各処理 プログラムは、アセンブリ言語で記述されている。処理時間はHITAC M-180 換算で、日本語文構造解析方式が合計で435ms/センテンス、不要語除去方式が 合計で65.1ms/センテンスである。

方式	処理プログラム名	ステップ数	所要メモリ	性 能	備考	
造日解本	文節構成語の認定 処理	3kS	60kバイト	2,640ms/ センテンス	使用計算機 HITAC 8350	
析語 方文 式構	文構造解析によるロ ール自動付与処理	12kS	120kバイト	650ms/ センテンス		
要語除去方4	文字列分割処理	I.9kS	13.3kバイト	8.8ms/ センテンス		
	文節認定処理	I.4kS	40.3kバイト	14.4ms/ センテンス	使用計算機 HITAC M-180	
	不要語除去処理	I.IkS	7.3kバイト	41.9ms/ センテンス	STATE E S	

注:略語説明 kS(キロステップ)

成、ステップ数、所要容量及び1センテンス当たりの平均処理時間は、**表2**に示すとおりである。

500件/日(10万件/年)のデータ更新を考えると、自動インデクシング処理時間はHITAC M-180換算で、日本語文構造解析方式ではCPU(中央処理装置)処理時間で約15.8分、実行時間で約1時間、不要語除去方式では、それぞれ約2.5分、約10分程度であり、大量のデータ更新には後者が適している(なお、各処理時間は、原理実験を目的としたプロトタイプによるものである)。

図5に、抽出結果の編集出力様式を示す。

表3,4に、抽出処理精度を示す。

日本語文構造解析方式における精度は、辞書類(各用語に関する品詞情報や意味分類情報などが登録されている。)の出現語いカバー率\*\*3, 文型表の出現構文カバー率に大きく依存しており、事前の語い調査によって収録された用語辞書や文型表(図2参照)をもとにした外電記事文データ、1,225センテンスの実験では、両カバー率90%、文節構成語決定精度\*\*<sup>1</sup>85~90%、ロール付与精度80~85%程度であった。

本方式は抽出されるキーワードに、主体、客体といった文中での位置づけや、場所、時などの意味情報が付加されるので的確な検索に役立つが、一般には、日本語文の表現形態の多様さから、辞書や文型表の追加・整備に追われることが予想され、分野限定形の利用方法となる。

不要語除去方式の「日立評論」タイトルと要旨,361センテンスに対する適用実験では,文字列分割精度98.5%,文節認定精度95.3%程度であった。本方式の実用化には,画一的なキーワードの定義に従わないキーワードや,不要語に対する例外辞書の充実及び正確な複合語分割技術が必要となる。今

※3) カバー率:(出現語中の辞書表収録延べ数) (出現語延べ数)

※4) 精 度: (正処理数) (対象総数)×100

記事番号: 1110225○ ア 文番号: 1 <入力原文>

770107カーター次期大統領、景気対策大綱を発表

<処理> ID作成区分 ( )

LK 番号 R キーワード 番号 R キーワード

番号 R キーワード

主文 1 ①カーター

2 ①次期大統領

3 ⑥景気

4 ⑥対策 5 ⑥大綱

6 5 発表

非パ 7 ③770107

#### (a) 日本語文構造解析方式の例

プロセス制御システムは複雑化、高級化する傾向にあり、その設計、エンジニヤリング、計装工事などに多大の時間を必要としている。この解決のために、 アナログ制御、直接ディジタル制御、シーケンス制御及び計算機制御を合理的に結合した総合計装システム、「ユニトロールシグマシリーズ」を開発した。 「ユニトロールシグマシリーズ」は、従来の電子式工業計器「ユニトロールEシリーズ」を総合計装システムの立場から見なおしてモジュール化するととも に、マイクロコンピュータを集合計器の形で導入し、性能の拡大、結合の合理化を図ったものである。本稿は、その全体の構成、方式について紹介する。

#### (b) 不要語除去方式の例

図 5 自動インデクシング適用例 日本語文構造解析方式の適用例では,同一主文中のロール(n)とキーワードが「対となって示される。ここで「非パ」と は、文型表によらないロール付与不要語除去方式では、アンダーライン部分がキーワードと認定されたことを示す。

日本語文構造解析方式による自動インデクシングプロトタ イプの処理精度 表 | の外電記事文に対して、出現語いカバー率、出現構 文力バー率がそれぞれ90%であるときの処理精度である。

項番		評		価		項			目		精度(%)
1		文	節	構	į j	戓	語		認	定	85~90
2	П		複	文	木	冓	造	7 11 13 14	認	定	94~98
3	-1		格	支	配	関	7	係	認	定	86~91
4	かけた		体	言	支	配	関	係	認	定	90~95
5	与		П	_	ル	付		与	精	度	80~85

表 4 不要語除去方式による自動インデクシングプロトタイプの処 理精度 表しの「日立評論」のタイトルや要旨に対する処理精度である。

項番			評	価	項	目	1 %	精度,語数
1		文	字	列 多	子割	精	度	98.5%
2	· cu	文	節	認	定	精	度	95.3%
3	+		+-	ワードと正	しく認定	されたも	の	4,384語
4	7		136語					
5	ř			/]\		計		4,520語
6	不		9,470語					
7	要		3,573語					
8	語			/]\		計		13,043語
9		97.0%						
10		キーワード認定適合率 (項3) (項3)+(項7)						
1.1	不要語除去精度 (項 6 ) (項 6 )+(項 7 )						72.6%	
12	キーワードテーブルに収録する必要のある語数							51語
13	不要語テーブルに収録する必要のある語数							1,972語

回の実験では、例外辞書を用意しない場合、キーワードとほ ぼ同数の不要語がキーワードと誤認された。しかし, 本方式 は,適用分野によって例外辞書の登録語は異なるが,抽出ア ルゴリズムに分野依存性がなく、また人による校正情報をも とにした辞書の自動登録が可能であり、汎用性、保守性に富 むシステムに拡張できると考えられる。

プロトタイプは、特定分野の日本語文を対象とする専用シ ステムとして試作され,自動インデクシング方式の実現性を 確認することができたが、実用化を図るためには、適用対象 に応じた柔軟な抽出方式を設定できる日本語処理機能のツー ル化、辞書作成技術の汎用化及び複合語分割処理技術の向上 が必要である。

更に、校正プロセスでは、自動抽出結果の校正と、それを もとに抽出方式や辞書類の更新が行なわれるが、これをシス テムの学習プロセスと考え、学習メカニズムの導入が必要で ある。

#### 5 結 言

日本語情報検索システムでの情報蓄積作業の省力化,付与 インデックスの品質の均質化及びOAの進展に伴う電子化され た文書を検索可能とすることを実現するツールとしての自動 インデクシング方式について、その適用対象、利用目的を考 慮して次の2方式を示した。

- (1) 日本語文構造解析方式
- (2) 不要語除去方式

更に, 両方式についてプロトタイプによる実験を行ない, それらの実現性を確認した。両方式の実用化には, 自動イン デクシングのための複合語分割, 同義語認定などの日本語処 理技術のツール化, 抽出処理及び辞書の拡充のための学習メ カニズムの研究, 更に, OAシステム適用のための小形化など が、今後の課題として残されている。これらの解決のために は、より高い処理精度を得るため処理手法のもつ目標精度に 対する効果の評価と, 意味解析, 談話解析の導入を検討して ゆく必要があると考える。これらの解決により、より使いや すい日本語情報検索システムの運用が可能になるとともに, OAシステムに有機的に結合された文書検索システムの建設が 可能になると期待している。

#### 参考文献

- 1) 福島,外:JICSTオンライン情報検索システムJOIS-IIの概 要,情報管理,23,2,118~131(昭55-5)
- 2) 大山: PATOLIS漢字オンライン特許情報検索システムとそ の利用状況,情報管理,23,2,132~144(昭55-5)
- 3) 菅:日本経済新聞社の記事検索システム「NEEDS-IR」につ いて、情報管理、21、1、63~70(昭53-4)
- 4) 絹川:情報検索のための日本語解析,情報処理,20,10, 907~910 (昭54-10)
- 5) 絹川、外:日本語文構造解析による自動インデクシング方式、 情報処理学会論文誌, 21, 3, 200~207 (昭56-5)